

Theoretical Concepts of Machine Learning Part 1

Sepp Hochreiter

Institute of Bioinformatics
Johannes Kepler University, Linz, Austria

Course



3 ECTS 2 SWS VO (class)

1.5 ECTS 1 SWS UE (exercise)

Basic Course: Master Bioinformatics

Elective Course: Major Comp. Engineering Master Comp. Science

Class: Thu 15:30-17:00 (S3 055)

Exercise: Thu 14:30-15:15 (S3 55)

VO: **three part-exams**

UE: weekly homework (evaluated)

Other Courses



	Lecture		Lecturer	
365,077	Machine Learning: Unsupervised Tech.	VL	Hochreiter	Mon 15:30-17:00/HS 18
365,078	Machine Learn.: Unsup. Tech. – G1	UE	Hochreiter	Mon 13:45-14:30/S2 053
365,095	Machine Learn.: Unsup. Tech. – G2+3	UE	Hochreiter	Mon 14:30-15:15/S2 053
365,041	Theoretical Concepts of Machine Learn.	VL	Hochreiter	Thu 15:30-17:00/S3 055
365,042	Theoretical Concepts of Machine Learn.	UE	Hochreiter	Thu 14:30-15:15/S3 055
365,081	Genome Analysis & Transcriptomics	KV	Regl	Fri 8:30-11:00/S2 053
365,082	Structural Bioinformatics	KV	Regl	Tue 8:30-11:00/HS 11
365,093	Deep Learning and Neural Networks	KV	Unterthiner	Thu 10:15-11:45/MT 226
365,090	Special Topics: Population genetics	KV	Klambauer	block
365,096	Special Topics: AI in Life Sciences	KV	Klambauer	block
365,079	Introduction to R	KV	Bodenhofer	Wed 15:30-17:00/MT 127
365,067	Master's Seminar	SE	Hochreiter	Mon 10:15-11:45/S3 318
365,080	Master's Thesis Seminar SS	SE	Hochreiter	Mon 10:15-11:45/S3 318
365,091	Bachelor's Seminar	SE	Hochreiter	-
365,019	Dissertantenseminar Informatik 3	SE	Hochreiter	Mon 10:15-11:45/S3 318
347,337	Bachelor Sem. Bio. Chemistry (+Thesis)	SE	Hochreiter	-
347,327	Sem. in Structural and Comp. Biochemistry	SE	Hochreiter	-
365,083	Projektpraktikum	PR	Hochreiter	-

Outline



1 Introduction

2 Generalization Error

3 Maximum Likelihood

4 Noise Models

5 Statistical Learning Theory

7 Optimization Techniques

8 Bayes Techniques

1 Introduction

1.1 Machine Learning Introduction

1.2 Course Specific Introduction

2 Generalization Error

2.1 Model Quality Criteria

2.2 Introducing the Generalization Error

2.3 Minimal Risk for a Gaussian Classification Task

3 Maximum Likelihood

3.1 Loss for Unsupervised Learning

3.2 Mean Squared Error, Bias, and Variance

3.3 Fisher Information Matrix, Cramer-Rao Lower Bound

3.4 Maximum Likelihood Estimator

3.5 Properties of Maximum Likelihood Estimator

3.6 Expectation Maximization

3.7 Maximum Entropy Estimation

4 Noise Models

4.1 Gaussian Noise

4.2 Laplace Noise and Minkowski Error

4.3 Binary Models

5 Statistical Learning Theory

- 5.1 Error Bounds for a Gaussian Classification Task
- 5.2 Empirical Risk Minimization
- 5.3 Error Bounds
- 5.4 Structural Risk Minimization
- 5.5 Margin as Complexity Measure
- 5.6 Average Error Bounds for SVMs

6 Theory of Kernels and Dot Products

- 6.1 Kernels, Dot Products, and Mercer's Theorem
- 6.2 Reproducing Kernel Hilbert Space

7 Optimization Techniques

- 7.1 Parameter Optimization and Error Minimization
- 7.2 On-line Optimization
- 7.3 Convex Optimization

8 Bayes Techniques

8.1 Likelihood, Prior, Posterior, Evidence

8.2 Maximum A Posteriori Approach

8.3 Posterior Approximation

8.4 Error Bars and Confidence Intervals

8.5 Hyper-parameter Selection: Evidence Framework

8.6 Hyper-parameter Selection: Integrate Out

8.7 Model Comparison

8.8 Posterior Sampling

- **Stat. Lern. Th.:** V. N. Vapnik; Statistical Learning Theory, Wiley & Sons, 1998
- **Stat. Lern. Th.:** Schölkopf, Smola; Learning with kernels, MIT Press, 2002
- **Estimation Th.:** S. M. Kay; Fundamentals of Statistical Signal Processing, Prentice Hall, 1993
- **Estimation Th.:** M. I. Jordan (ed.); Learning in Graphical Models, MIT Press, 1998
- **ML:** Duda, Hart, Stork; Pattern Classification; Wiley & Sons, 2001
- **ML:** C. M. Bishop; Neural Networks for Pattern Recognition, Oxford Univ. Press, 1995
- **ML:** C. M. Bishop; Pattern Recognition and Machine Learning, Springer, 2006

Chapter 1

Introduction

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

- part of curriculum “master of science in bioinformatics”
- part of curriculum “computer science” (major CE, major int. sys.)
- Machine learning major research topic: Google, Microsoft, Amazon, Facebook, AltaVista, Zalando, and many more
- Applications: computer vision (image recognition), speech recognition, recommender systems, analysis of Big Data, information retrieval
- Mining the web: search engines, social networks, videos, music
- Machine learning applications in biology and medicine:
 - microarrays, sequencing
 - alternative splicing, nucleosome positions, gene regulation
 - single nucleotide polymorphisms / variants (SNPs, SNVs)
 - copy number variations (CNVs)
 - diseases: Alzheimer, Parkinson, cancer, multiples sclerosis, schizophrenia or alcohol dependence

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower Bound and Efficiency

3.4 Maximum Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

This course introduces theoretical concepts machine learning methods.

maximum likelihood estimator is motivated by estimation theory:

- bias
- efficient
- Cramer-Rao lower bound
- Fisher information matrix
- consistent estimator

Consistent: more data lead to better results

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower Bound and Efficiency

3.4 Maximum Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy Estimation

Statistical learning theory:

- empirical risk minimization
- complexity of model classes: VC-dim., growth, annealed entropy
- bounds on the generalization error: Vapnik, Chernoff, covering #
- structural risk minimization

Optimization

- gradient-based
- convex optimization

Bayes framework:

- estimate the posterior
- derive error bounds for model predictions
- optimize hyperparameters
 - integrating out the posterior
 - evidence framework

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

We define quality criteria for selected models in order to pin down a goal for model selection, i.e. learning.

maximum likelihood:

- Unbiased
- Efficient

maximum likelihood in supervised learning: noise models

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower Bound and Efficiency

3.4 Maximum Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

A central question in machine learning is:

Does learning from examples help in the future?

Learning on the training data is called “**empirical risk minimization**” (ERM) in statistical learning theory.

Complexity is restricted: ERM helps

Under mild conditions the convergence is uniform and even fast, i.e. exponentially.

To measure the complexity of the model class we will introduce the VC-dimension (**Vapnik-Chervonenkis**).

theoretical upper bounds on the generalization error

→ “**structural risk minimization**”

We introduce basic techniques for minimizing the error that is techniques for model selection for a parameterized model class.

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

- does learning from examples help in the future?
- “empirical risk minimization” (ERM)
- complexity is restricted and dynamics fixed
- “learning helps”: more training examples improve the model
- converges to the best model for all future data
- convergence is fast
- complexity of a model class: VC-dimension (Vapnik-Chervonenkis)
- “structural risk minimization” (SRM): complexity and model quality
- bounds on the generalization error

Chapter 2

Generalization Error

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower Bound and Efficiency

3.4 Maximum Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

- quality criteria \rightarrow goal for model selection / learning
- approximations
- unsupervised learning: maximum likelihood
- concepts: bias and variance, efficient estimator, Fisher information
- supervised learning considered in an unsupervised framework:
error model

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower Bound and Efficiency

3.4 Maximum Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy Estimation

- learning equivalent to model selection
- quality criteria: future data is optimally processed
- other concepts: visualization, modeling, data compression
- Kohonen networks= self-organizing maps (SOMs):
no scalar quality criterion (potential function)
- advantage quality criteria: - comparison of different models
- quality during learning known
- supervised quality criteria: rate of misclassifications or squared error
- unsupervised criteria: - likelihood
- ratio of between and within cluster distance
- independence of the components
- information content

Generalization Error: Definition



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- performance of a model on future data: **generalization error**

- error on **one** example: **loss** or **error**

- expected loss: **risk** or **generalization error**

Generalization Error: Definition



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Training set: $X = \{\mathbf{x}^1, \dots, \mathbf{x}^l\} \in \mathbb{R}^d$

Label or target value: $y^i \in \mathbb{R}$

Simple: $\mathbf{z} = (\mathbf{x}, y)$ and $\mathbf{z} \in Z = \mathbb{R}^{d+1}$

Training set: $\{\mathbf{z}^1, \dots, \mathbf{z}^l\}$

Matrix notation for training inputs: $\mathbf{X} = (\mathbf{x}^1, \dots, \mathbf{x}^l)^T$

Vector notation for labels: $\mathbf{y} = (y^1, \dots, y^l)^T$

Matrix notation for training set: $\mathbf{Z} = (\mathbf{z}^1, \dots, \mathbf{z}^l)$

Generalization Error: Definition



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

The loss function

quadratic loss:
$$L(y, g(\mathbf{x}; w)) = (y - g(\mathbf{x}; w))^2$$

zero-one loss:
$$L(y, g(\mathbf{x}; w)) = \begin{cases} 0 & \text{for } y = g(\mathbf{x}; w) \\ 1 & \text{for } y \neq g(\mathbf{x}; w) \end{cases}$$

Generalization error:

$$R(g(\cdot; w)) = E_{\mathbf{z}} (L(y, g(\mathbf{x}; w))) = \int_{\mathbf{z}} L(y, g(\mathbf{x}; w)) p(\mathbf{z}) d\mathbf{z}$$

Generalization Error: Definition



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

y is a function of \mathbf{x} (target function: $y = f(\mathbf{x})$) plus noise:

$$y = f(\mathbf{x}) + \epsilon$$

$$p(y | \mathbf{x}) = p_n(y - f(\mathbf{x}))$$

$$p(\mathbf{z}) = p(\mathbf{x}) p(y | \mathbf{x}) = p(\mathbf{x}) p_n(y - f(\mathbf{x}))$$

Now the risk can be computed as

$$R(g(\cdot; w)) = \int_{\mathbf{Z}} L(y, g(\mathbf{x}; w)) p(\mathbf{x}) p_n(y - f(\mathbf{x})) dz =$$

$$\int_{\mathbf{X}} p(\mathbf{x}) \int_{\mathbb{R}} L(y, g(\mathbf{x}; w)) p_n(y - f(\mathbf{x})) dy d\mathbf{x}$$

Generalization Error: Definition



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

↓

$$R(g(\mathbf{x}; w)) = \mathbb{E}_{y|\mathbf{x}} (L(y, g(\mathbf{x}; w))) = \int_{\mathbb{R}} L(y, g(\mathbf{x}; w)) p_n(y - f(\mathbf{x})) dy$$

Generalization Error: Definition



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

The noise-free case is $y = f(\mathbf{x})$

$$R(g(\mathbf{x}; w)) = L(f(\mathbf{x}), g(\mathbf{x}; w)) = L(y, g(\mathbf{x}; w))$$

simplifies to:

$$R(g(\cdot; w)) = \int_{\mathbf{X}} p(\mathbf{x}) L(f(\mathbf{x}), g(\mathbf{x}; w)) d\mathbf{x}$$

Generalization Error: Estimation



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- $p(\mathbf{z})$ is unknown
- especially $p(\mathbf{y}|\mathbf{x})$
- risk cannot be computed
- practical applications: approximation of the risk
- model performance estimation for the user

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Test set approximation:

$$R(g(.; w)) = \mathbb{E}_{\mathbf{z}} (L(y, g(\mathbf{x}; w)))$$

expectation can be approximated using

$$R(g(.; w)) \approx \frac{1}{m} \sum_{i=l+1}^{l+m} L(y^i, g(\mathbf{x}^i; w))$$

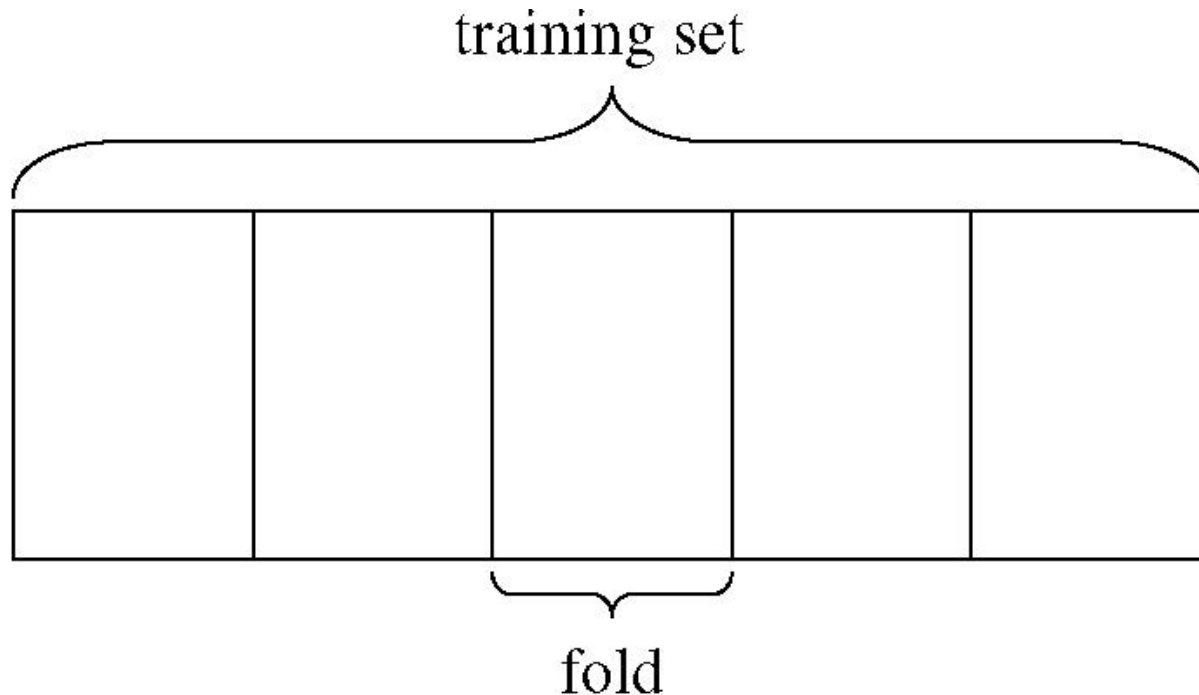
with test set: $\{\mathbf{z}^{l+1}, \dots, \mathbf{z}^{l+m}\}$

Generalization Error: Estimation

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- not enough data for test set (needed for training)
- cross-validation

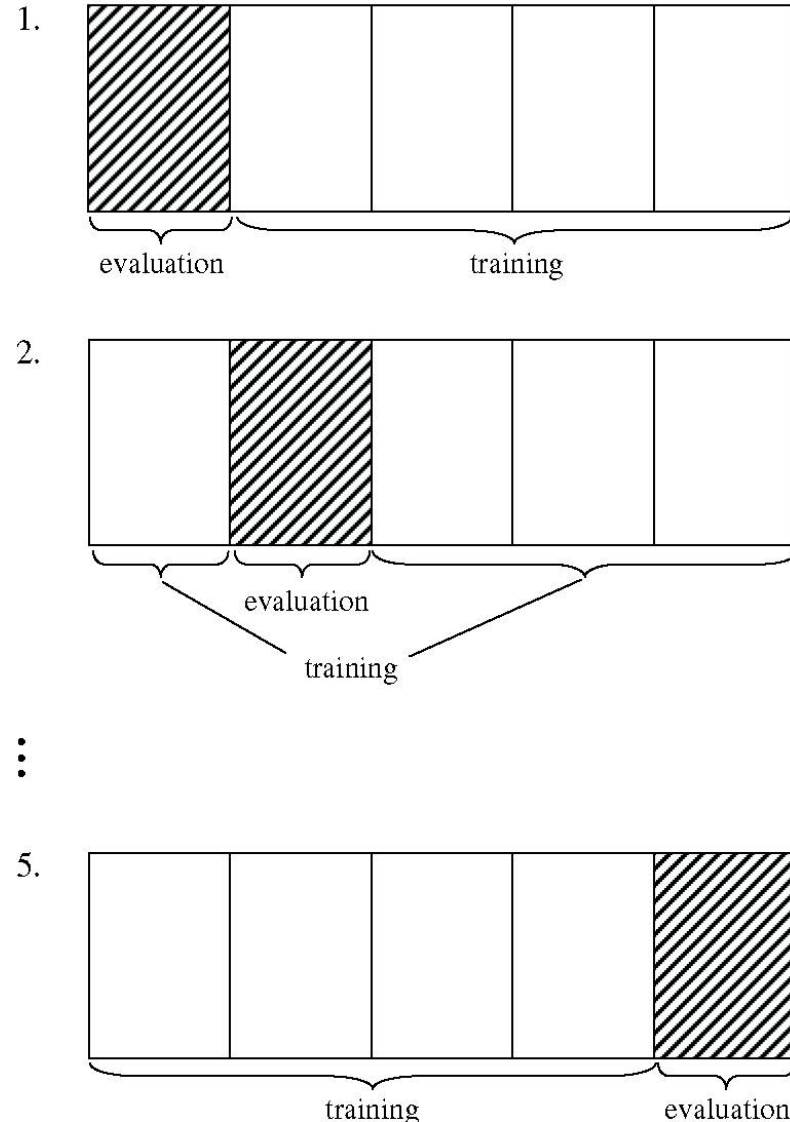
Cross-validation folds:



Generalization Error: Estimation

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

n-fold cross-validation
(here 5-fold):



Generalization Error: Estimation



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

$$R_{n-cv}(Z_l) = \frac{1}{n} \sum_{j=1}^n \underbrace{\frac{n}{l} \sum_{z \in Z_{l/n}^j} \left(L \left(y, g \left(\mathbf{x}; w_j \left(Z_l \setminus Z_{l/n}^j \right) \right) \right) \right)}_{R_{n-cv,j}(Z_l)}$$

cross-validation is an almost unbiased estimator for the generalization error:

$$\mathbb{E}_{Z_{l(1-1/n)}} \left(L \left(g \left(\cdot; w \left(Z_{l(1-1/n)} \right) \right) \right) \right) = \mathbb{E}_{Z_l} \left(R_{n-cv} \left(Z_l \right) \right)$$

Generalization error on trainings size without one fold $l - l/n$ can be estimated by cross-validation on training data l by n -fold cross-validation

Generalization Error: Estimation



1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter Estimation

3.2 MSE, Bias, & Variance

3.3 Cramer-Rao Lower Bound and Efficiency

3.4 Maximum Likelihood Estimation

3.5 Properties Estimator

3.5.1 Invariant

3.5.2 MLE is Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent for Zero CRLB

3.6 Expectation Maximization

3.7 Maximum Entropy Estimation

- advantage: test examples only once used (better than multiple dividing the data into training and test set)
- disadvantage:
 - training sets are overlapping
 - one fold on same model → test examples dependent
 - these dependencies → cv has high variance (one outlier influences all estimates)
- special case: **leave-one-out cross-validation** (LOO-CV)
 - l -fold cross-validation, where each fold is one example
 - test examples to not use the same model
 - training sets are maximal overlapping

Minimal Risk for a Gaussian Classification Task



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Class $y = 1$ data points are drawn according to

$$p(\mathbf{x} \mid y = 1) \propto \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

and class $y = -1$ according to

$$p(\mathbf{x} \mid y = -1) \propto \mathcal{N}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}_{-1})$$

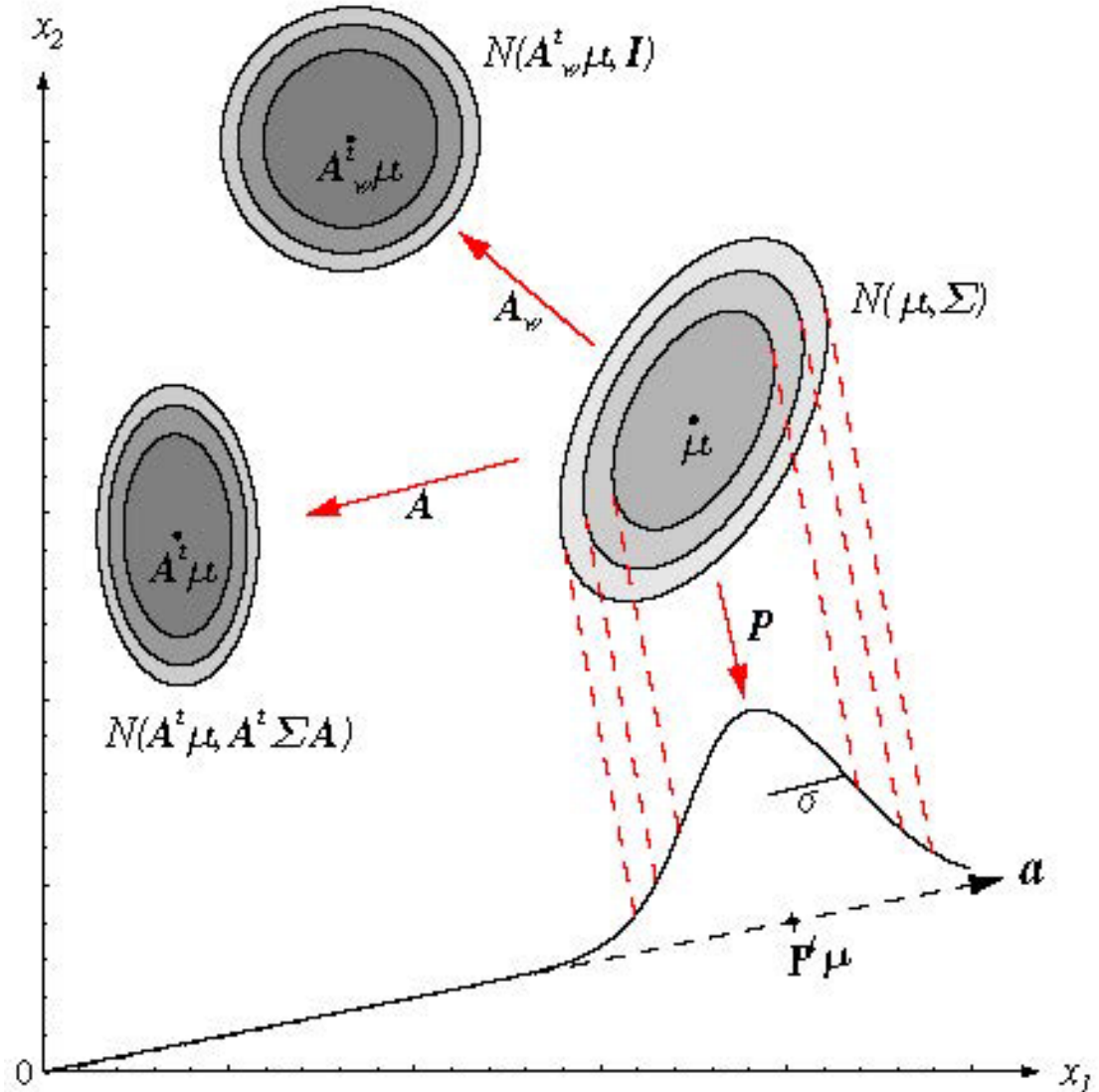
where the Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ has density

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

Minimal Risk for a Gaussian Classification Task

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Linear transformations of Gaussians lead to Gaussians



Minimal Risk for a Gaussian Classification Task



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- probability of observing a point from class $y=1$ at \mathbf{x} : $p(\mathbf{x}, y = 1)$
- probability of observing a point from class $y=-1$ at \mathbf{x} : $p(\mathbf{x}, y = -1)$
- Conditional probability: $p(\mathbf{x}, y = 1) = p(\mathbf{x} | y = 1) p(y = 1)$

- probability of observing a point at \mathbf{x} :

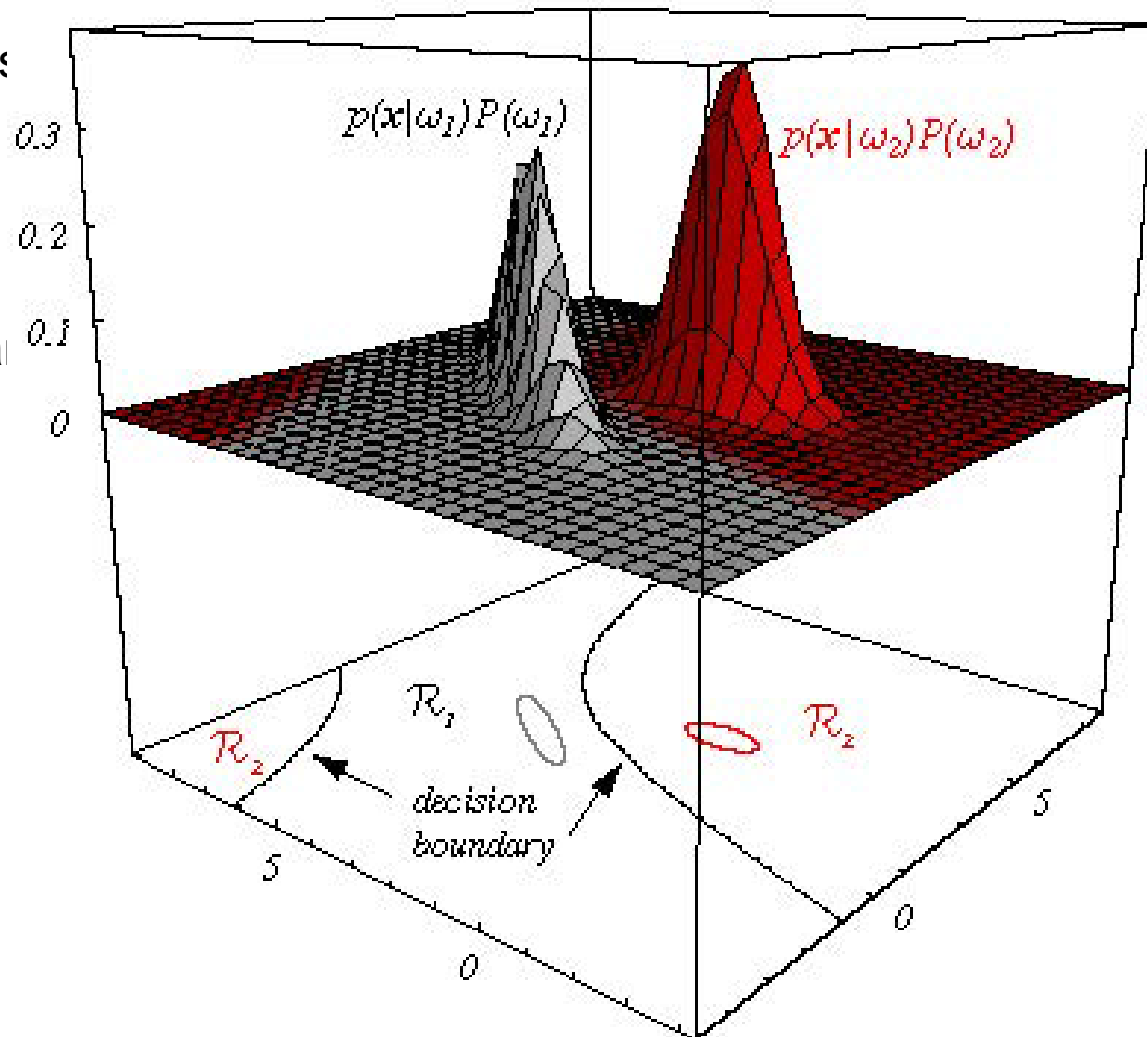
$$p(\mathbf{x}) = p(\mathbf{x}, y = 1) + p(\mathbf{x}, y = -1)$$

y is “integrated out” - here “summed out”

Minimal Risk for a Gaussian Classification Task

1 Introduction
2 Generalization Error
2.1 Model Quality
2.2 Gen. Error
2.2.1 Definition
2.2.2 Estimation
2.3 Minimal Risk Example
3 Maximum Likelihood
3.1 Loss Unsupervised
3.1.1 Projections
3.1.2 Generative
3.1.3 Parameter Estimation
3.2 MSE, Bias, & Variance
3.3 Cramer-Rao Lower Bound and Efficiency
3.4 Maximum Likelihood Estimation
3.5 Properties Estimator
3.5.1 Invariant
3.5.2 MLE is Asymptotically Unbiased and Efficient
3.5.3 MLE is consistent for Zero CRLB
3.6 Expectation Maximization
3.7 Maximum Entropy Estimation

- two-dimensional classification task
- data for each class from a Gaussian (black: class 1, red: class -1)
- optimal discriminant functions are two hyperbolas



Minimal Risk for a Gaussian Classification Task



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- Bayes rule for probability of \mathbf{x} belonging to class $y = 1$:

$$p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | y = 1) p(y = 1)}{p(\mathbf{x})}$$

regions of predicted class $y = 1$:

$$X_1 = \{\mathbf{x} | g(\mathbf{x}) > 0\}$$

regions of predicted class $y = -1$:

$$X_{-1} = \{\mathbf{x} | g(\mathbf{x}) < 0\} .$$

loss function:

$$L(y, g(\mathbf{x}; \mathbf{w})) = \begin{cases} 0 & \text{for } y g(\mathbf{x}; \mathbf{w}) > 0 \\ 1 & \text{for } y g(\mathbf{x}; \mathbf{w}) < 0 \end{cases}$$

Minimal Risk for a Gaussian Classification Task



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

$$\text{Risk: } R(g(.; w)) = \int_{\mathcal{Z}} L(y, g(\mathbf{x}; w)) p(\mathbf{z}) d\mathbf{z}$$

Loss function contributions:

$$X_1 : p(\mathbf{x}, y = -1)$$

$$X_{-1} : p(\mathbf{x}, y = 1)$$

$$R(g(.; w)) = \int_{X_1} p(\mathbf{x}, y = -1) d\mathbf{x} + \int_{X_{-1}} p(\mathbf{x}, y = 1) d\mathbf{x} =$$

$$\int_{X_1} p(y = -1 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} + \int_{X_{-1}} p(y = 1 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} =$$

$$\int_{\mathcal{X}} \left\{ \begin{array}{ll} p(y = -1 | \mathbf{x}) & \text{for } g(\mathbf{x}) > 0 \\ p(y = 1 | \mathbf{x}) & \text{for } g(\mathbf{x}) < 0 \end{array} \right\} p(\mathbf{x}) d\mathbf{x}$$

Minimal Risk for a Gaussian Classification Task



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Optimal discriminant (see later) function:

$$g(\mathbf{x}; w) \begin{cases} > 0 & \text{for } p(y = 1 | \mathbf{x}) > p(y = -1 | \mathbf{x}) \\ < 0 & \text{for } p(y = -1 | \mathbf{x}) > p(y = 1 | \mathbf{x}) \end{cases}$$

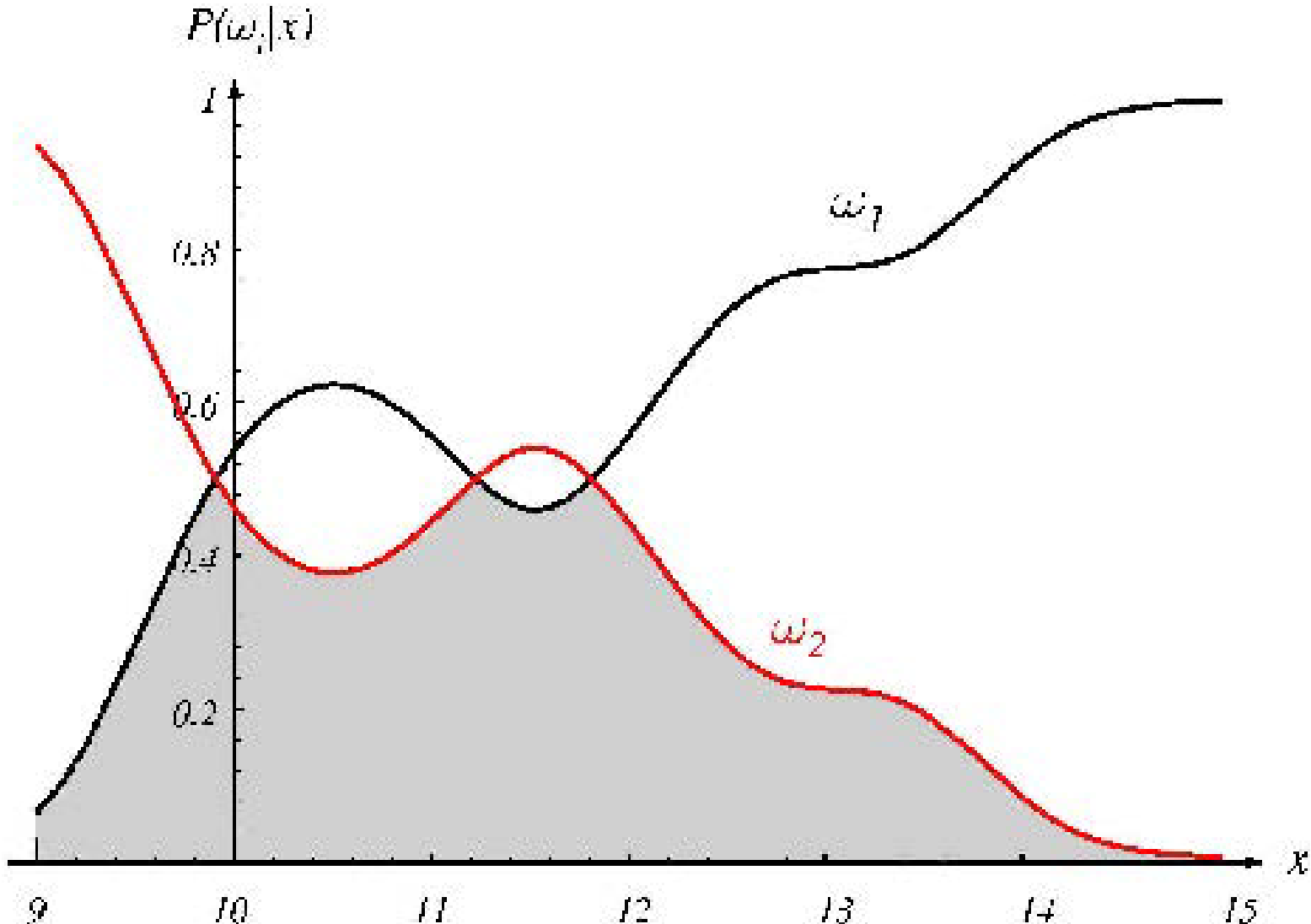
at each position \mathbf{x} take smallest value

The minimal risk is

$$R_{\min} = \int_X \min\{p(\mathbf{x}, y = -1), p(\mathbf{x}, y = 1)\} d\mathbf{x} = \int_X \min\{p(y = -1 | \mathbf{x}), p(y = 1 | \mathbf{x})\} p(\mathbf{x}) d\mathbf{x}$$

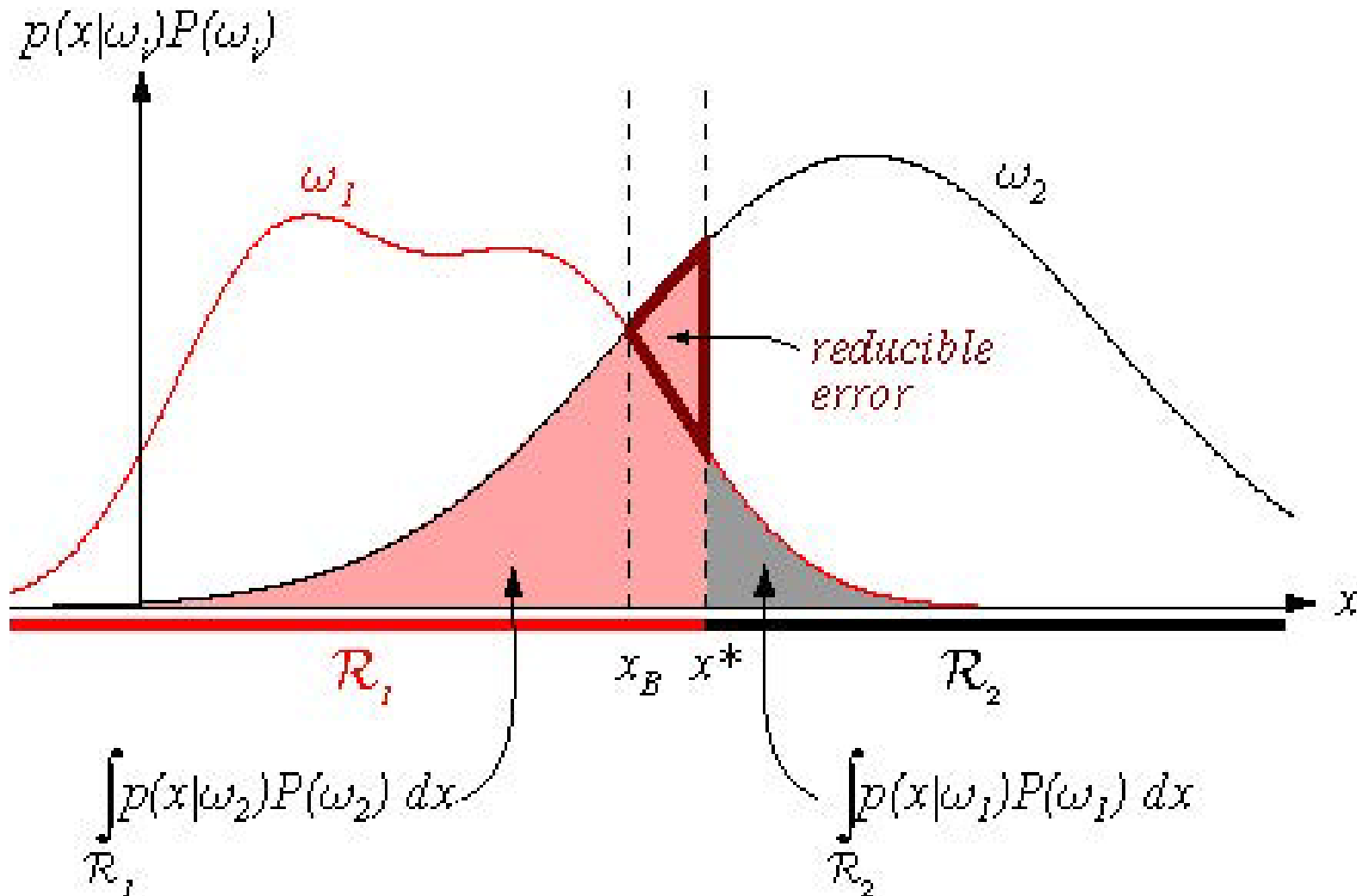
Minimal Risk for a Gaussian Classification Task

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation



Minimal Risk for a Gaussian Classification Task

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation



Minimal Risk for a Gaussian Classification Task



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- **discriminant function** g : $g(\mathbf{x}) > 0$ then \mathbf{x} is assigned to $y = 1$
 $g(\mathbf{x}) < 0$ then \mathbf{x} is assigned to $y = -1$

- **classification function** $\hat{y}(\mathbf{x})$: $\hat{y}(\mathbf{x}) = \text{sign}(g(\mathbf{x}))$

- optimal discriminant functions (minimal risk):

$$g(\mathbf{x}) = p(y = 1 | \mathbf{x}) - p(y = -1 | \mathbf{x})$$

or

$$g(\mathbf{x}) = \ln p(y = 1 | \mathbf{x}) - \ln p(y = -1 | \mathbf{x}) = \ln \frac{p(\mathbf{x} | y = 1)}{p(\mathbf{x} | y = -1)} + \ln \frac{p(y = 1)}{p(y = -1)}$$

Minimal Risk for a Gaussian Classification Task



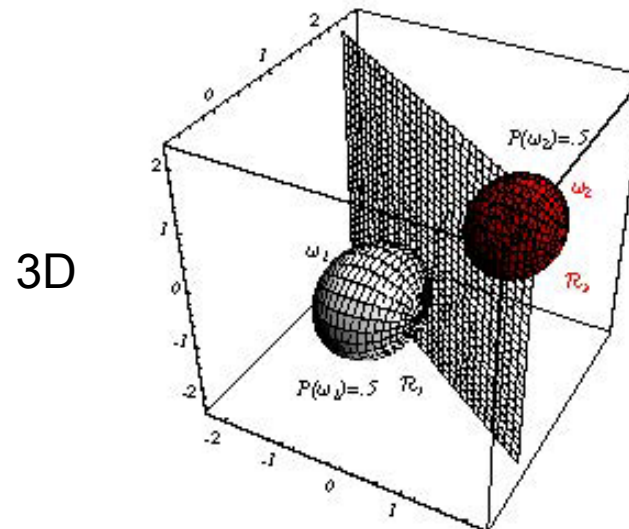
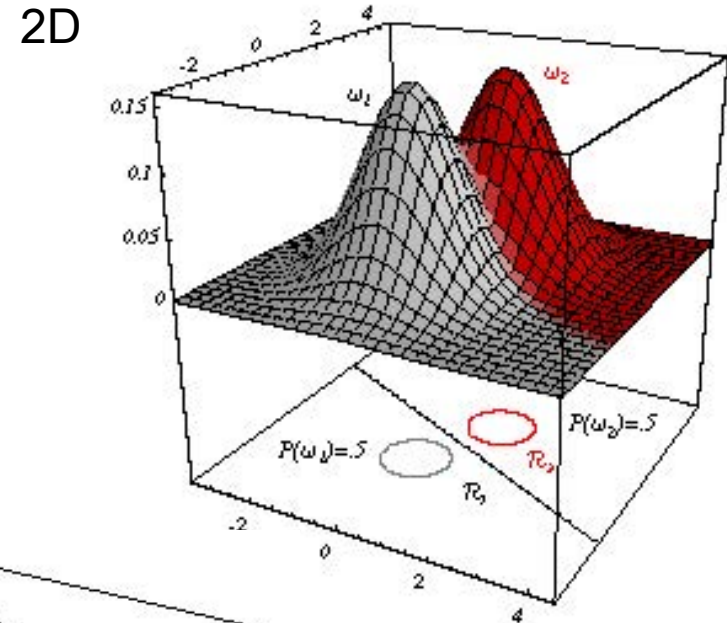
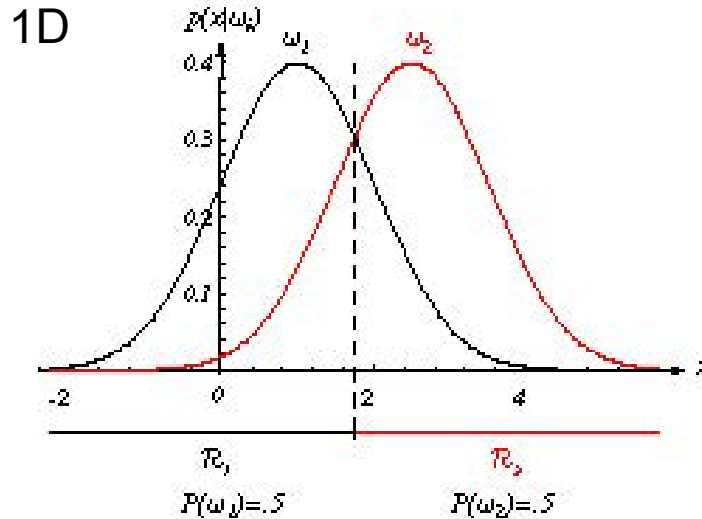
- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

For Gaussians:

$$\begin{aligned}
 g(\mathbf{x}) &= -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - \frac{d}{2} \ln 2\pi - \\
 &\frac{1}{2} \ln |\boldsymbol{\Sigma}_1| + \ln p(y = 1) + \\
 &\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \frac{d}{2} \ln 2\pi + \frac{1}{2} \ln |\boldsymbol{\Sigma}_2| - \ln p(y = -1) = \\
 &-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_1| + \ln p(y = 1) + \\
 &\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \frac{1}{2} \ln |\boldsymbol{\Sigma}_2| - \ln p(y = -1) = \\
 &-\frac{1}{2} \mathbf{x}^T \underbrace{(\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})}_A \mathbf{x} + \mathbf{x}^T \underbrace{(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2)}_w - \\
 &\frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 - \frac{1}{2} \ln |\boldsymbol{\Sigma}_1| + \frac{1}{2} \ln |\boldsymbol{\Sigma}_2| + \\
 &\ln p(y = 1) - \ln p(y = -1) = \\
 &-\frac{1}{2} \mathbf{x}^T A \mathbf{x} + w^T \mathbf{x} + b
 \end{aligned}$$

Minimal Risk for a Gaussian Classification Task

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation



Minimal Risk for a Gaussian Classification Task

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

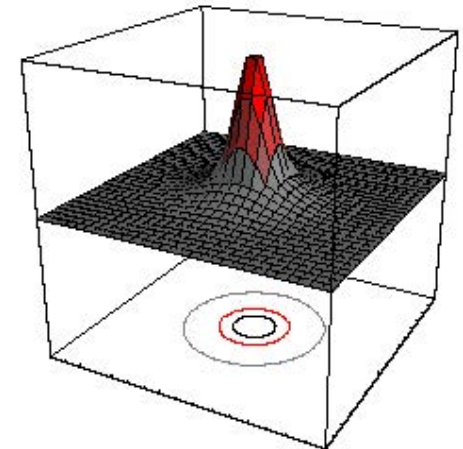
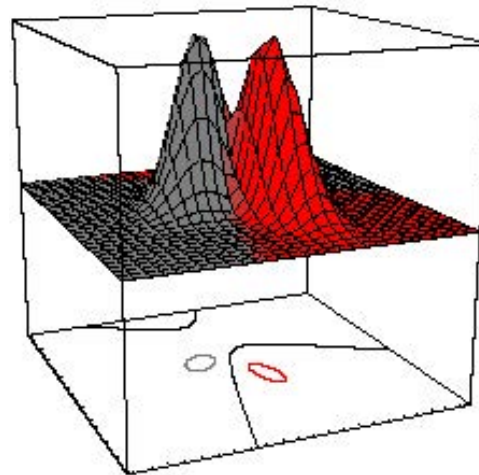
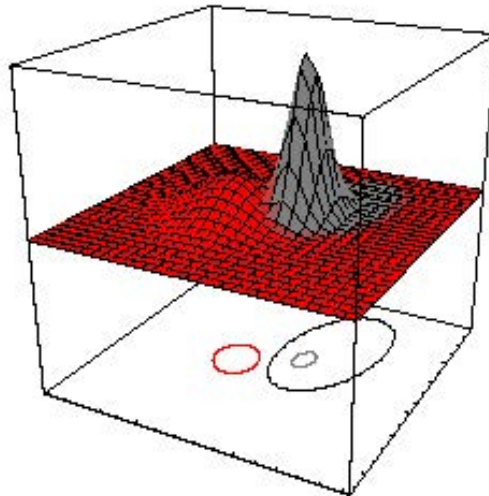
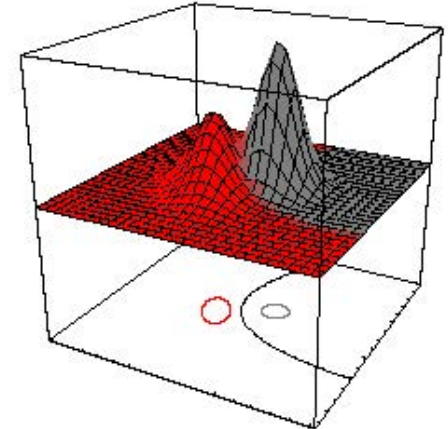
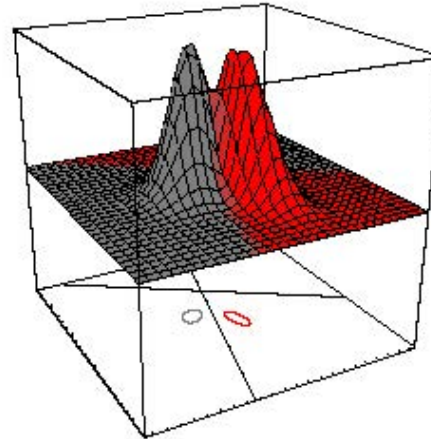
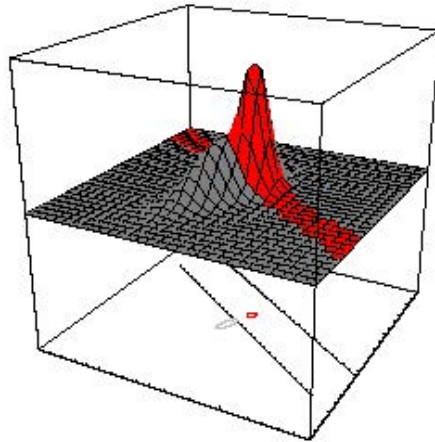
for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation



Chapter 3

Maximum Likelihood

Maximum Likelihood



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- One of the major objectives if learning generative models
- It has certain theoretical properties
- Theoretical concepts like efficient estimator or biased estimator are introduced
- Even supervised methods can be viewed as special case of maximum likelihood

Loss for Unsupervised Learning



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

First we consider different loss functions which are used for unsupervised learning

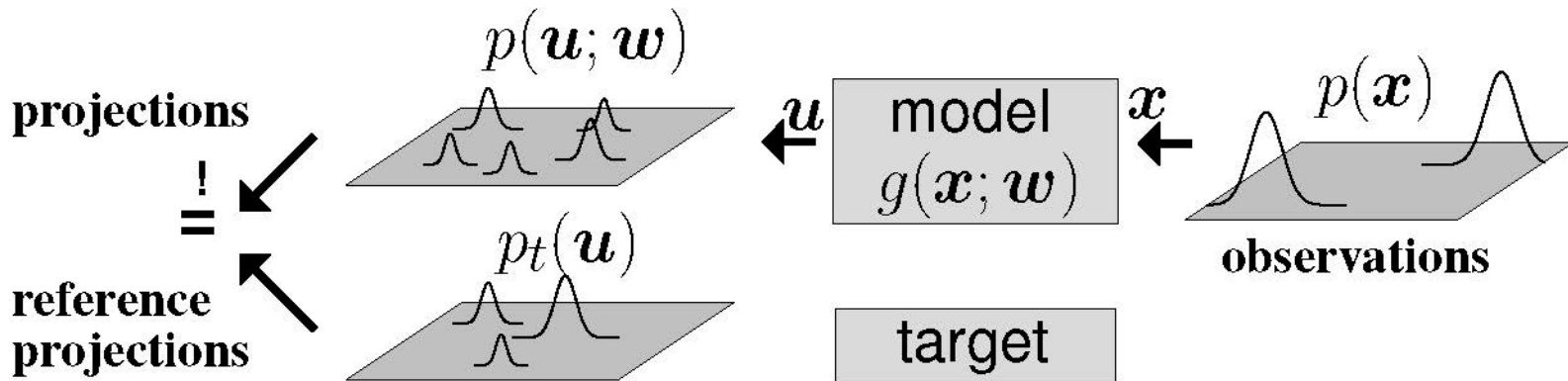
Generative approaches → maximum likelihood

Projection methods → desired property

Loss for Unsupervised Learning

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

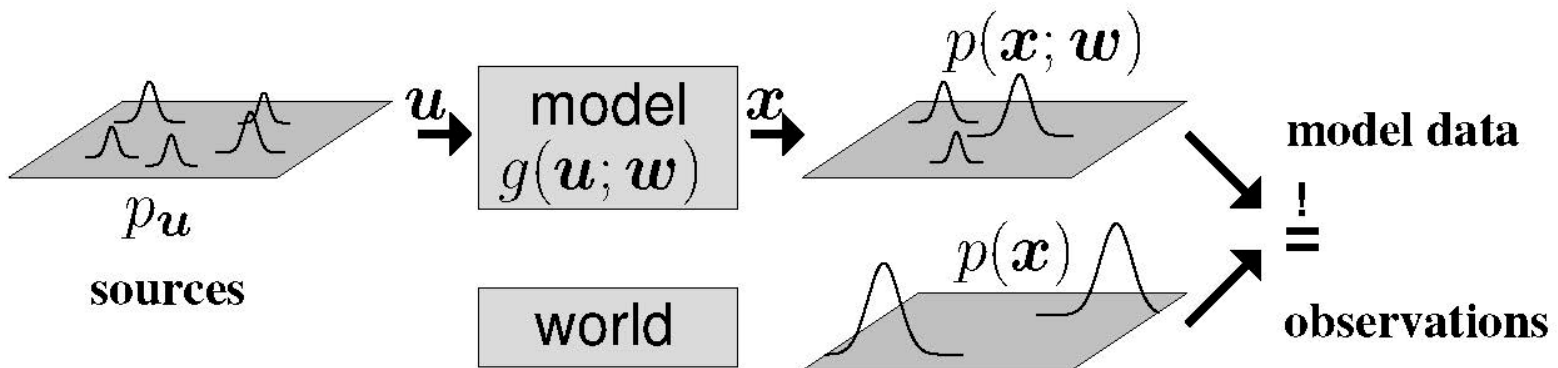
- data projection into another space with desired requirements



Loss for Unsupervised Learning

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

“**generative model**”: model simulates the world and produces the same data as the world



Loss for Unsupervised Learning



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- parameterized model known
- task: estimate actual parameters
- loss: difference between true and estimated parameter
- evaluate estimator: expected loss

Mean Squared Error, Bias, and Variance



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Theoretical concepts of parameter estimation

- training data: $\{\mathbf{x}\} = \{\mathbf{x}^1, \dots, \mathbf{x}^l\}$

simply $\mathbf{X} = (\mathbf{x}^1, \dots, \mathbf{x}^l)^T$ (the matrix of training data)

- true parameter vector: \mathbf{w}
- estimate of \mathbf{w} : $\hat{\mathbf{w}}$

Mean Squared Error, Bias, and Variance



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- **unbiased** estimator:
$$\mathbb{E}_{\mathbf{X}} \hat{\mathbf{w}} = \mathbf{w}$$

on average (over training set) the true parameter is obtained

- **bias**:
$$b(\hat{\mathbf{w}}) = \mathbb{E}_{\mathbf{X}} \hat{\mathbf{w}} - \mathbf{w}$$

- **variance**:
$$\text{var}(\hat{\mathbf{w}}) = \mathbb{E}_{\mathbf{X}} \left((\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}))^T (\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}})) \right)$$

- **mean squared error** (MSE, different to supervised loss):

$$\text{mse}(\hat{\mathbf{w}}) = \mathbb{E}_{\mathbf{X}} \left((\hat{\mathbf{w}} - \mathbf{w})^T (\hat{\mathbf{w}} - \mathbf{w}) \right)$$

expected squared error between the estimated and true parameter

Mean Squared Error, Bias, and Variance



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

$$\text{mse}(\hat{w}) = \mathbb{E}_{\mathbf{X}} \left((\hat{w} - w)^T (\hat{w} - w) \right) =$$

$$\mathbb{E}_{\mathbf{X}} \left(\underbrace{((\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w})) + (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w))^T}_{\text{zero}} \right)$$

$$\mathbb{E}_{\mathbf{X}} \left(((\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w})) + (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w)) \right) =$$

$$\mathbb{E}_{\mathbf{X}} \left((\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w}))^T (\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w})) -$$

$$2 (\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w}))^T (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w) +$$

$$(\mathbb{E}_{\mathbf{X}}(\hat{w}) - w)^T (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w) \right) =$$

Only \hat{w} depends on \mathbf{X}

$$\mathbb{E}_{\mathbf{X}} \left((\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w}))^T (\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w})) \right) +$$

$$(\mathbb{E}_{\mathbf{X}}(\hat{w}) - w)^T (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w) =$$

$$\text{var}(\hat{w}) + b^2(\hat{w})$$

$$\mathbb{E}_{\mathbf{X}} \left((\hat{w} - \mathbb{E}_{\mathbf{X}}(\hat{w}))^T (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w) \right) =$$

$$(\mathbb{E}_{\mathbf{X}}(\hat{w}) - \mathbb{E}_{\mathbf{X}}(\hat{w}))^T (\mathbb{E}_{\mathbf{X}}(\hat{w}) - w) = 0$$

Mean Squared Error, Bias, and Variance



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Averaging reduces variance – each of the N subsets has l/N examples which gives l examples in total

Average is $\hat{w}_{a_N} = \frac{1}{N} \sum_{p=1}^N \hat{w}_i$ where

$$\hat{w}_i = \hat{w}_i(\mathbf{X}_i) \quad \mathbf{X}_i = \left\{ \mathbf{x}^{(i-1)l/N+1}, \dots, \mathbf{x}^{il/N} \right\}$$

Unbiased:
$$\mathbb{E}_{\mathbf{X}}(\hat{w}_{a_N}) = \frac{1}{N} \sum_{p=1}^N \mathbb{E}_{\mathbf{X}_i} \hat{w}_i = \frac{1}{N} \sum_{p=1}^N w = w$$

Variance:
$$\text{covar}_{\mathbf{X}}(\hat{w}_{a_N}) = \frac{1}{N^2} \sum_{p=1}^N \text{covar}_{\mathbf{X}_i}(\hat{w}_i) =$$
$$\frac{1}{N^2} \sum_{p=1}^N \text{covar}_{\mathbf{X}, l/N}(\hat{w}) = \frac{1}{N} \text{covar}_{\mathbf{X}, l/N}(\hat{w})$$

Mean Squared Error, Bias, and Variance



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- averaging: training sets X_i are independent, therefore covariance between them vanishes
- **Minimal Variance Unbiased** (MVU) estimator: construct from all **unbiased** estimators the one with minimal variance
- MVU estimator does not always exist
- methods to check whether a given estimator is a MVU

Fisher Information Matrix, Cramer-Rao Lower Bound, and Efficiency



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- We will find a lower bound for the variance of an unbiased estimator: **Cramer-Rao Lower Bound** (that is a lower bound for the MSE)

- We need the **Fisher information matrix** \mathbf{I}_F :

$$\mathbf{I}_F(\mathbf{w}) : [\mathbf{I}_F(\mathbf{w})]_{ij} = \mathbb{E}_{p(\mathbf{x};\mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x};\mathbf{w})}{\partial w_i} \frac{\partial \ln p(\mathbf{x};\mathbf{w})}{\partial w_j} \right)$$

$$\mathbb{E}_{p(\mathbf{x};\mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x};\mathbf{w})}{\partial w_i} \frac{\partial \ln p(\mathbf{x};\mathbf{w})}{\partial w_j} \right) = \int \frac{\partial \ln p(\mathbf{x};\mathbf{w})}{\partial w_i} \frac{\partial \ln p(\mathbf{x};\mathbf{w})}{\partial w_j} p(\mathbf{x};\mathbf{w}) d\mathbf{x}$$

Fisher Information Matrix, Cramer-Rao Lower Bound, and Efficiency



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

If $p(\mathbf{x}; \mathbf{w})$ satisfies $\forall \mathbf{w} : \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right) = \mathbf{0}$

then the Fisher information matrix is

$$\mathbf{I}_F(\mathbf{w}) : \mathbf{I}_F(\mathbf{w}) = - \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}} \right)$$

Fisher information: information of observation \mathbf{x} about parameter \mathbf{w} upon which the parameterized density function $p(\mathbf{x}; \mathbf{w})$ of \mathbf{x} depends

Fisher Information Matrix, Cramer-Rao Lower Bound, and Efficiency



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Theorem 1 (Cramer-Rao Lower Bound (CRLB))

Assume that

$$\forall \mathbf{w} : \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right) = \mathbf{0}$$

and that the estimator $\hat{\mathbf{w}}$ is unbiased.

Then,

$$\text{covar}(\hat{\mathbf{w}}) - \mathbf{I}_F^{-1}(\mathbf{w})$$

is positive definite:

$$\text{covar}(\hat{\mathbf{w}}) - \mathbf{I}_F^{-1}(\mathbf{w}) \geq \mathbf{0} .$$

An unbiased estimator attains the bound in that $\text{covar}(\hat{\mathbf{w}}) = \mathbf{I}_F^{-1}(\mathbf{w})$ if and only if

$$\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} = \mathbf{A}(\mathbf{w}) (\mathbf{g}(\mathbf{x}) - \mathbf{w})$$

for some function \mathbf{g} and square matrix $\mathbf{A}(\mathbf{w})$.

In this case the MVU estimator is

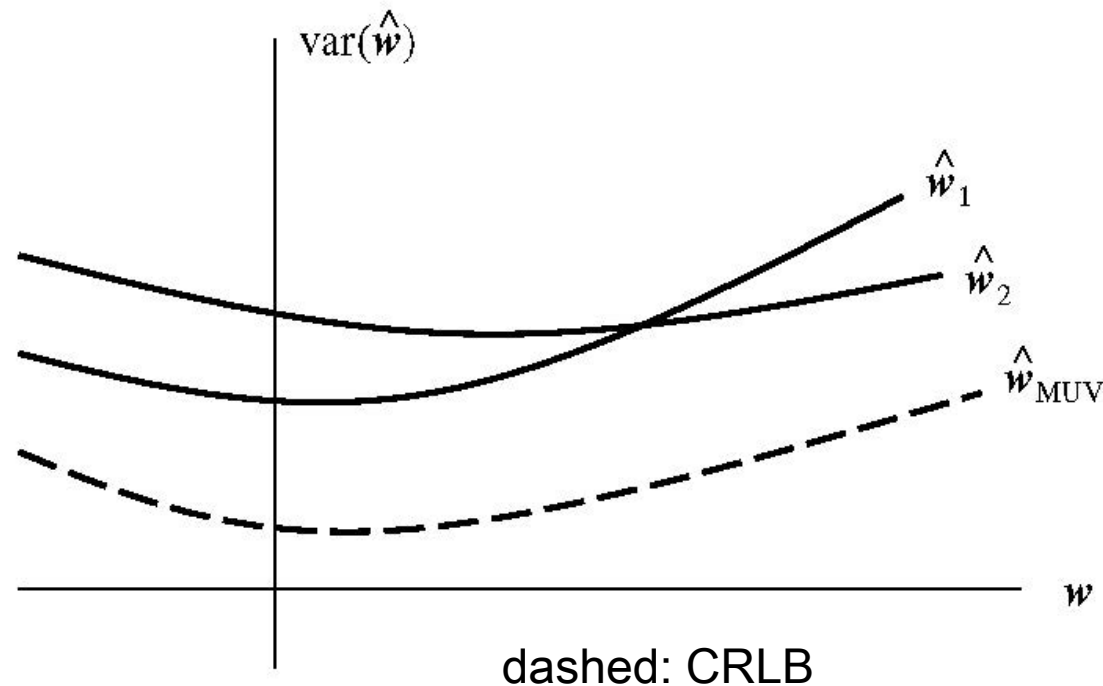
$$\hat{\mathbf{w}} = \mathbf{g}(\mathbf{x}) \quad \text{with} \quad \text{covar}(\hat{\mathbf{w}}) = \mathbf{A}^{-1}(\mathbf{w}) .$$

Fisher Information Matrix, Cramer-Rao Lower Bound, and Efficiency

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

$$\text{var}(\hat{w}_i) = [\text{covar}(\hat{w})]_{ii} \geq [\mathbf{I}_F^{-1}(\mathbf{w})]_{ii}$$

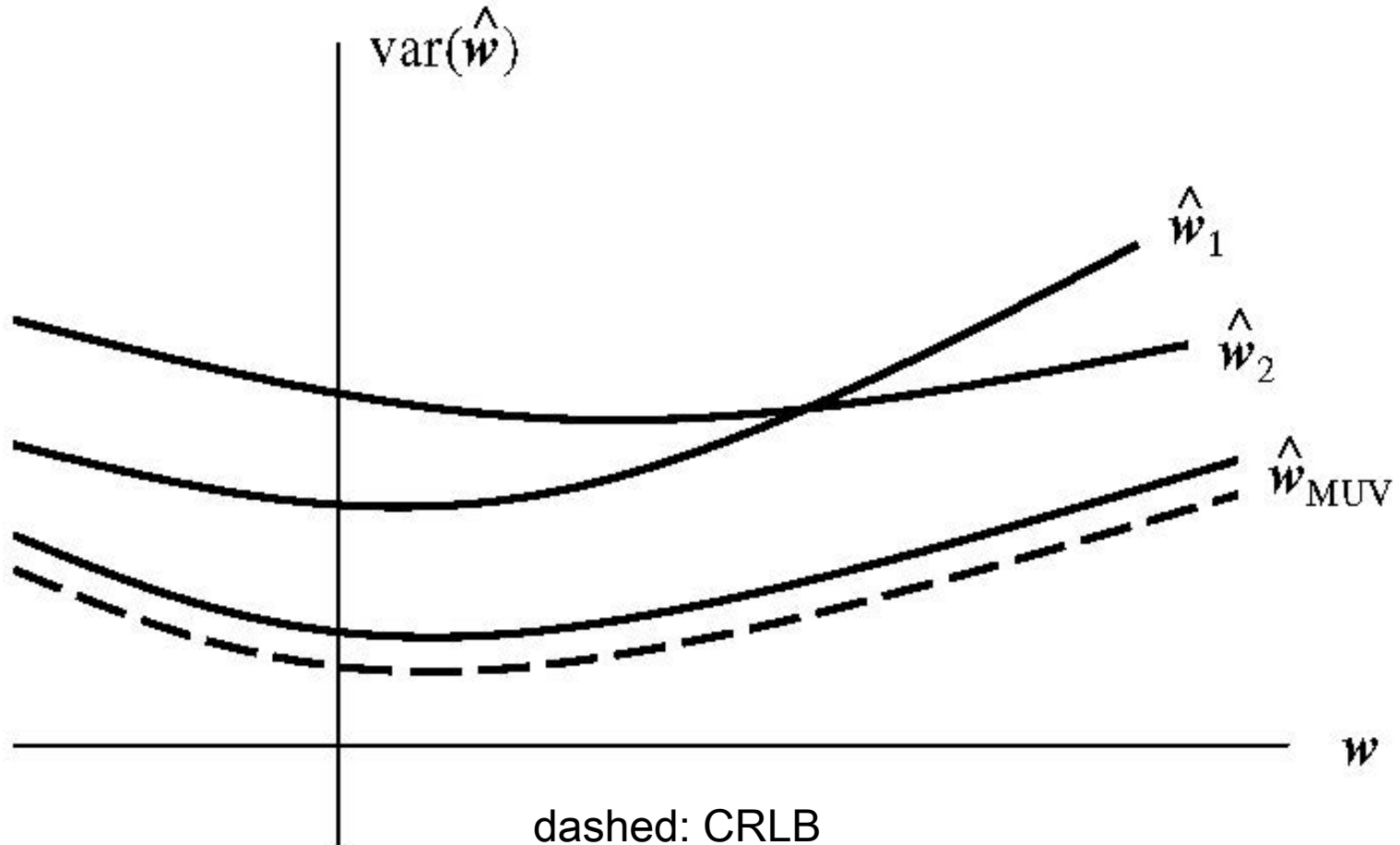
- **efficient** estimator: reaches the CRLB (efficiently uses the data)
- MVU estimator can be efficient but need not



dashed: CRLB

Fisher Information Matrix, Cramer-Rao Lower Bound, and Efficiency

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation



Maximum Likelihood Estimator



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- MVU estimator is unknown or does not exist
- Maximum Likelihood Estimator (MLE)
- MLE can be applied to a broad range of problems
- MLE approximates the MVU estimator for large data sets
- MLE is even asymptotically efficient and unbiased
- MLE does everything right and this efficiently (enough data)

Maximum Likelihood Estimator



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

The likelihood \mathcal{L} of the data set $\{\mathbf{x}\} = \{\mathbf{x}^1, \dots, \mathbf{x}^l\}$:

$$\mathcal{L}(\{\mathbf{x}\}; \mathbf{w}) = p(\{\mathbf{x}\}; \mathbf{w})$$

probability of the model $p(\mathbf{x}; \mathbf{w})$ to produce the data

iid (independent identical distributed) data:

$$\mathcal{L}(\{\mathbf{x}\}; \mathbf{w}) = p(\{\mathbf{x}\}; \mathbf{w}) = \prod_{i=1}^l p(\mathbf{x}^i; \mathbf{w})$$

Negative log-likelihood:

$$-\ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}) = -\sum_{i=1}^l \ln p(\mathbf{x}^i; \mathbf{w})$$

Maximum Likelihood Estimator



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- likelihood is based on finite many densities values which have zero measure: problem?
- assume instead of $p(\mathbf{x}; \mathbf{w})$ the volume element $p(\mathbf{x}^i; \mathbf{w}) d\mathbf{x}$ (region around \mathbf{x}^i)

- MLE popular:
 - simple use
 - properties

Properties of Maximum Likelihood Estimator



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

MLE:

- invariant under parameter change
- asymptotically unbiased and efficient → asymptotically optimal
- consistent for zero CRLB

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Theorem 1 (Parameter Change Invariance)

Let g be a function changing the parameter w into parameter u : $u = g(w)$, then

$$\hat{u} = g(\hat{w}) ,$$

where the estimators are MLE.

If g changes w into different u then $\hat{u} = g(\hat{w})$ maximizes the likelihood function

$$\max_{w:u=g(w)} p(\{x\}; w) .$$

Properties of Maximum Likelihood Estimator



1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

The maximum likelihood estimator is **asymptotically unbiased**:

$$\mathbb{E}_{p(\mathbf{x}; \mathbf{w})}(\hat{\mathbf{w}}) \xrightarrow{l \rightarrow \infty} \mathbf{w}$$

The maximum likelihood estimator is **asymptotically efficient**:

$$\text{covar}(\hat{\mathbf{w}}) \xrightarrow{l \rightarrow \infty} \text{CRLB}$$

Theorem 1 (MLE Asymptotic Properties)

If $p(\mathbf{x}; \mathbf{w})$ satisfies

$$\forall \mathbf{w} : \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right) = \mathbf{0}$$

then the MLE which maximizes $p(\{\mathbf{x}\}; \mathbf{w})$ is asymptotically distributed according to

$$\hat{\mathbf{w}} \stackrel{l \rightarrow \infty}{\propto} \mathcal{N}(\mathbf{w}, \mathbf{I}_F^{-1}(\mathbf{w})) ,$$

where $\mathbf{I}_F(\mathbf{w})$ is the Fisher information matrix evaluated at the unknown parameter \mathbf{w} .

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Properties of Maximum Likelihood Estimator



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- practical applications: finite examples \rightarrow MLE performance unknown

- Example: general linear model $\mathbf{x} = \mathbf{A}\mathbf{w} + \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon} \propto \mathcal{N}(\mathbf{0}, \mathbf{C})$

MLE is $\hat{\mathbf{w}} = (\mathbf{A}^T \mathbf{C}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{C}^{-1} \mathbf{x}$
which is efficient and MUV

$$\hat{\mathbf{w}} \propto \mathcal{N}\left(\mathbf{w}, (\mathbf{A}^T \mathbf{C}^{-1} \mathbf{A})^{-1}\right).$$

Note the noise covariance must be known

Properties of Maximum Likelihood Estimator



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- **consistent:** $\hat{w} \xrightarrow{l \rightarrow \infty} w$

for large training sets the estimator approaches the true value (difference to unbiased \rightarrow variance decreases)

- more formal definition for consistency as

$$\lim_{l \rightarrow \infty} p(|\hat{w} - w| > \epsilon) = 0$$

Thus, the MLE is consistent if the CRLB is zero

Expectation Maximization

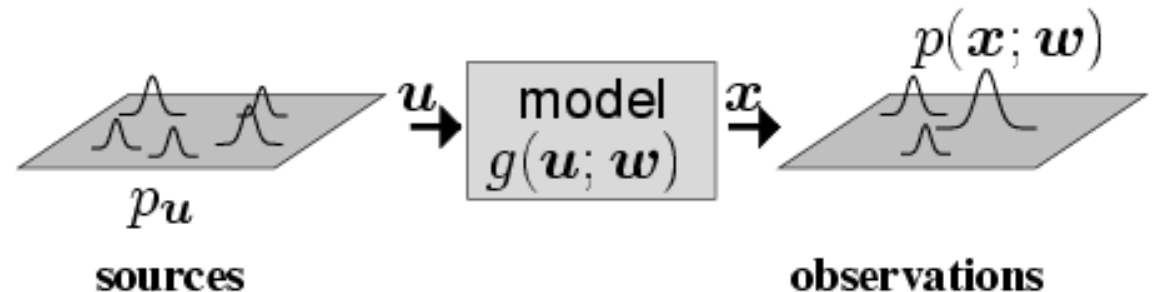
- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

likelihood can be optimized by gradient descent methods

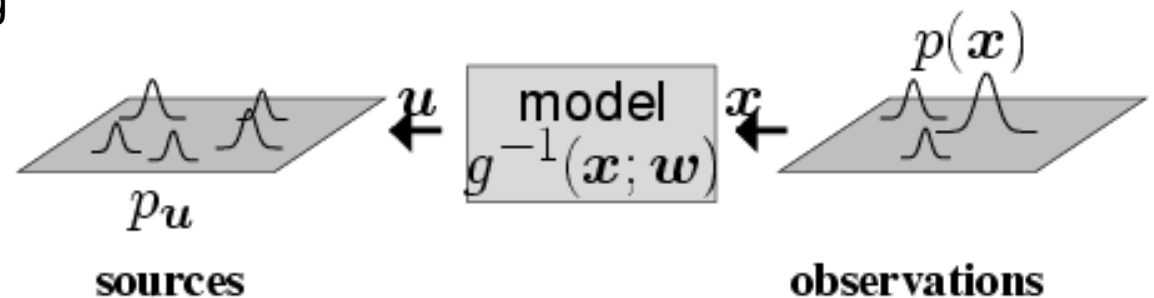
likelihood cannot be computed analytically:

- hidden states
- many-to-one output mapping
- non-linearities

Generative Model



Likelihood



$$p(x; w) = \int_U p_u(u) \delta(x = g(u; w)) du$$

Expectation Maximization



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

- hidden variables, latent variables, unobserved variables \mathbf{u}
- likelihood is determined by all \mathbf{u} mapped to \mathbf{x}

$$p(\mathbf{x}; \boldsymbol{w}) = \int_U p_{\mathbf{u}}(\mathbf{u}) \delta(\mathbf{x} = g(\mathbf{u}; \boldsymbol{w})) d\mathbf{u}$$

Expectation Maximization



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Expectation Maximization (EM) algorithm:

- joint probability $p(\mathbf{x}, \mathbf{u}; \mathbf{w})$ is easier to compute than likelihood
- estimate $p(\mathbf{u} | \mathbf{x}; \mathbf{w})$ by $Q(\mathbf{u} | \mathbf{x})$

$$\ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}) = \ln p(\{\mathbf{x}\}; \mathbf{w}) = \ln \int_U p(\{\mathbf{x}\}, \mathbf{u}; \mathbf{w}) d\mathbf{u} =$$

$$\ln \int_U \frac{Q(\mathbf{u} | \{\mathbf{x}\})}{Q(\mathbf{u} | \{\mathbf{x}\})} p(\{\mathbf{x}\}, \mathbf{u}; \mathbf{w}) d\mathbf{u} \geq \text{Jensen's inequality}$$

$$\int_U Q(\mathbf{u} | \{\mathbf{x}\}) \ln \frac{p(\{\mathbf{x}\}, \mathbf{u}; \mathbf{w})}{Q(\mathbf{u} | \{\mathbf{x}\})} d\mathbf{u} =$$

$$\int_U Q(\mathbf{u} | \{\mathbf{x}\}) \ln p(\{\mathbf{x}\}, \mathbf{u}; \mathbf{w}) d\mathbf{u} -$$

$$\int_U Q(\mathbf{u} | \{\mathbf{x}\}) \ln Q(\mathbf{u} | \{\mathbf{x}\}) d\mathbf{u} =$$

$$\mathcal{F}(Q, \mathbf{w})$$

Expectation Maximization



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

EM algorithm is an iteration between "E"-step and "M"-step:

E-step:

$$Q_{k+1} = \arg \max_Q \mathcal{F}(Q, \mathbf{w}_k)$$

M-step:

$$\mathbf{w}_{k+1} = \arg \max_{\mathbf{w}} \mathcal{F}(Q_{k+1}, \mathbf{w})$$

Expectation Maximization



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

After E-step:

$$\begin{aligned} Q_{k+1}(\mathbf{u} \mid \{\mathbf{x}\}) &= p(\mathbf{u} \mid \{\mathbf{x}\}; \mathbf{w}_k) \\ \mathcal{F}(Q_{k+1}, \mathbf{w}_k) &= \ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}_k) \end{aligned}$$

Proof: $p(\mathbf{u}, \{\mathbf{x}\}; \mathbf{w}_k) = p(\mathbf{u} \mid \{\mathbf{x}\}; \mathbf{w}_k) p(\{\mathbf{x}\}; \mathbf{w}_k)$

$$\begin{aligned} \mathcal{F}(Q, \mathbf{w}) &= \int_U Q(\mathbf{u} \mid \{\mathbf{x}\}) \ln \frac{p(\{\mathbf{x}\}, \mathbf{u}; \mathbf{w})}{Q(\mathbf{u} \mid \{\mathbf{x}\})} d\mathbf{u} = \\ &= \int_U Q(\mathbf{u} \mid \{\mathbf{x}\}) \ln \frac{p(\mathbf{u} \mid \{\mathbf{x}\}; \mathbf{w})}{Q(\mathbf{u} \mid \{\mathbf{x}\})} d\mathbf{u} + \ln p(\{\mathbf{x}\}; \mathbf{w}) = \\ &= \int_U Q(\mathbf{u} \mid \{\mathbf{x}\}) \ln \frac{p(\mathbf{u} \mid \{\mathbf{x}\}; \mathbf{w})}{Q(\mathbf{u} \mid \{\mathbf{x}\})} d\mathbf{u} + \ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}) \end{aligned}$$

Kullback-Leibler divergence:

$$D_{\text{KL}}(Q \parallel p) = \int_U Q(\mathbf{u} \mid \{\mathbf{x}\}) \ln \frac{Q(\mathbf{u} \mid \{\mathbf{x}\})}{p(\mathbf{u} \mid \{\mathbf{x}\}; \mathbf{w})} d\mathbf{u} \geq 0$$

Zero for: $Q(\mathbf{u} \mid \{\mathbf{x}\}) = p(\mathbf{u} \mid \{\mathbf{x}\}; \mathbf{w})$

Expectation Maximization



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

EM increases the lower bound in both steps

Beginning of the M-step: $\mathcal{F}(Q_{k+1}, \mathbf{w}_k) = \ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}_k)$

E-step does not change the parameters

$$\ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}_k) = \mathcal{F}(Q_{k+1}, \mathbf{w}_k) \leq \mathcal{F}(Q_{k+1}, \mathbf{w}_{k+1}) \leq \mathcal{F}(Q_{k+2}, \mathbf{w}_{k+1}) = \ln \mathcal{L}(\{\mathbf{x}\}; \mathbf{w}_{k+1})$$

EM algorithm:

- hidden Markov models
- mixture of Gaussians
- factor analysis
- independent component analysis

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

maximum entropy probability distribution:

- maximal entropy given a class of distributions
- minimal prior assumptions
- physical systems converge to maximal entropy configurations
- most likely observed solution
- connection: statistical mechanics and information theory

principle of maximum entropy first expounded by E.T. Jaynes in 1957

1 Introduction

2 Generalization Error

2.1 Model Quality

2.2 Gen. Error

2.2.1 Definition

2.2.2 Estimation

2.3 Minimal Risk

Example

3 Maximum Likelihood

3.1 Loss Unsupervised

3.1.1 Projections

3.1.2 Generative

3.1.3 Parameter

Estimation

3.2 MSE, Bias, &

Variance

3.3 Cramer-Rao Lower

Bound and Efficiency

3.4 Maximum

Likelihood Estimation

3.5 Properties

Estimator

3.5.1 Invariant

3.5.2 MLE is

Asymptotically

Unbiased and Efficient

3.5.3 MLE is consistent

for Zero CRLB

3.6 Expectation

Maximization

3.7 Maximum Entropy

Estimation

Entropy

$$H = - \sum_{k \geq 1} p_k \log p_k$$

$$p_k \log p_k = 0 \text{ for } p_k = 0$$

$$H = - \int_{-\infty}^{\infty} p(x) \log p(x) dx$$

$$p(x) \log p(x) = 0 \text{ for } p(x) = 0$$

Examples:

- normal distribution: given mean and standard deviation
- uniform distribution: supported in the interval $[a, b]$
- exponential distribution: given mean in $[0, \infty]$

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Theorem 1 (Boltzmann's Theorem: Discrete)

Suppose $S = \{x_1, x_2, \dots\}$ is a (finite or infinite) discrete subset of the reals and n functions f_1, \dots, f_n and n numbers a_1, \dots, a_n are given. Let be C the class of all discrete random variables X which are supported on S and which satisfy the n conditions:

$$E(f_j(X)) = a_j \quad \text{for } j = 1, \dots, n$$

If there exists a member of C which assigns positive probability to all members of S and if there exists a maximum entropy distribution for C , then this distribution has the following shape:

$$\Pr(X = x_k) = c \exp \left(\sum_{j=1}^n \lambda_j f_j(x_k) \right) \quad \text{for } k = 1, 2, \dots,$$

where the constants c and λ_j have to be determined so that the sum of the probabilities is 1 and the above conditions for the expected values are satisfied.

Conversely, if constants c and λ_j as above can be found, then the above distribution is indeed the maximum entropy distribution for class C .

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Theorem 1 (Boltzmann's Theorem: Continuous)

Suppose S is a closed subset of the real numbers \mathbb{R} and n measurable functions f_1, \dots, f_n and n numbers a_1, \dots, a_n are given. Let be C the class of all continuous random variables which are supported on S and which satisfy the n expected value conditions:

$$E(f_j(X)) = a_j \quad \text{for } j = 1, \dots, n$$

If there is a member in C whose density function is positive everywhere in S , and if there exists a maximal entropy distribution for C , then its probability density $p(x)$ has the following shape:

$$p(x) = c \exp \left(\sum_{j=1}^n \lambda_j f_j(x) \right) \quad \text{for all } x \in S ,$$

where the constants c and λ_j have to be determined so that the integral of $p(x)$ over S is 1 and the above conditions for the expected values are satisfied.

Conversely, if constants c and λ_j like this can be found, then $p(x)$ is indeed the density of the (unique) maximum entropy distribution for class C .

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk
- Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Not all classes of distributions contain a maximum entropy distribution:

- arbitrarily large entropy: distributions with mean

- entropies of a class are bounded from above but not attained: distributions with mean zero, second moment one, and third moment one

Maximum Entropy: Discrete Solution



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Constraints:
$$\sum_{i=1}^n p(x_i) f_k(x_i) = F_k \quad k = 1, \dots, m$$

$$\sum_{i=1}^n p(x_i) = 1$$

Solution, the **Gibbs distribution**

$$p(x_i) = \frac{1}{Z(\lambda_1, \dots, \lambda_m)} \exp(\lambda_1 f_1(x_i) + \dots + \lambda_m f_m(x_i))$$

with **partition function**

$$Z(\lambda_1, \dots, \lambda_m) = \sum_{i=1}^n \exp(\lambda_1 f_1(x_i) + \dots + \lambda_m f_m(x_i))$$

The **Lagrange multipliers** are determined by the equation system

$$F_k = \frac{\partial}{\partial \lambda_k} \log Z(\lambda_1, \dots, \lambda_m)$$

Maximum Entropy: Continuous Solution



- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

Instead of the entropy, we use the Kullback-Leibler divergence:

$$D_{\text{KL}}(p \parallel m) = - \int p(x) \log \frac{p(x)}{m(x)} dx$$

$m(x)$ is proportional to the limiting density of discrete points and is assumed to be known

$$\text{Constraints: } \int p(x) f_k(x) dx = F_k \quad k = 1, \dots, m \quad \int p(x) dx = 1$$

Solution, is **Gibbs distribution**:

$$p(x) = \frac{1}{Z(\lambda_1, \dots, \lambda_m)} m(x) \exp(\lambda_1 f_1(x) + \dots + \lambda_m f_m(x))$$

with **partition function**

$$Z(\lambda_1, \dots, \lambda_m) = \int m(x) \exp(\lambda_1 f_1(x) + \dots + \lambda_m f_m(x)) dx$$

and the equation system for the **Lagrange multipliers**:

$$F_k = \frac{\partial}{\partial \lambda_k} \log Z(\lambda_1, \dots, \lambda_m)$$

- 1 Introduction
- 2 Generalization Error
 - 2.1 Model Quality
 - 2.2 Gen. Error
 - 2.2.1 Definition
 - 2.2.2 Estimation
 - 2.3 Minimal Risk Example
- 3 Maximum Likelihood
 - 3.1 Loss Unsupervised
 - 3.1.1 Projections
 - 3.1.2 Generative
 - 3.1.3 Parameter Estimation
 - 3.2 MSE, Bias, & Variance
 - 3.3 Cramer-Rao Lower Bound and Efficiency
 - 3.4 Maximum Likelihood Estimation
 - 3.5 Properties Estimator
 - 3.5.1 Invariant
 - 3.5.2 MLE is Asymptotically Unbiased and Efficient
 - 3.5.3 MLE is consistent for Zero CRLB
 - 3.6 Expectation Maximization
 - 3.7 Maximum Entropy Estimation

The invariant measure function $m(x)$ is actually the prior density function representing the lack of relevant information. It can be determined by the principle of transformation groups or marginalization theory.

If x takes values only in (a, b) , then the maximum entropy probability density function is $p(x) = Z \cdot m(x)$, $a < x < b$ where Z is a normalization constant.

Chapter 4

Noise Models

4 Noise Models

4.1 Gaussian Noise

4.2 Laplace Noise

4.3 Binary Models

4.3.1 Cross-Entropy

4.3.2 Logistic

Regression

4.3.3 Log. Regression

Convex

4.3.4 Softmax

4.3.5 Softmax Convex

5 Statistical Learning Theory

5.1 Error Bound

Example

5.2 Empirical Risk

Minimization

5.2.1 Complexity:

Finite Number of

Functions

5.2.2 Complexity: VC-

Dimension

5.3 Error Bounds

5.4 Structural Risk

Minimization

5.5 Margin

5.6 Average Bounds

6 Kernels and Dot

Products

6.1 Mercer's Theorem

6.2 Reproducing

Kernel Hilbert Space

- connecting unsupervised and supervised learning
- quality measure
- noise on the targets
- apply maximum likelihood

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy Regression
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- Gaussian target noise
- linear model

$$\mathbf{s} = \mathbf{X} \mathbf{w}$$

$$\mathbf{y} = \mathbf{s} + \boldsymbol{\epsilon} = \mathbf{X} \mathbf{w} + \boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \propto \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

$$\mathcal{L}((\mathbf{y}, \mathbf{X}); \mathbf{w}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{y} - \mathbf{X} \mathbf{w})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{w})\right)$$

log-likelihood:

$$\ln \mathcal{L}((\mathbf{y}, \mathbf{X}); \mathbf{w}) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{y} - \mathbf{X} \mathbf{w})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{w})$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- minimize $(\mathbf{y} - \mathbf{X} \mathbf{w})^T \Sigma^{-1} (\mathbf{y} - \mathbf{X} \mathbf{w})$

least square criterion

linear least square estimator

derivative with respect to \mathbf{w} :

$$- 2\mathbf{X}^T \Sigma^{-1} \mathbf{y} + 2\mathbf{X}^T \Sigma^{-1} \mathbf{X} \mathbf{w}$$

Setting the derivative to zero (Wiener-Hopf equations):

$$\hat{\mathbf{w}} = (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{y}$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Noise covariance matrix gives the noise for each measure
In most cases we have the same noise for each observation:

$$\Sigma^{-1} = \frac{1}{\sigma} \mathbf{I}$$

We obtain

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$: pseudo inverse or Moore-Penrose inverse

minimal value: $\frac{1}{\sigma} \mathbf{y}^T \left(\mathbf{I} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) \mathbf{y}$

Laplace Noise and Minkowski Error



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Laplace noise assumption:

$$\|y - g(x; w)\|_1$$

$$p(y - g(x; w)) = \frac{\beta}{2} \exp(-\beta |y - g(x; w)|)$$

More general Minkowski error:

$$|y - g(x; w)|^r$$

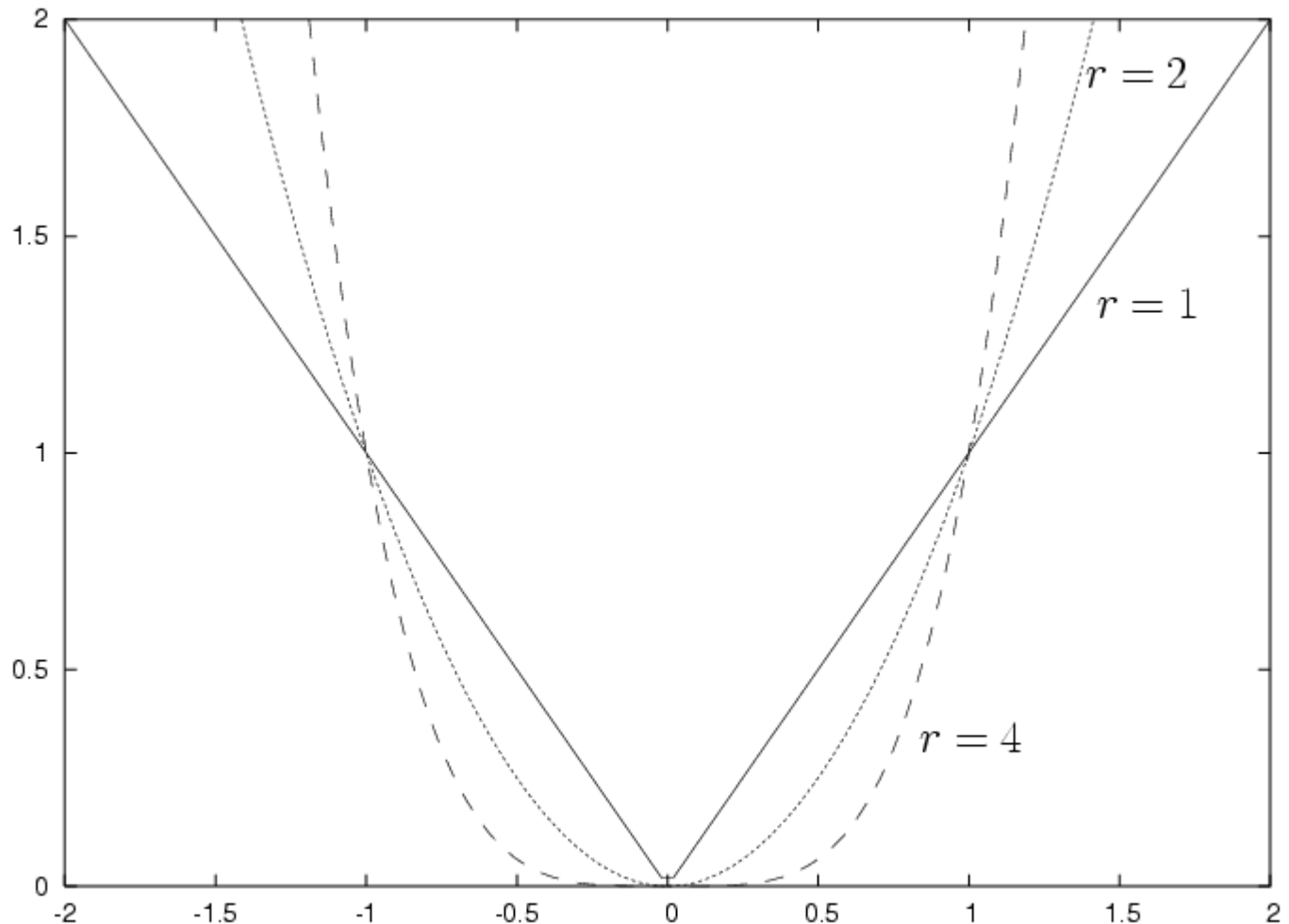
$$p(y - g(x; w)) = \frac{r \beta^{1/r}}{2 \Gamma(1/r)} \exp(-\beta |y - g(x; w)|^r),$$

gamma function

$$\Gamma(a) = \int_0^{\infty} u^{a-1} e^{-u} du, \quad \Gamma(n) = (n-1)!$$

Laplace Noise and Minkowski Error

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space



4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic
Regression
4.3.3 Log. Regression
Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning
Theory
5.1 Error Bound
Example
5.2 Empirical Risk
Minimization
5.2.1 Complexity:
Finite Number of
Functions
5.2.2 Complexity: VC-
Dimension
5.3 Error Bounds
5.4 Structural Risk
Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot
Products
6.1 Mercer's Theorem
6.2 Reproducing
Kernel Hilbert Space

- noise considerations do not hold for binary target
- classification not treated

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

classification problem with K classes:

$$g_k(\mathbf{x}; \mathbf{w}) = p(\mathbf{y} = \mathbf{e}_k \mid \mathbf{x})$$

$$\mathbf{y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_K\}$$

If \mathbf{x} is in the k th class then $\mathbf{y} = (0, \dots, 0, \underset{\substack{\text{position } k \\ \downarrow}}{\mathbf{1}}, 0, \dots, 0)$

Likelihood:

$$\mathcal{L}(\{\mathbf{z}\}; \mathbf{w}) = p(\{\mathbf{z}\}; \mathbf{w}) = \prod_{i=1}^l \prod_{k=1}^K p(\mathbf{y}^i = \mathbf{e}_k \mid \mathbf{x}^i; \mathbf{w})^{[\mathbf{y}^i]_k} p(\mathbf{x}^i)$$

$$\prod_{k=1}^K p(\mathbf{y}^i = \mathbf{e}_k \mid \mathbf{x}^i; \mathbf{w})^{[\mathbf{y}^i]_k} = p(\mathbf{y}^i = \mathbf{e}_r \mid \mathbf{x}^i; \mathbf{w}) \text{ for } \mathbf{y}^i = \mathbf{e}_r$$

Cross-Entropy



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

The log-likelihood:

$$\ln \mathcal{L}(\{z\}; \mathbf{w}) = \sum_{k=1}^K \sum_{i=1}^l [\mathbf{y}^i]_k \ln p(\mathbf{y}^i = \mathbf{e}_k | \mathbf{x}^i; \mathbf{w}) + \sum_{i=1}^l \ln p(\mathbf{x}^i)$$

loss function:

$$\sum_{k=1}^K \sum_{i=1}^l [\mathbf{y}^i]_k \ln p(\mathbf{y}^i = \mathbf{e}_k | \mathbf{x}^i; \mathbf{w})$$

cross entropy (Kullback-Leibler)

$$[\mathbf{y}^i]_k \sim p(\mathbf{y}^i = \mathbf{e}_k)$$

$$[\mathbf{y}^i]_k = \begin{cases} 1 & \text{for } \mathbf{y}^i = \mathbf{e}_k \\ 0 & \text{for } \mathbf{y}^i \neq \mathbf{e}_k \end{cases}$$

$$\sum_{k=1}^K \sum_{i=1}^l [\mathbf{y}^i]_k \ln \frac{p(\mathbf{y}^i = \mathbf{e}_k | \mathbf{x}^i; \mathbf{w})}{[\mathbf{y}^i]_k} - \sum_{k=1}^K \sum_{i=1}^l [\mathbf{y}^i]_k \ln [\mathbf{y}^i]_k$$

Logistic Regression



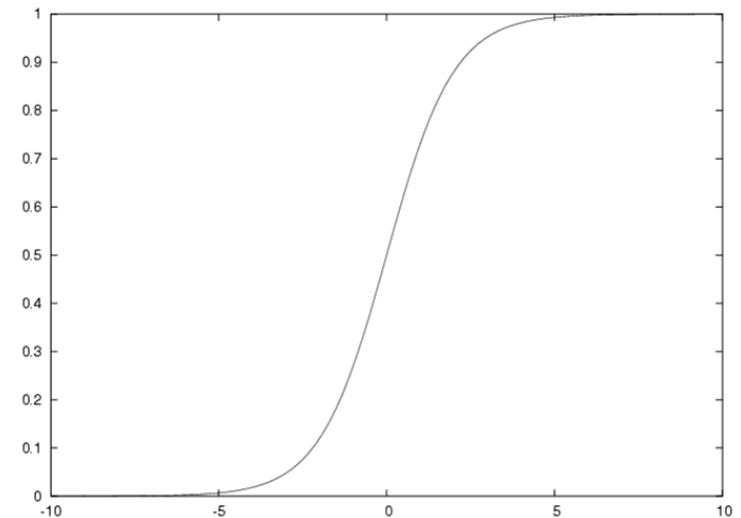
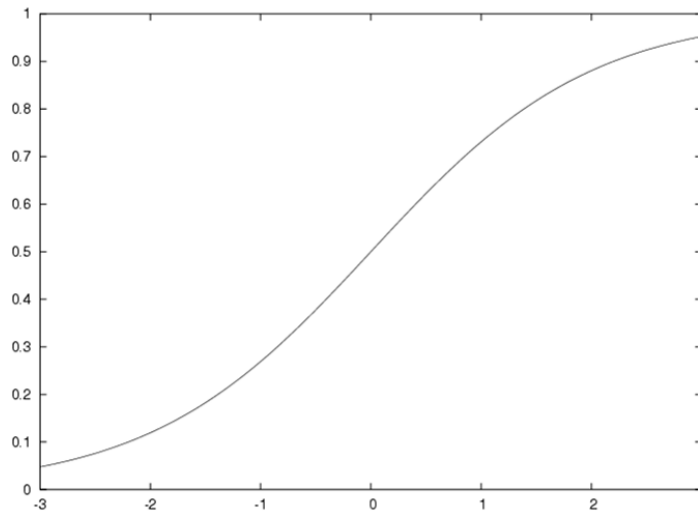
- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

a function g mapping x onto \mathbb{R} can be transformed into a probability:

$$p(y = 1 \mid \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-g(\mathbf{x}; \mathbf{w})}}$$



Logistic Regression



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

If follows:

$$g(\mathbf{x}; \mathbf{w}) = \ln \left(\frac{p(y = 1 | \mathbf{x})}{1 - p(y = 1 | \mathbf{x})} \right) .$$

log-likelihood:

$$\begin{aligned} \ln \mathcal{L}(\{\mathbf{z}\}; \mathbf{w}) &= \sum_{i=1}^l \ln p(\mathbf{z}_i; \mathbf{w}) = \sum_{i=1}^l \ln p(y^i, \mathbf{x}^i; \mathbf{w}) = \\ &= \sum_{i=1}^l \ln p(y^i | \mathbf{x}^i; \mathbf{w}) + \sum_{i=1}^l \ln p(\mathbf{x}^i) \end{aligned}$$

maximum likelihood maximizes

$$\sum_{i=1}^l \ln p(y^i | \mathbf{x}^i; \mathbf{w})$$

Logistic Regression



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

derivative of the log-likelihood:

$$\frac{\partial}{\partial w_j} \sum_{i=1}^l \ln p(y^i | \mathbf{x}^i; \mathbf{w}) = \sum_{i=1}^l (p(y = 1 | \mathbf{x}^i; \mathbf{w}) - y^i) \frac{\partial g(\mathbf{x}^i; \mathbf{w})}{\partial w_j}$$

similar to the derivative of the quadratic loss function in the regression:

$$\frac{1}{2} \frac{\partial}{\partial w_j} \sum_{i=1}^l (g(\mathbf{x}^i; \mathbf{w}) - y^i)^2 = \sum_{i=1}^l (g(\mathbf{x}^i; \mathbf{w}) - y^i) \frac{\partial g(\mathbf{x}^i; \mathbf{w})}{\partial w_j}$$

$(g(\mathbf{x}^i; \mathbf{w}) - y^i)$ instead of $(p(y = 1 | \mathbf{x}^i; \mathbf{w}) - y^i)$

$$p(y = 1 | \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-g(\mathbf{x}; \mathbf{w})}}$$

Logistic Regression



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

linear case: $g(\mathbf{x}^i; \mathbf{w}) = \mathbf{w}^T \mathbf{x}^i$

labels: $y \in \{+1, -1\}$

$$\frac{\partial L}{\partial w_j} = - \sum_{i=1}^l y^i x_{ij} (1 - p(y^i | \mathbf{x}^i; \mathbf{w})) \quad (\text{We drop } i \text{ from } \mathbf{x} !)$$

second order derivatives with $H_{jk} = \frac{\partial L}{\partial w_j \partial w_k} = \sum_{i=1}^l (y^i)^2 x_{ij} x_{ik} p(y^i | \mathbf{x}; \mathbf{w}) (1 - p(y^i | \mathbf{x}; \mathbf{w}))$

H as the Hessian

$$\rho_{ij} = x_{ij} \sqrt{p(y^i | \mathbf{x}; \mathbf{w}) (1 - p(y^i | \mathbf{x}; \mathbf{w}))}$$

bilinear form of the Hessian with a vector:

$$\mathbf{a}^T \mathbf{H} \mathbf{a} = \sum_{i=1}^l \sum_{j=1}^d \sum_{k=1}^d x_{ij} x_{ik} a_j a_k p(y^i | \mathbf{x}; \mathbf{w}) (1 - p(y^i | \mathbf{x}; \mathbf{w})) =$$

$$\sum_{i=1}^l \sum_{j=1}^d a_j x_{ij} \sqrt{p(y^i | \mathbf{x}; \mathbf{w}) (1 - p(y^i | \mathbf{x}; \mathbf{w}))} \sum_{k=1}^d a_k x_{ik} \sqrt{p(y^i | \mathbf{x}; \mathbf{w}) (1 - p(y^i | \mathbf{x}; \mathbf{w}))} =$$

$$\sum_{i=1}^l (\mathbf{a}^T \boldsymbol{\rho}_i) (\mathbf{a}^T \boldsymbol{\rho}_i) = \sum_{i=1}^l (\mathbf{a}^T \boldsymbol{\rho}_i)^2 \geq 0$$

→ Hessian is positive definite

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Multi-class: logistic regression can be generalized by **Softmax**

We assume K classes with $y \in \{1, \dots, K\}$

$$p(y = k \mid \mathbf{x}; g_1, \dots, g_K, \mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{e^{g_k(\mathbf{x}; \mathbf{w}_k)}}{\sum_{j=1}^K e^{g_j(\mathbf{x}; \mathbf{w}_j)}}$$

→ multinomial distribution across the classes

$$\text{Objective: } L = - \sum_{i=1}^l \ln p(y = y^i \mid \mathbf{x}^i; \mathbf{w}) = \sum_{i=1}^l \ln \left(\sum_{j=1}^K e^{g_j(\mathbf{x}^i; \mathbf{w}_j)} \right) - g_{y^i}(\mathbf{x}^i; \mathbf{w}_{y^i})$$

In the following we set $p(y = k \mid \mathbf{x}; g_1, \dots, g_K, \mathbf{w}_1, \dots, \mathbf{w}_K) = p(k \mid \mathbf{x}; \mathbf{W})$

Derivatives:

$$\frac{\partial L}{\partial w_{kn}} = \sum_{i=1}^l \frac{\partial g_k(\mathbf{x}^i; \mathbf{w}_k)}{\partial w_{kn}} p(k \mid \mathbf{x}^i; \mathbf{W}) - \delta_{y^i=k} \sum_{i=1}^l \frac{\partial g_k(\mathbf{x}^i; \mathbf{w}_k)}{\partial w_{kn}}$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

We show that linear Softmax is strictly convex

linear case: $g_k(\mathbf{x}^i; \mathbf{w}_k) = \mathbf{w}_k^T \mathbf{x}^i$ $g(\mathbf{x}^i; \mathbf{W}) = \mathbf{W}^T \mathbf{x}^i$

derivatives: $\frac{\partial L}{\partial w_{kn}} = \sum_{i=1}^l x_{in} p(k | \mathbf{x}^i; \mathbf{W}) - \delta_{y^i=k} \sum_{i=1}^l x_{in}$

second order derivatives:

$$\frac{\partial p(v | \mathbf{x}^i; \mathbf{W})}{\partial w_{vm}} = x_{im} p(k | \mathbf{x}^i; \mathbf{W}) (1 - p(k | \mathbf{x}^i; \mathbf{W}))$$

$$\frac{\partial p(k | \mathbf{x}^i; \mathbf{W})}{\partial w_{vm}} = x_{im} p(k | \mathbf{x}^i; \mathbf{W}) p(v | \mathbf{x}^i; \mathbf{W})$$

$$H_{kn,vm} = \frac{\partial L}{\partial w_{kn} \partial w_{vm}} =$$

$$\sum_{i=1}^l x_{in} x_{im} p(k | \mathbf{x}^i; \mathbf{W}) (\delta_{k=v} (1 - p(k | \mathbf{x}^i; \mathbf{W})) - (1 - \delta_{k=v}) p(v | \mathbf{x}^i; \mathbf{W}))$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

We consider the bilinear form

$$\begin{aligned}
 \mathbf{a}^T \mathbf{H} \mathbf{a} &= \sum_{k,n} \sum_{v,m} \sum_i a_{kn} a_{vm} x_{in} x_{im} p(k | \mathbf{x}^i; \mathbf{W}) (\delta_{k=v} (1 - p(k | \mathbf{x}^i; \mathbf{W})) - (1 - \delta_{k=v}) p(v | \mathbf{x}^i; \mathbf{W})) = \\
 &= \sum_{k,n} \sum_i a_{kn} x_{in} p(k | \mathbf{x}^i; \mathbf{W}) \sum_m x_{im} \left(a_{km} - \sum_v a_{vm} p(v | \mathbf{x}^i; \mathbf{W}) \right) = \\
 &= \sum_i \sum_n x_{in} \sum_k a_{kn} p(k | \mathbf{x}^i; \mathbf{W}) \sum_m x_{im} \left(a_{km} - \sum_v a_{vm} p(v | \mathbf{x}^i; \mathbf{W}) \right) = \\
 &= \sum_i - \left\{ \left(\sum_n x_{in} \sum_k a_{kn} p(k | \mathbf{x}^i; \mathbf{W}) \right) \left(\sum_m x_{im} \sum_v a_{vm} p(v | \mathbf{x}^i; \mathbf{W}) \right) \right\} + \\
 &= \left\{ \sum_n x_{in} \sum_k a_{kn} p(k | \mathbf{x}^i; \mathbf{W}) \sum_m x_{im} a_{km} \right\} = \\
 &= \sum_i - \left\{ \left(\sum_n x_{in} \sum_k a_{kn} p(k | \mathbf{x}^i; \mathbf{W}) \right)^2 \right\} + \left\{ \sum_k p(k | \mathbf{x}^i; \mathbf{W}) \left(\sum_n x_{in} a_{kn} \right) \left(\sum_m x_{im} a_{km} \right) \right\} = \\
 &= \sum_i - \left\{ \left(\sum_n x_{in} \sum_k a_{kn} p(k | \mathbf{x}^i; \mathbf{W}) \right)^2 \right\} + \left\{ \sum_k p(k | \mathbf{x}^i; \mathbf{W}) \left(\sum_n x_{in} a_{kn} \right)^2 \right\}
 \end{aligned}$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

If for each summand of the sum over i

$$\sum_k p(k | \mathbf{x}^i; \mathbf{W}) \left(\sum_n x_{in} a_{kn} \right)^2 - \left(\sum_k p(k | \mathbf{x}^i; \mathbf{W}) \sum_n x_{in} a_{kn} \right)^2 \geq 0$$

holds, then the Hessian is positive semidefinite.

Both sums are an expectation with density $p(k | \mathbf{x}^i; \mathbf{W})$ of the value $\sum_n x_{in} a_{kn}$ or its square.

The second moment minus the expectation squared is the variance. Therefore the left hand side of inequality is a second central moment, which is larger than zero.

→ the Hessian is positive semidefinite

Chapter 5

Statistical Learning Theory

4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory

5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products
6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

- Does learning help for future tasks?
- Explains a model which explains the training data also new data?
- Yes, if complexity is bounded
- VC-dimension as complexity measure
- *statistical learning theory*: bounds for the generalization error (future)
- bounds comprise training error and complexity
- **structural risk minimization** minimizes both terms simultaneously

4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory

5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products
6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

- statistical learning theory:
 - the uniform law of large numbers (empirical risk minimization)
 - complexity constrained models (structural risk minimization)
- error bound on the mean squared error: bias-variance formulation
 - bias is training error = empirical risk
 - variance is model complexity
high complexity \rightarrow more models \rightarrow more solutions \rightarrow large variance

Error Bounds for a Gaussian Classification Task



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

We revisit the Gaussian classification task

$$\begin{aligned} R_{\min} &= \int_{\mathbf{X}} \min\{p(\mathbf{x}, y = -1), p(\mathbf{x}, y = 1)\} d\mathbf{x} \\ &= \int_{\mathbf{X}} \min\{p(\mathbf{x} | y = -1) p(y = -1), p(\mathbf{x} | y = 1) p(y = 1)\} d\mathbf{x} \end{aligned}$$

$$\forall a, b > 0 : \forall 0 \leq \beta \leq 1 : \min\{a, b\} \leq a^\beta b^{1-\beta}$$

$$\begin{aligned} \forall 0 \leq \beta \leq 1 : R_{\min} &\leq (p(y = 1))^\beta (p(y = -1))^{1-\beta} \\ &\int_{\mathbf{X}} (p(\mathbf{x} | y = 1))^\beta (p(\mathbf{x} | y = -1))^{1-\beta} d\mathbf{x} \end{aligned}$$

Error Bounds for a Gaussian Classification Task



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Gaussian assumption:

$$\int_{\mathbf{x}} (p(\mathbf{x} | y = 1))^{\beta} (p(\mathbf{x} | y = -1))^{1-\beta} d\mathbf{x} = \exp(-v(\beta))$$

where

$$v(\beta) = \frac{\beta(1-\beta)}{2} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T (\beta \boldsymbol{\Sigma}_1 + (1-\beta) \boldsymbol{\Sigma}_2)^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \frac{1}{2} \ln \frac{|\beta \boldsymbol{\Sigma}_1 + (1-\beta) \boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|^{\beta} |\boldsymbol{\Sigma}_2|^{1-\beta}}$$

Chernoff bound: maximizing $v(\beta)$ with respect to β

Bhattacharyya bound: $\beta = \frac{1}{2}$

$$v(1/2) = \frac{1}{4} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \frac{1}{2} \ln \frac{|\frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2}|}{\sqrt{|\boldsymbol{\Sigma}_1| |\boldsymbol{\Sigma}_2|}}$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

empirical risk minimization (ERM) principle states:

if the training set is explained by the model then the model generalizes to future examples

restrict the complexity of the model class

**empirical risk minimization (ERM):
minimize error on training set**

Complexity: Finite Number of Functions



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- intuition why complexity matters
- complexity is just the number M of functions in model class
- difference training error (**empirical risk**) and test error (**risk**)

empirical risk:
$$R_{\text{emp}}(g, \mathbf{Z}) = \frac{1}{l} \sum_{i=1}^l L(y^i, g(x^i))$$

finite set of functions $\{g_1, \dots, g_M\}$

worst case (learning chooses unknown function):

$$\max_{j=1, \dots, M} \|R_{\text{emp}}(g_j, l) - R(g_j)\|$$

Complexity: Finite Number of Functions



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

$$p \left(\max_{j=1, \dots, M} \|R_{\text{emp}}(g_j, l) - R(g_j)\| > \epsilon \right) \leq$$

$$\sum_{j=1}^M p \left(\|R_{\text{emp}}(g_j, l) - R(g_j)\| > \epsilon \right) \leq$$

$$M 2 \exp(-2 \epsilon^2 l) = 2 \exp \left(\left(\frac{\ln M}{l} - 2 \epsilon^2 \right) l \right) = \delta$$

$$\max_{j=1, \dots, M} f_j > \epsilon$$

$$f_1 > \epsilon \text{ OR } f_2 > \epsilon \dots$$

$p(a \text{ OR } b) \leq p(a) + p(b)$: union bound

distance of average and expectation: Chernoff inequality

$$\text{(for each } j) p(\mu_l - s > \epsilon) < \exp(-2 \epsilon^2 l)$$

where μ_l is the empirical mean of the true value s for l trials

we obtain complexity term $\epsilon(l, M, \delta) = \sqrt{\frac{\ln M + \ln(2/\delta)}{2l}}$

Complexity: Finite Number of Functions

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

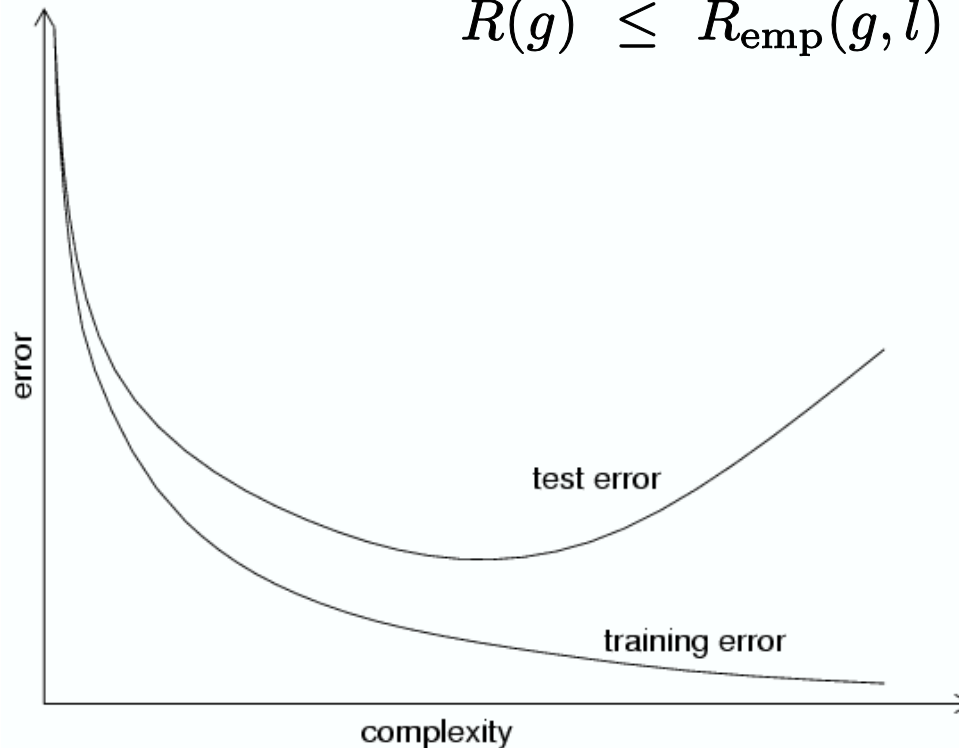
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Theorem 1 (Finite Set Error Bound)

With probability of at least $(1 - \delta)$ over possible training sets with l elements and for M possible functions we have

$$R(g) \leq R_{\text{emp}}(g, l) + \epsilon(l, M, \delta) .$$



Complexity: Finite Number of Functions

4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory

5.1 Error Bound Example
5.2 Empirical Risk Minimization

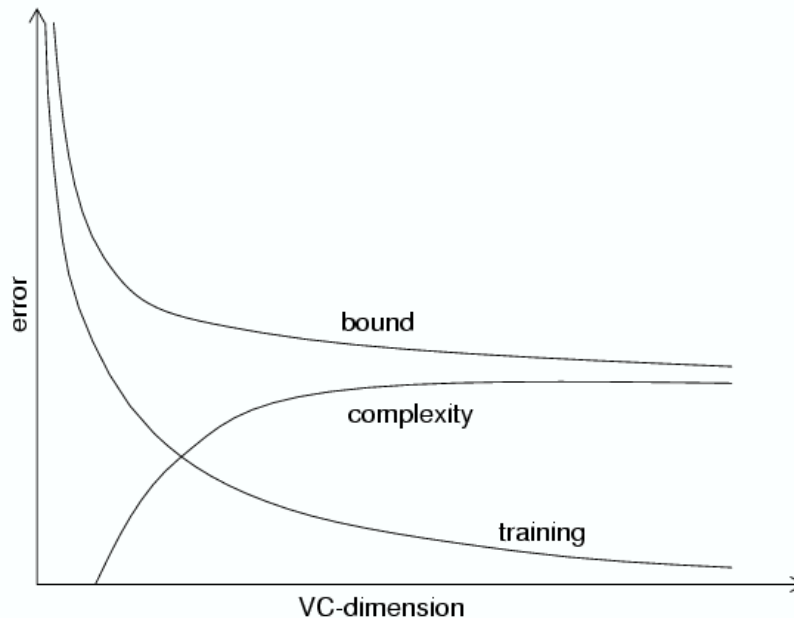
5.2.1 Complexity: Finite Number of Functions

5.2.2 Complexity: VC-Dimension

5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products

6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space



$\epsilon(l, M, \delta)$ should converge to zero as l increases, therefore

$$\frac{\ln M}{l} \xrightarrow{l \rightarrow \infty} 0$$

$$\epsilon(l, M, \delta) = \sqrt{\frac{\ln M + \ln(2/\delta)}{2l}}$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- we want apply the previous bound for infinite function classes
- idea: on training set only finite number of functions is different
- example: all discriminant functions g giving the same classification function $\text{sign } g(\cdot)$
- parametric models $g(\cdot; \mathbf{w})$ with parameter vector \mathbf{w}
- does minimizing the parameter on the training set convergence to the best solution with increasing training set?
- empirical risk minimization (ERM): consistent or not?
- do we select better models with larger training sets?

Complexity: VC-Dimension



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- parameter which minimizes the empirical risk for l training examples:

$$\hat{w}_l = \arg \min_w R_{\text{emp}}(g(\cdot; w), l)$$

- ERM is *consistent* if

$$R(g(\cdot; \hat{w}_l)) \xrightarrow{l \rightarrow \infty} \inf_w R(g(\cdot; w)) \quad (\text{overfitting})$$

convergence in probability

$$R_{\text{emp}}(g(\cdot; \hat{w}_l), l) \xrightarrow{l \rightarrow \infty} \inf_w R(g(\cdot; w)) \quad (\text{underfitting})$$

Empirical risk and expected risk converge to minimal risk

Complexity: VC-Dimension

4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory

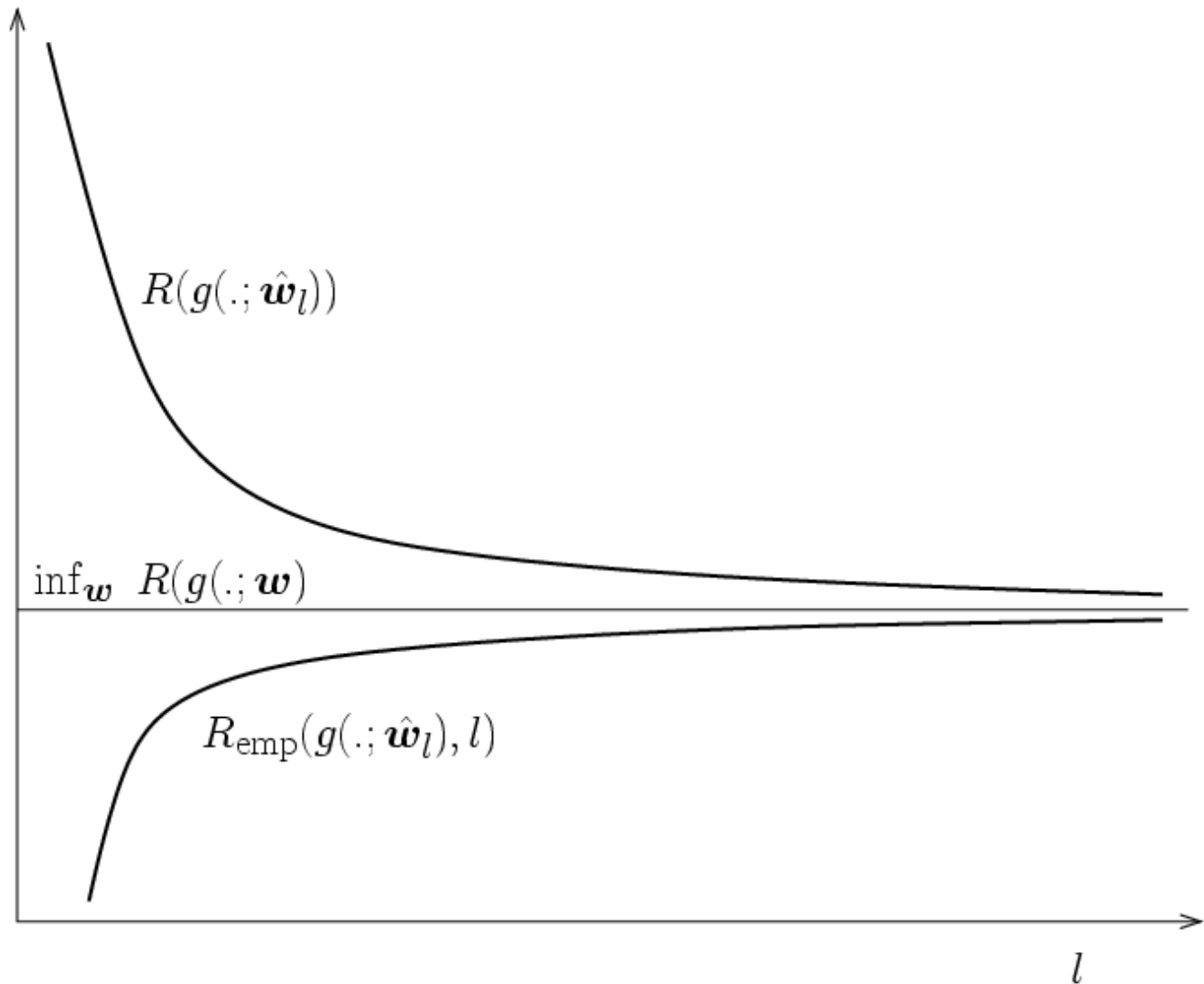
5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions

5.2.2 Complexity: VC-Dimension

5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products

6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space



Complexity: VC-Dimension



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

ERM is **strictly consistent** if for all

$$\Lambda(c) = \left\{ \mathbf{w} \mid z = (\mathbf{x}, y), \int L(y, g(\mathbf{x}; \mathbf{w})) p(\mathbf{z}) d\mathbf{z} \geq c \right\}$$

$$\inf_{\mathbf{w} \in \Lambda(c)} R_{\text{emp}}(g(\cdot; \mathbf{w}), l) \xrightarrow{l \rightarrow \infty} \inf_{\mathbf{w} \in \Lambda(c)} R(g(\cdot; \mathbf{w}))$$

holds (convergence in probability)

Instead of “strictly consistent” we write “consistent”

maximum likelihood is consistent for a set of densities if

$$0 < a \leq p(\mathbf{x}; \mathbf{w}) \leq A < \infty$$

$\exists \mathbf{w}_1 :$

$$\inf_{\mathbf{w}} \frac{1}{l} \sum_{i=1}^l (-\ln p(\mathbf{x}_i; \mathbf{w})) \xrightarrow{l \rightarrow \infty} \inf_{\mathbf{w}} \int_X (-\ln p(\mathbf{x}; \mathbf{w})) p(\mathbf{x}; \mathbf{w}_1) d\mathbf{x}$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- Under what conditions is the ERM consistent?
- New concepts and new capacity measures:
 - points to be shattered
 - annealed entropy
 - entropy (new definition)
 - growth function
 - VC-dimension

Possibilities to label the input data x^i by binary labels $y^i \in \{-1, 1\}$
 $2^l \rightarrow$ **shattering** the input data

complexity of a model class: number different labelings
how many points can be shattered

Complexity: VC-Dimension

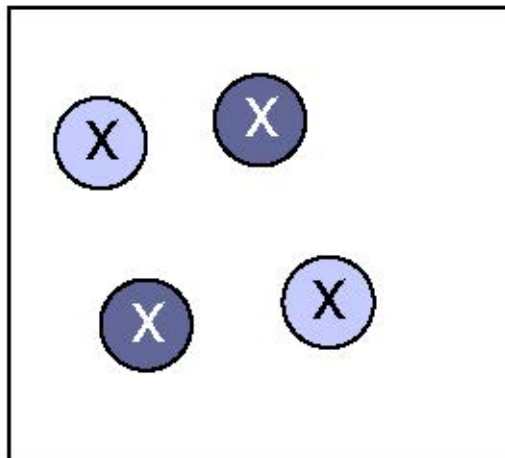
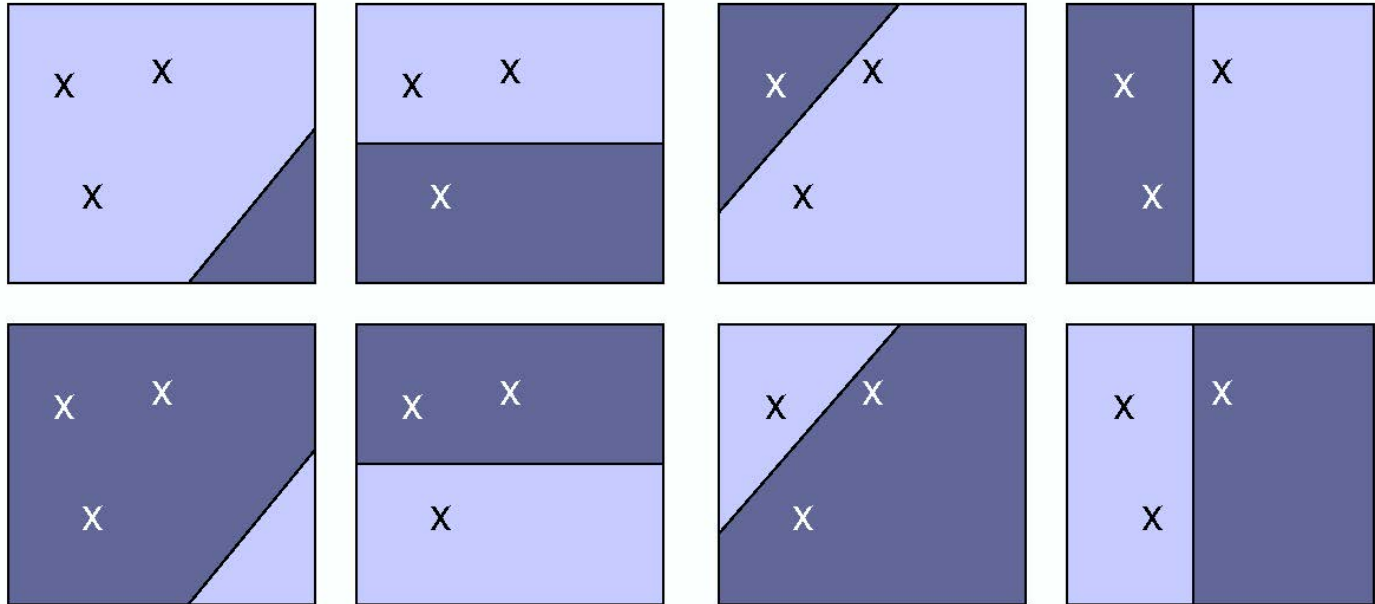
4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic
Regression
4.3.3 Log. Regression
Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning
Theory
5.1 Error Bound
Example
5.2 Empirical Risk
Minimization
5.2.1 Complexity:
Finite Number of
Functions

5.2.2 Complexity: VC-
Dimension

5.3 Error Bounds
5.4 Structural Risk
Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot
Products
6.1 Mercer's Theorem
6.2 Reproducing
Kernel Hilbert Space



Note, that each “x” is placed in a circle around its position independent of the other “x”. Therefore each constellation represents a set with non-zero probability mass.

Complexity: VC-Dimension



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- number of points a function class can shatter: *VC-dimension* (later)
- function class \mathcal{F}
- **shattering coefficient:** $N_{\mathcal{F}}(\mathbf{x}^1, \dots, \mathbf{x}^l)$ (# labeling class can shatter)
- **entropy of a function class:** $H_{\mathcal{F}}(l) = \mathbb{E}_{(\mathbf{x}^1, \dots, \mathbf{x}^l)} \ln N_{\mathcal{F}}(\mathbf{x}^1, \dots, \mathbf{x}^l)$

- **annealed entropy of a function class:**

$$H_{\mathcal{F}}^{\text{ann}}(l) = \ln \mathbb{E}_{(\mathbf{x}^1, \dots, \mathbf{x}^l)} N_{\mathcal{F}}(\mathbf{x}^1, \dots, \mathbf{x}^l)$$

- **growth function of a function class:**

$$G_{\mathcal{F}}(l) = \ln \sup_{(\mathbf{x}^1, \dots, \mathbf{x}^l)} N_{\mathcal{F}}(\mathbf{x}^1, \dots, \mathbf{x}^l)$$

$$H_{\mathcal{F}}(l) \underset{\text{Jensen}}{\leq} H_{\mathcal{F}}^{\text{ann}}(l) \underset{\text{supremum}}{\leq} G_{\mathcal{F}}(l)$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

ERM is **consistent**:

Theorem 1 (Sufficient Condition for Consistency of ERM)

If

$$\lim_{l \rightarrow \infty} \frac{H_{\mathcal{F}}(l)}{l} = 0$$

then ERM is consistent.

ERM has **fast rate of convergence** (exponential convergence):

$$p \left(\sup_{\mathbf{w}} |R(g(\cdot; \mathbf{w})) - R_{\text{emp}}(g(\cdot; \hat{\mathbf{w}}_l), l)| > \epsilon \right) < b \exp(-c \epsilon^2 l)$$

Theorem 1 (Sufficient Condition for Fast Rate)

If

$$\lim_{l \rightarrow \infty} \frac{H_{\mathcal{F}}^{\text{ann}}(l)}{l} = 0$$

then ERM has a fast rate of convergence.

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

theorems valid for a given probability measure on the observations
probability measure enters the formulas via the expectation

Theorem 1 (Consistency of ERM for Any Probability)

The condition

$$\lim_{l \rightarrow \infty} \frac{G_{\mathcal{F}}(l)}{l} = 0$$

*is necessary and sufficient for the ERM to be consistent
and also ensures a fast rate of convergence.*

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- VC (Vapnik-Chervonenkis) dimension d_{VC} is the largest integer for which $G_{\mathcal{F}}(l) = l \ln 2$ holds

$$d_{VC} = \max_l \{l \mid G_{\mathcal{F}}(l) = l \ln 2\} .$$

If the maximum does not exist: $d_{VC} = \infty$

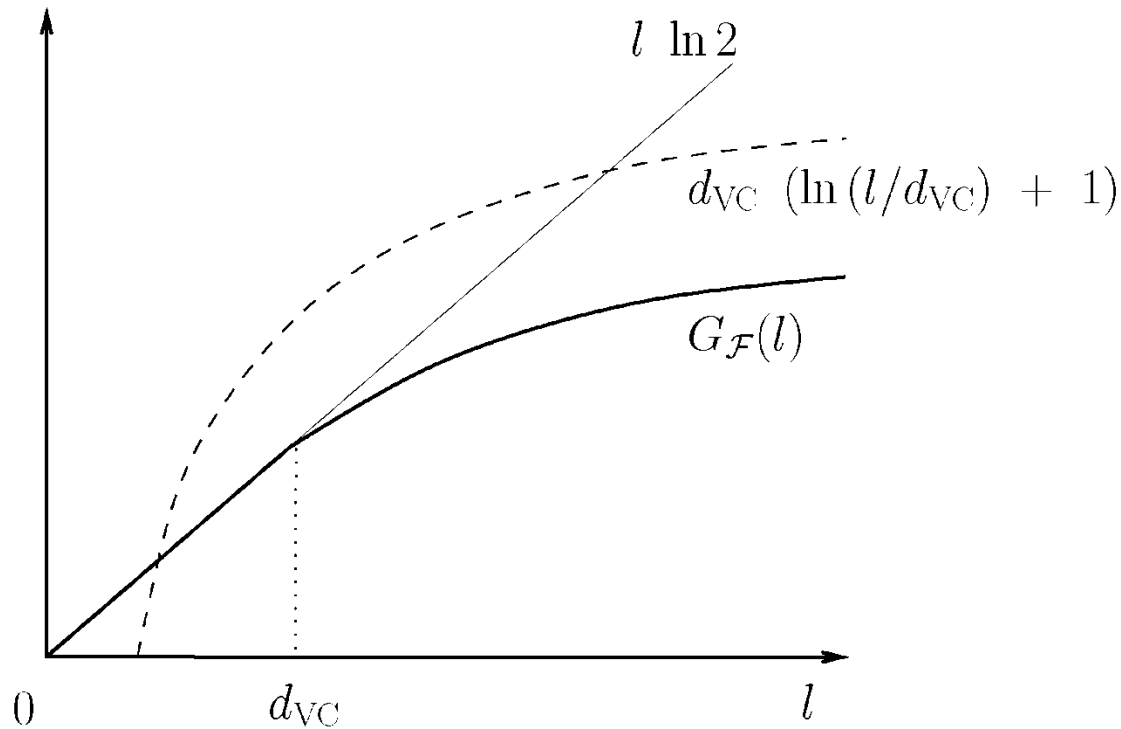
- VC-dimension is the maximum number of vectors that can be shattered by the function class

Complexity: VC-Dimension

Theorem 1 (VC-Dimension Bounds the Growth Function)

The growth function is bounded by

$$G_{\mathcal{F}}(l) \begin{cases} = l \ln 2 & \text{if } l \leq d_{\text{VC}} \\ \leq d_{\text{VC}} \left(1 + \ln \frac{l}{d_{\text{VC}}}\right) & \text{if } l > d_{\text{VC}} \end{cases} .$$



4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory
5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products
6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

Complexity: VC-Dimension



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- function class with finite VC-dim.: consistent and converges fast

- Linear functions in d -dimensional of the input space:

$$g(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} \quad d_{VC} = d$$

$$g(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} + b \quad d_{VC} = d + 1$$

- Nondecreasing nonlinear one-dimensional functions

$$\sum_{i=1}^k |a_i x^i| \operatorname{sign} x + a_0 \quad d_{VC} = 1$$

- Nonlinear one-dimensional functions:

$$\sin(\omega z) \quad d_{VC} = \infty$$

Complexity: VC-Dimension



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- Neural Networks:

$$d_{VC} \leq 2 W \log_2(e M)$$

M are the number of units, W is the number of weights, e is the base of the natural logarithm (Baum & Haussler 89, Shawe-Taylor & Anthony 91)

$$d_{VC} \leq 2 W \log_2(24 e W D)$$

inputs restricted to $[-D; D]$

Bartlett & Williamson (1996)

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- idea of deriving the error bounds: set of distinguishable functions
cardinality given by $N_{\mathcal{F}}$

- trick of two half-samples and their difference (“symmetrization”):

$$p \left(\sup_{\mathbf{w}} \left| \frac{1}{l} \sum_{i=1}^l L(y^i, g(\mathbf{x}^i; \mathbf{w})) - \frac{1}{l} \sum_{i=l+1}^{2l} L(y^i, g(\mathbf{x}^i; \mathbf{w})) \right| > \epsilon - \frac{1}{l} \right) \geq \frac{1}{2} p \left(\sup_{\mathbf{w}} \left| \frac{1}{l} \sum_{i=1}^l L(y^i, g(\mathbf{x}^i; \mathbf{w})) - R(g(\cdot; \mathbf{w})) \right| > \epsilon \right)$$

therefore in the following we use $2l$

→ l example used for complexity definition and l for empirical error

minimal possible risk:

$$\mathbf{w}_0 = \arg \min_{\mathbf{w}} R(g(\cdot; \mathbf{w}))$$

$$R_{\min} = \min_{\mathbf{w}} R(g(\cdot; \mathbf{w})) = R(g(\cdot; \mathbf{w}_0))$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Theorem 1 (Error Bound)

With probability of at least $(1 - \delta)$ over possible training sets with l elements, the parameter \mathbf{w}_l which minimizes the empirical risk we have

$$R(g(\cdot; \mathbf{w}_l)) \leq R_{\text{emp}}(g(\cdot; \mathbf{w}_l), l) + \sqrt{\epsilon(l, \delta)}.$$

With probability of at least $(1 - 2\delta)$ the difference between the optimal risk and the risk of \mathbf{w}_l is bounded by

$$R(g(\cdot; \mathbf{w}_l)) - R_{\min} < \sqrt{\epsilon(l, \delta)} + \sqrt{\frac{-\ln \delta}{l}}$$

Here $\epsilon(l, \delta)$ can be defined for a specific probability as

$$\epsilon(l, \delta) = \frac{8}{l} (H_{\mathcal{F}}^{\text{ann}}(2l) + \ln(4/\delta))$$

or for any probability as

$$\epsilon(l, \delta) = \frac{8}{l} (G_{\mathcal{F}}(2l) + \ln(4/\delta))$$

where the later can be expressed through the VC-dimension d_{VC}

$$\epsilon(l, \delta) = \frac{8}{l} (d_{\text{VC}} (\ln(2l/d_{\text{VC}}) + 1) + \ln(4/\delta))$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- complexity measure depend on the ratio $\frac{d_{VC}}{l}$
- The bound above is from Anthony and Bartlett whereas an older bound from Vapnik is

$$R(g(\cdot; \mathbf{w}_l)) \leq R_{\text{emp}}(g(\cdot; \mathbf{w}_l), l) + \frac{\epsilon(l, \delta)}{2} \left(1 + \sqrt{1 + \frac{R_{\text{emp}}(g(\cdot; \mathbf{w}_l), l)}{\epsilon(l, \delta)}} \right)$$

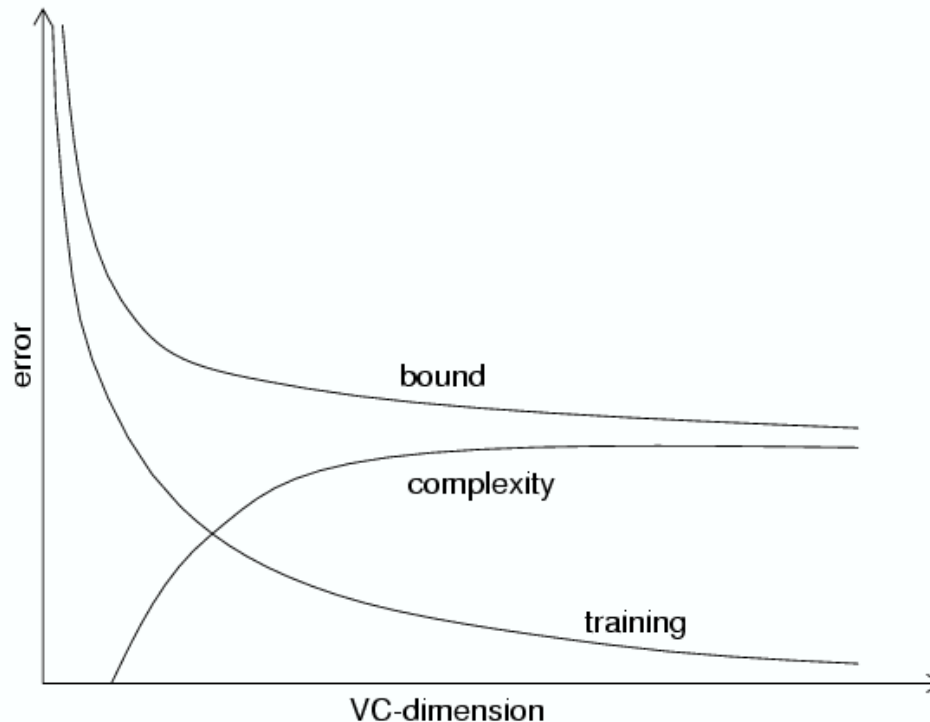
- complexity term decreases with $\frac{1}{\sqrt{l}}$
- zero empirical risk then the bound on the risk decreases with $\frac{1}{\sqrt{l}}$
- later: *expected* risk decreases with $\frac{1}{l}$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- bound on the risk $R \leq R_{\text{emp}} + \text{complexity}$
- bound is similar to the bias-variance formulation
 - bias corresponds to empirical risk
 - variance corresponds to complexity



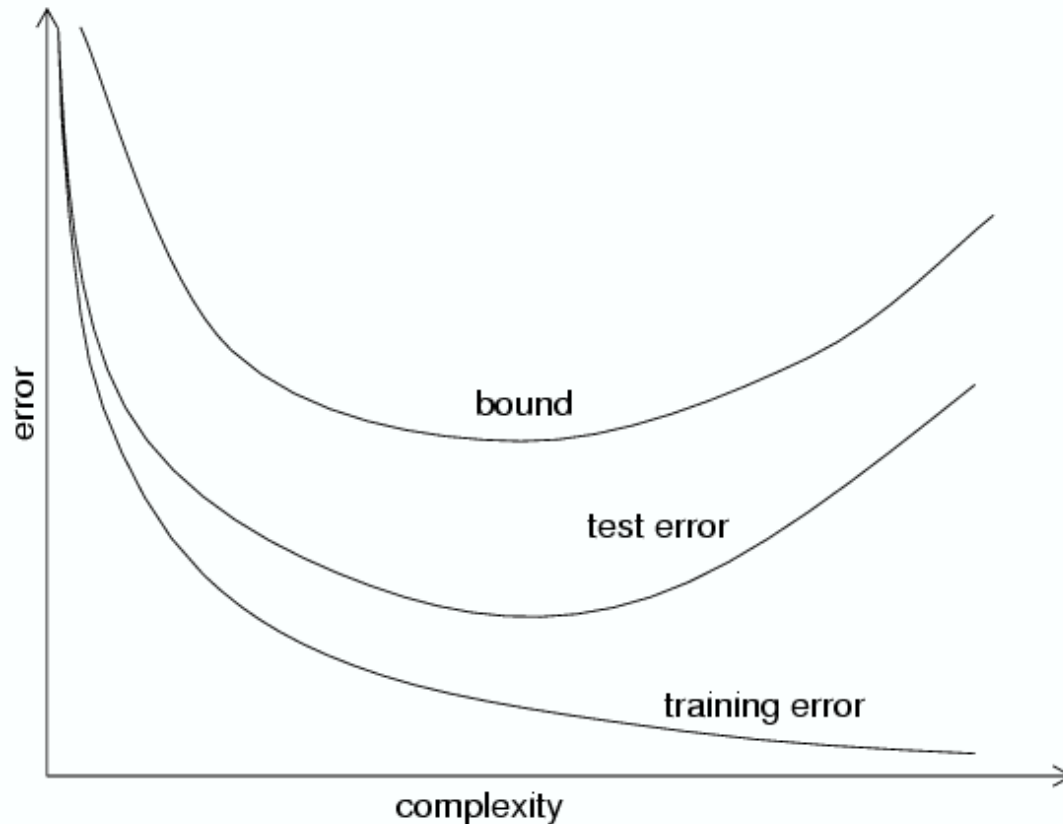
Error Bounds

4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory
5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products
6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

- In many practical cases the bound is not useful: not tight
- However in many practical cases the minimum of the bound is close to the minimum of the test error



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- regression: instead of the shattering coefficient **covering number** (covering of the functions with distance epsilon)

- growth function is then: $G(\varepsilon, \mathcal{F}, l) = \ln \sup_{\mathbf{X}} \mathcal{N}(\varepsilon, \mathcal{F}, \mathbf{X}_\infty)$

bounds on the generalization error:

$$R(g(\cdot; \mathbf{w}_l)) \leq R_{\text{emp}}(g(\cdot; \mathbf{w}_l), l) + \sqrt{\epsilon(\varepsilon, l, \delta)}$$

where

$$\epsilon(\varepsilon, l, \delta) = \frac{36}{l} (\ln(12 l) + G(\varepsilon/6, \mathcal{F}, l) - \ln \delta)$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

The Structural Risk Minimization (SRM) principle minimizes the guaranteed risk that is a bound on the risk instead of the empirical risk alone

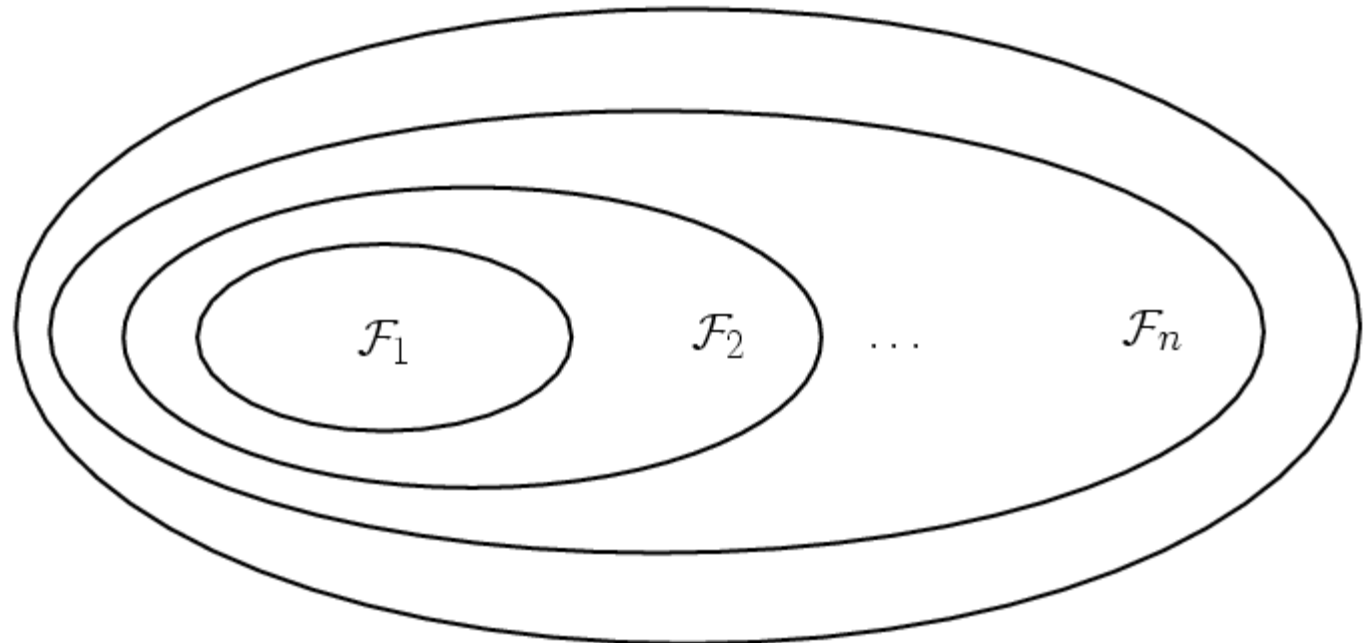
Structural Risk Minimization

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- nested set of function classes: $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}_n \subset \dots$
where class \mathcal{F}_n possesses VC-dimension d_{VC}^n
and $d_{\text{VC}}^1 \leq d_{\text{VC}}^2 \leq \dots \leq d_{\text{VC}}^n \leq \dots$



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- Example for SRM: **minimum description length**

- sender transmits a model (once) and the inputs and errors
- receiver has to recover the labels

goal: minimize transmission costs (description length)

$$\text{transmission costs} = R_{\text{emp}} + \text{complexity}$$

- Is the SRM principle consistent? How fast does it converge?

SRM is consistent !!

asymptotic rate of convergence:

$$r(l) = |R_{\min}^n - R_{\min}| + \sqrt{\frac{d_{\text{VC}}^n \ln l}{l}},$$

where R_{\min}^n is the minimal risk of the function class \mathcal{F}_n

$$p \left(\lim_{l \rightarrow \infty} \sup r^{-1}(l) \left| R \left(g \left(\cdot ; \mathbf{w}_l^{\mathcal{F}_n} \right) \right) - R_{\min} \right| < \infty \right) = 1$$

Structural Risk Minimization



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- $|R_{\min}^n - R_{\min}| \xrightarrow{l \rightarrow \infty} 0$

If the optimal solution belongs to some class \mathcal{F}_n then the convergence rate is

$$r(l) = O\left(\sqrt{\frac{\ln l}{l}}\right)$$

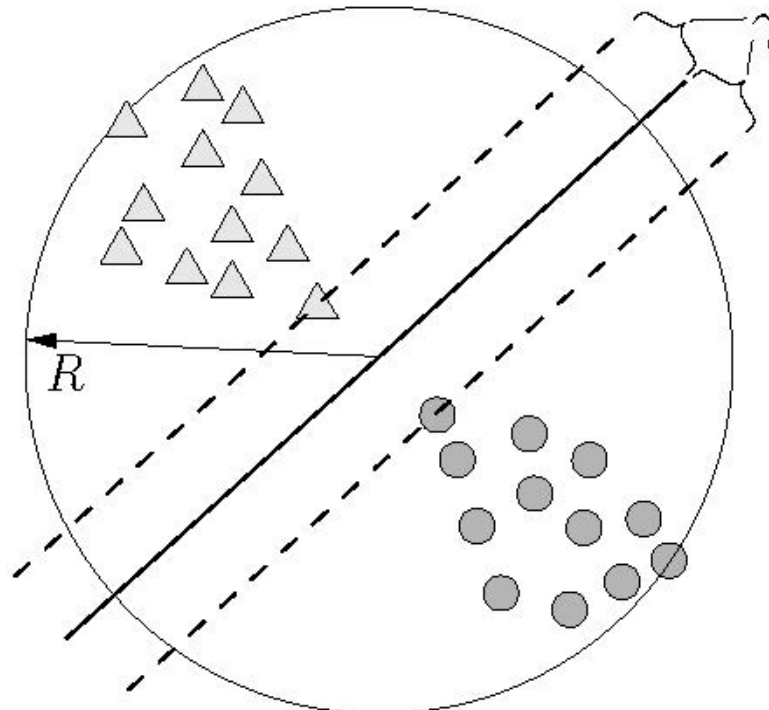
Margin as Complexity Measure

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

- VC-dimension: restrictions on the class of function
- most famous:
 - zero isoline of the discriminant function has minimal distance γ (**margin**) to all training data points which are contained in a sphere with radius R

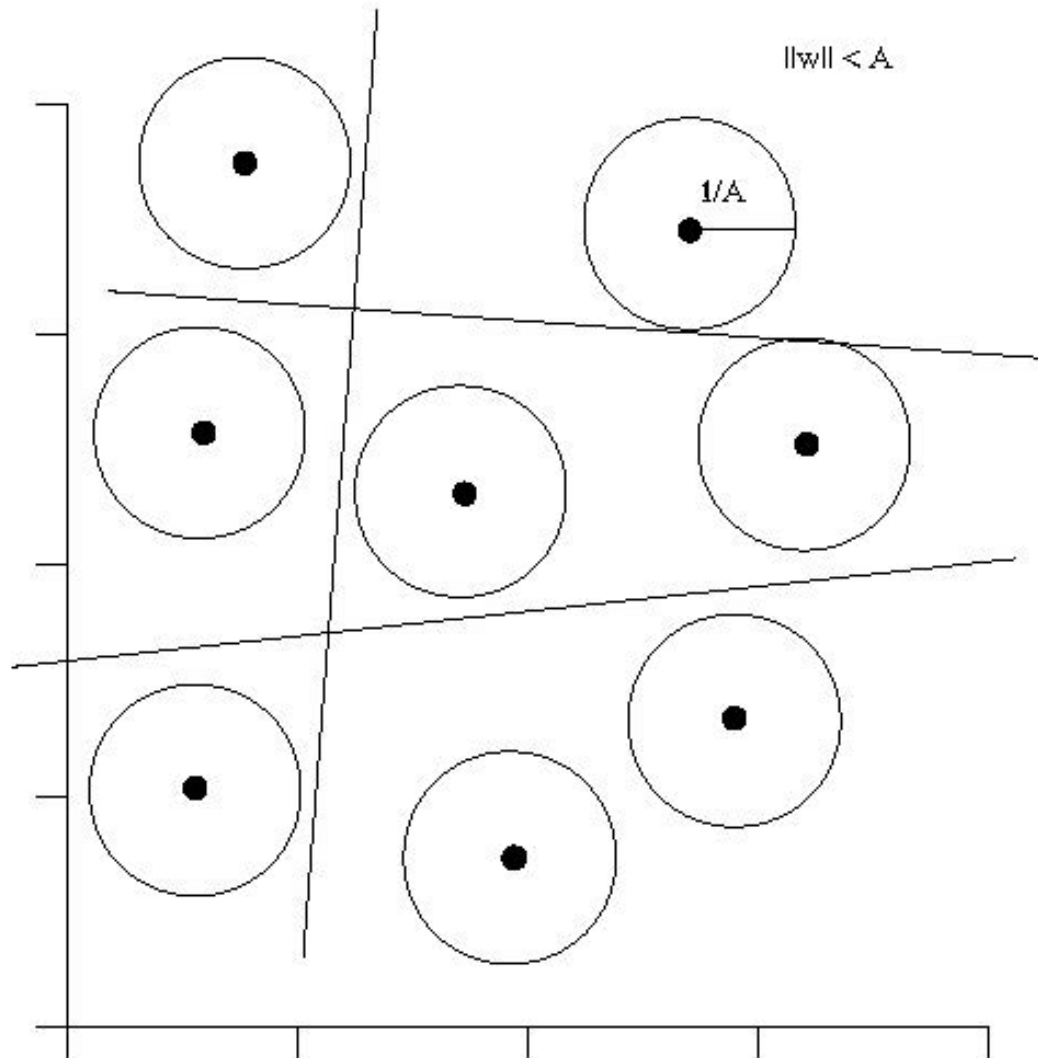


Margin as Complexity Measure

- 4 Noise Models
- 4.1 Gaussian Noise
- 4.2 Laplace Noise
- 4.3 Binary Models
- 4.3.1 Cross-Entropy
- 4.3.2 Logistic Regression
- 4.3.3 Log. Regression Convex
- 4.3.4 Softmax
- 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
- 5.1 Error Bound Example
- 5.2 Empirical Risk Minimization
- 5.2.1 Complexity: Finite Number of Functions
- 5.2.2 Complexity: VC-Dimension
- 5.3 Error Bounds
- 5.4 Structural Risk Minimization
- 5.5 Margin
- 5.6 Average Bounds

- 6 Kernels and Dot Products
- 6.1 Mercer's Theorem
- 6.2 Reproducing Kernel Hilbert Space



Margin as Complexity Measure



4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory
5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products
6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

- linear discriminant functions $w^T x + b$
- classification function $\text{sign}(w^T x + b)$
- scaling w and b does not change classification function
- classification function: one representative discriminant function
- **canonical form** w.r.t. the training data \mathbf{X} :

$$\min_{i=1, \dots, l} |w^T x^i + b| = 1$$

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Theorem 1 (Margin bounds VC-dimension)

The class of classification functions $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$, where the discriminant function $\mathbf{w}^T \mathbf{x} + b$ is in its canonical form versus \mathbf{X} which is contained in a sphere of radius R , and where $\|\mathbf{w}\| \leq \frac{1}{\gamma}$ satisfy

$$d_{\text{VC}} \leq \frac{R^2}{\gamma^2} .$$

Margin as Complexity Measure



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

If at least one data point exists for which the discriminant function is positive and at least one data point exists for which it is negative, then we can optimize b and rescale $\|w\|$ in order to obtain the smallest $\|w\|$

This gives the tightest bound and smallest VC-dimension

After optimizing b and rescaling w we have points for which

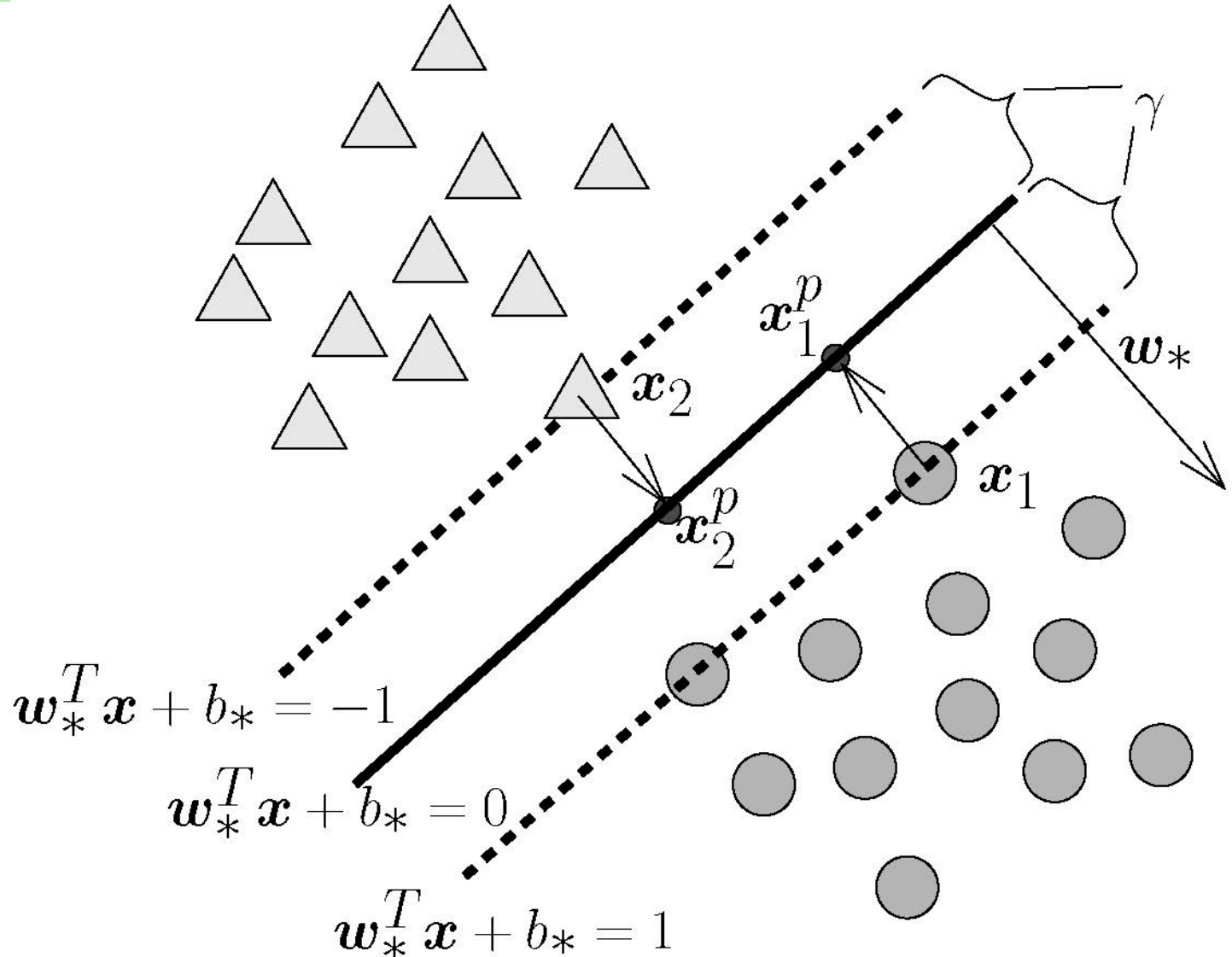
$$w^T x^1 + b = 1 \quad w^T x^2 + b = -1$$

Margin as Complexity Measure

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy Regression
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space



Margin as Complexity Measure



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy Regression
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

After this optimization: the distance of x^1 and x^2 to the boundary function is $\frac{1}{\|w_*\|} = \gamma$

Theorem 1 (Margin Error Bound)

The classification functions $\text{sign}(w^T x + b)$ are restricted to $\|w\| \leq \frac{1}{\gamma}$ and $\|x\| < R$. Let ν be the fraction of training examples which have a margin (distance to $w^T x + b = 0$) smaller than $\frac{\rho}{\|w\|}$.

With probability at least of $(1 - \delta)$ of drawing l examples, the probability to misclassify a new example is bounded from above by

$$\nu + \sqrt{\frac{c}{l} \left(\frac{R^2}{\rho^2 \gamma^2} \ln^2 l + \ln(1/\delta) \right)},$$

where c is a constant.

Average Error Bounds for SVMs

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

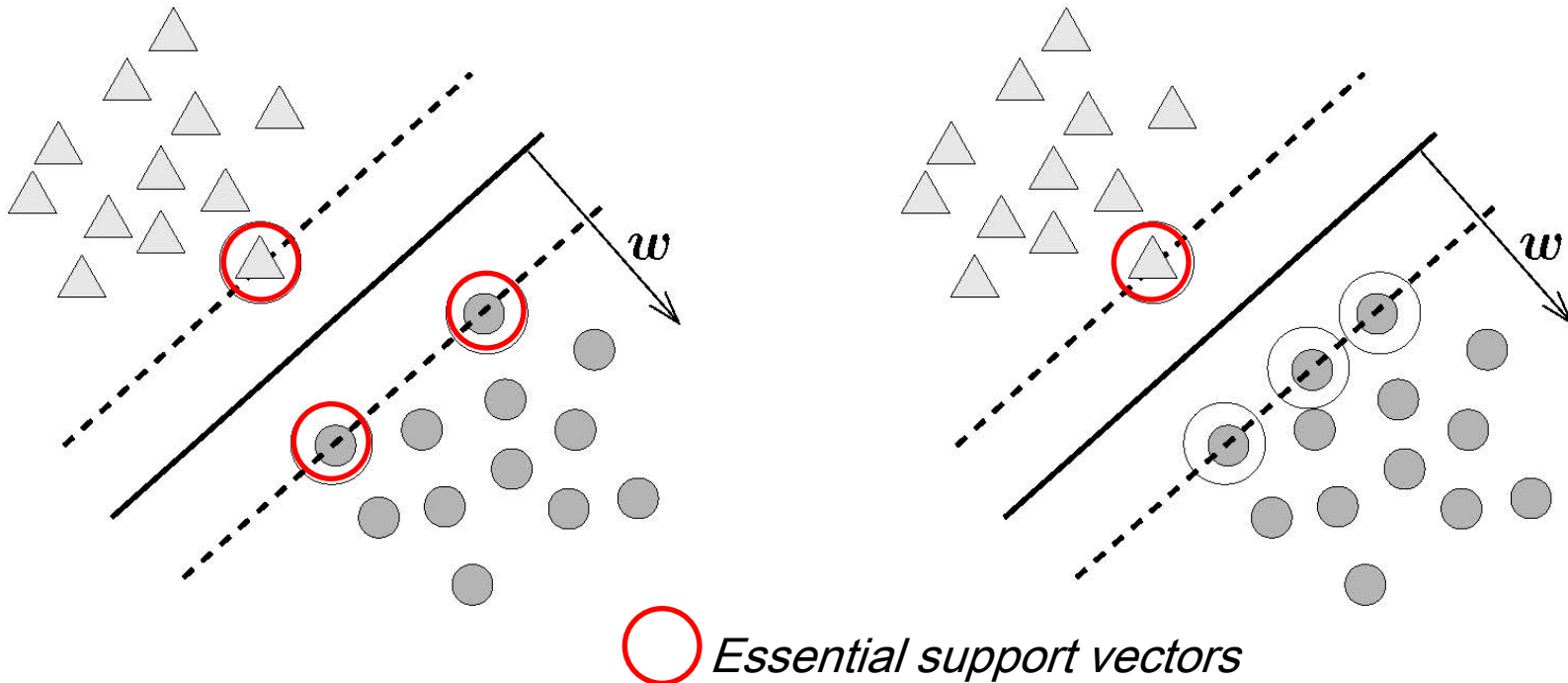
5 Statistical Learning Theory

- 5.1 Error Bound Example
- 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
- 5.3 Error Bounds
- 5.4 Structural Risk Minimization
- 5.5 Margin
- 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Leave-One-Out Cross Validation (LOO CV) estimator is almost unbiased (see Section 2.2.2).

Essential support vectors: support vectors for which the solution changes if they are removed from the training set



Average Error Bounds for SVMs



4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory

5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products

6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

number of essential support vectors: $k_l \leq d + 1$

d is the dimension of the space

radius of the sphere which contains all essential support vectors: r_l

expected risk $ER(g(., w_l))$, where the expectation is taken over the training set of size l and a test data point

Average Error Bounds for SVMs



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds

- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Theorem 1 (Average Bounds for SVMs)

For the expected risk with above definitions

$$ER(g(.; \mathbf{w}_l)) \leq \frac{E k_{l+1}}{l + 1} \quad ER(g(.; \mathbf{w}_l)) \leq \frac{d + 1}{l + 1}$$

$$ER(g(.; \mathbf{w}_l)) \leq \frac{E \left(\frac{r_{l+1}}{\gamma_{l+1}} \right)^2}{l + 1}$$

$$ER(g(.; \mathbf{w}_l)) \leq \frac{E \min \left\{ k_{l+1}, \left(\frac{r_{l+1}}{\gamma_{l+1}} \right)^2 \right\}}{l + 1}$$

$$ER(g(.; \mathbf{w}_l)) \leq \frac{E \left((k_{l+1}^*)^2 \sum_{i^*} \alpha_{i^*} + m \right)}{l + 1}$$

$$C \leq r_l^{-2} : ER(g(.; \mathbf{w}_l)) \leq \frac{E \left((k_{l+1})^2 \sum_i \alpha_i \right)}{l + 1}$$

i^* are SVs with $0 < \alpha_{i^*} < C$ and m is number SVs with $\alpha_i = C$

Average Error Bounds for SVMs



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

bounds are based on the fact that the leave-one-out cross validation is almost unbiased $\rightarrow \frac{1}{l+1}$

these average bounds are of the order $\frac{1}{l}$

worst case bounds (Vapnik bounds) are of the order $\frac{1}{\sqrt{l}}$

variance of the expected risk?

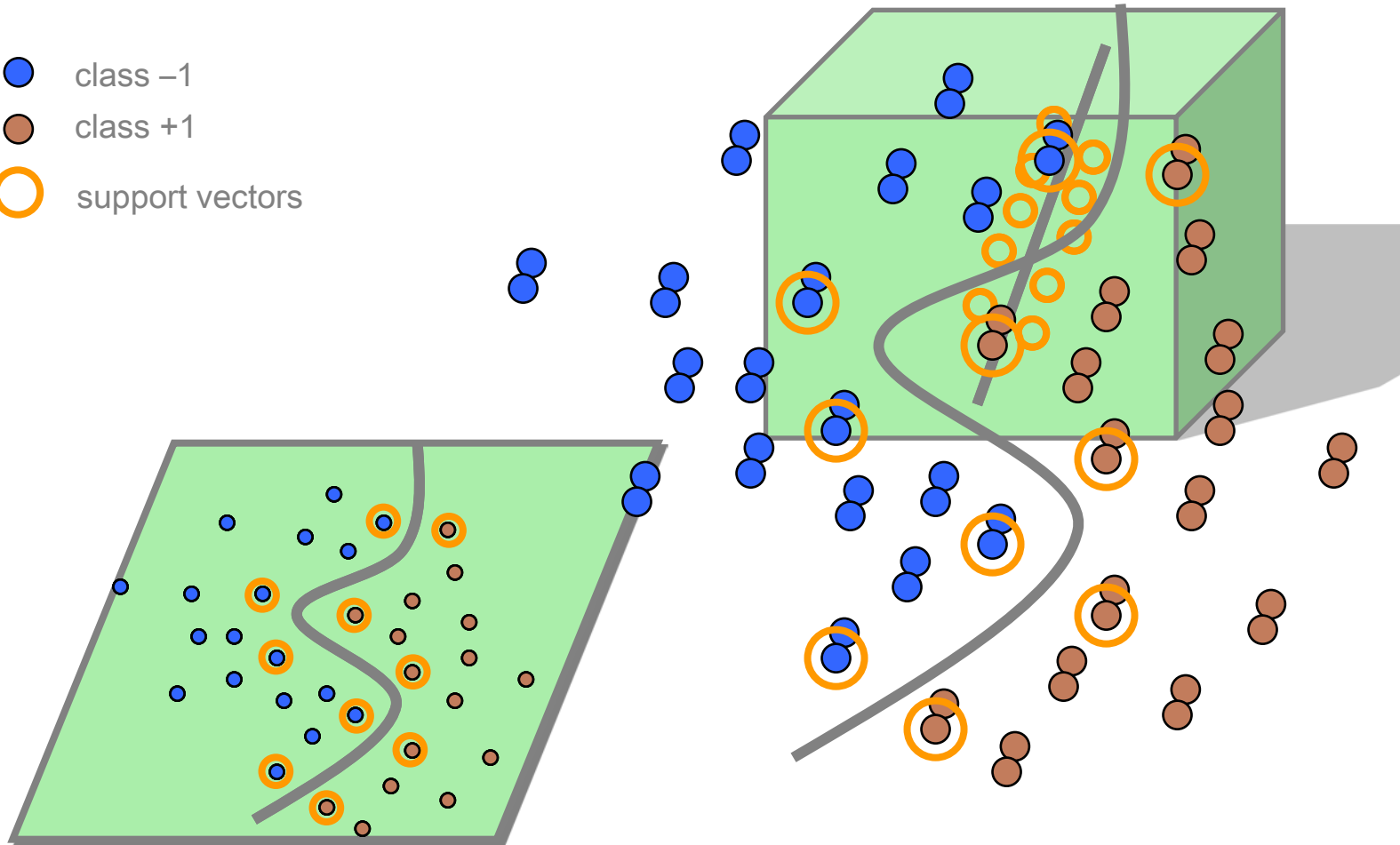
Chapter 6

Theory of Kernels and Dot Products

Support Vector Machine

non-linear support vector machine

- class -1
- class +1
- support vectors



4 Noise Models
4.1 Gaussian Noise
4.2 Laplace Noise
4.3 Binary Models
4.3.1 Cross-Entropy
4.3.2 Logistic Regression
4.3.3 Log. Regression Convex
4.3.4 Softmax
4.3.5 Softmax Convex

5 Statistical Learning Theory
5.1 Error Bound Example
5.2 Empirical Risk Minimization
5.2.1 Complexity: Finite Number of Functions
5.2.2 Complexity: VC-Dimension
5.3 Error Bounds
5.4 Structural Risk Minimization
5.5 Margin
5.6 Average Bounds

6 Kernels and Dot Products
6.1 Mercer's Theorem
6.2 Reproducing Kernel Hilbert Space

Support Vector Machine

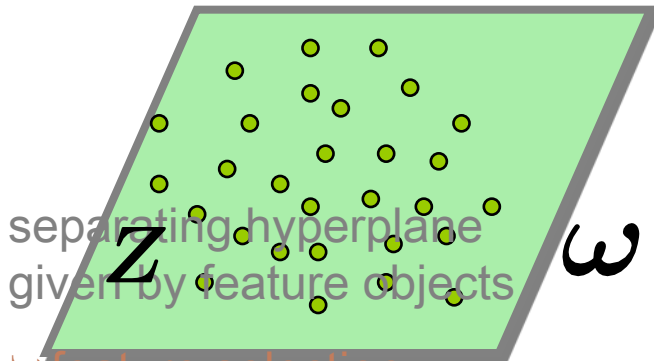
- 4 Noise Models
- 4.1 Gaussian Noise
- 4.2 Laplace Noise
- 4.3 Binary Models
- 4.3.1 Cross-Entropy
- 4.3.2 Logistic Regression
- 4.3.3 Log. Regression Convex
- 4.3.4 Softmax
- 4.3.5 Softmax Convex

- 5 Statistical Learning Theory
- 5.1 Error Bound Example
- 5.2 Empirical Risk Minimization

- 5.2.1 Complexity: Finite Number of Functions
- 5.2.2 Complexity: VC-Dimension
- 5.3 Error Bounds
- 5.4 Structural Risk Minimization
- 5.5 Margin
- 5.6 Average Bounds

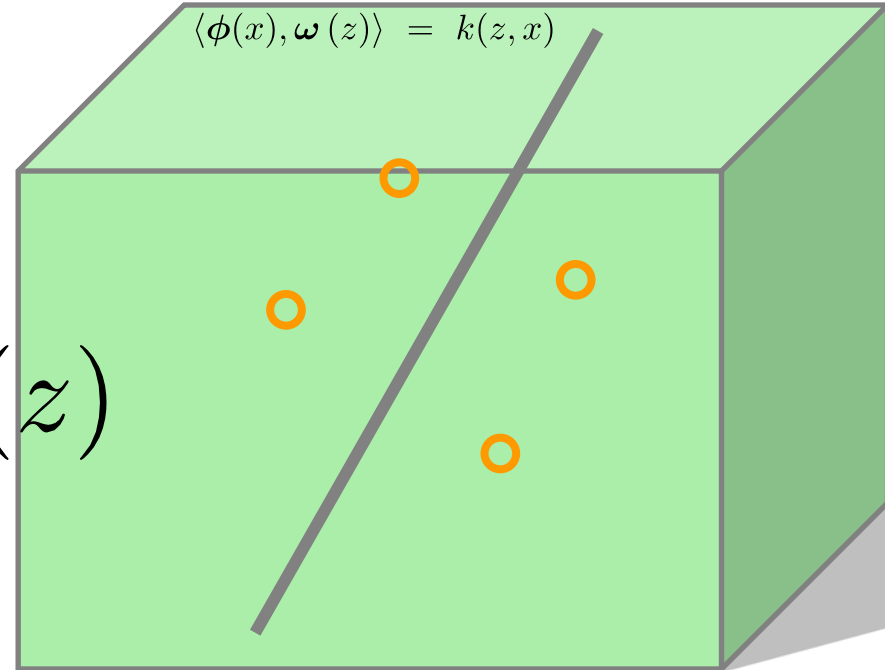
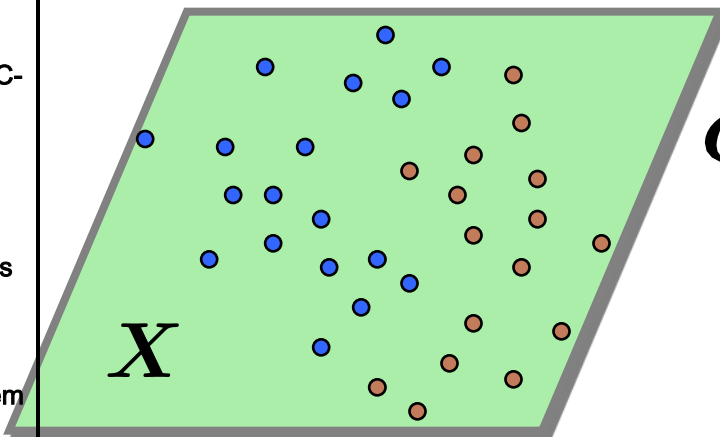
- 6 Kernels and Dot Products
- 6.1 Mercer's Theorem
- 6.2 Reproducing Kernel Hilbert Space

class -1 ● feature objects
class +1 ○ support vectors



separating hyperplane given by feature objects

feature selection



$\phi(x)$

Kernels, Dot Products, and Mercer's Theorem

- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Example for a kernel: $\Phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2} x_1 x_2)$

$\mathbf{x}^1 = (1, 1), \mathbf{x}^2 = (1, -1), \mathbf{x}^3 = (-1, 1), \mathbf{x}^4 = (-1, -1)$

with labels $y^1 = -1, y^2 = 1, y^3 = 1, y^4 = -1$

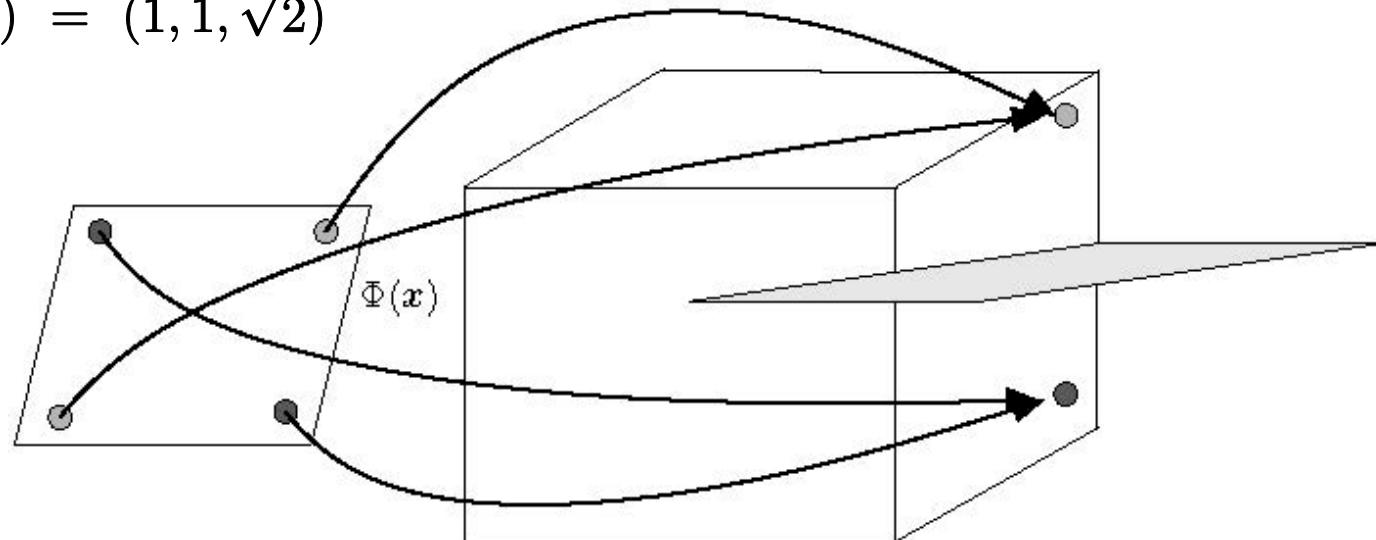
not separable in 2D but in feature space

$$\Phi(\mathbf{x}^1) = (1, 1, \sqrt{2})$$

$$\Phi(\mathbf{x}^2) = (1, 1, -\sqrt{2})$$

$$\Phi(\mathbf{x}^3) = (1, 1, -\sqrt{2})$$

$$\Phi(\mathbf{x}^4) = (1, 1, \sqrt{2})$$



Kernels, Dot Products, and Mercer's Theorem



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

$$\mathbf{x}_{mn} \equiv (\mathbf{x}^m)_n$$

$$\begin{aligned} \Phi^T(\mathbf{x}^i)\Phi(\mathbf{x}^j) &= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{i2} x_{j1} x_{j2} = \\ &= (x_{i1} x_{j1} + x_{i2} x_{j2})^2 = \left((\mathbf{x}^i)^T \mathbf{x}^j \right)^2 \end{aligned}$$

Another example:

$$\begin{aligned} \Phi(\mathbf{x}) &= \left(x_1^3, x_2^3, \sqrt{3} x_1^2 x_2, \sqrt{3} x_2^2 x_1, \right. \\ &\quad \left. \sqrt{3} x_1^2, \sqrt{3} x_2^2, \sqrt{6} x_1 x_2, \sqrt{3} x_1, \sqrt{3} x_2 \right) \end{aligned}$$

$$\begin{aligned} \Phi^T(\mathbf{x}^i)\Phi(\mathbf{x}^j) &= \\ &= x_{i1}^3 x_{j1}^3 + x_{i2}^3 x_{j2}^3 + \\ &+ 3 x_{i1}^2 x_{i2} x_{j1}^2 x_{j2} + 3 x_{i2}^2 x_{i1} x_{j2}^2 x_{j1} + \\ &+ 3 x_{i1}^2 x_{j1}^2 + 3 x_{i2}^2 x_{j2}^2 + \\ &+ 6 x_{i1} x_{i2} x_{j1} x_{j2} + 3 x_{i1} x_{j1} + 3 x_{i2} x_{j2} = \\ &= \left((\mathbf{x}^i)^T \mathbf{x}^j + 1 \right)^3 - 1 \end{aligned}$$

Kernels, Dot Products, and Mercer's Theorem



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

kernel k : a function which produces a scalar out of two vectors

$$k(\mathbf{x}^i, \mathbf{x}^j) = \left((\mathbf{x}^i)^T \mathbf{x}^j + 1 \right)^3$$

Certain kernels represent the mapping of vectors into a feature space and a dot product in this space.

The following theorem characterizes functions which build a dot product in some space.

Kernels, Dot Products, and Mercer's Theorem



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Theorem 1 (Mercer)

Let the kernel k be symmetric and from $L_2(X \times X)$ defining a Hilbert-Schmidt operator

$$T_k(f)(\mathbf{x}) = \int_X k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}') d\mathbf{x}' .$$

If T_k is positive definite, i.e. for all $f \in L_2(X)$

$$\int_{X \times X} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \geq 0 ,$$

then T_k has eigenvalues $\lambda_j \geq 0$ with associated eigenfunctions $\psi_j \in L_2(X)$. Further

$$(\lambda_1, \lambda_2, \dots) \in \ell_1$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_j \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{x}') ,$$

where ℓ_1 is the space of vectors with finite one-norm and the last sum converges absolutely and uniformly for almost all \mathbf{x} and \mathbf{x}' .

Kernels, Dot Products, and Mercer's Theorem



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

The sum may be an infinite sum for which the eigenvalues converge to zero. In this case the feature space is an infinite dimensional space.

„for almost all“ means „except for a set with zero measure“:
single points may lead to an absolute and uniform convergence

That the convergence is „absolutely and uniformly“ is important because the sum can be resorted and derivative and sum can be exchanged.

If X is a compact interval $[a, b]$ and k is continuous then positive definiteness of T_k is equivalent to positive definiteness of k . A kernel k is *positive semi-definite* if for all l , all $\mathbf{x}^1, \dots, \mathbf{x}^l$, and all $\alpha_i, 1 \leq i \leq l$

$$\sum_{i,j=1,1}^{l,l} \alpha_i \alpha_j k(\mathbf{x}^i, \mathbf{x}^j) \geq 0$$

Using the gram matrix \mathbf{K} with $K_{ij} = k(\mathbf{x}^i, \mathbf{x}^j)$ and the vector $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_l)$ this is

$$\boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \geq 0$$

Reproducing Kernel Hilbert Space



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Reproducing kernel Hilbert spaces (RKHSs) are an important theoretical tool for proofing properties of kernel methods.

X is a set and H a Hilbert space of complex-valued functions on X .

H is a **reproducing kernel Hilbert space** if every linear evaluation map

$$L_x : f \mapsto f(x)$$

from H to the complex numbers is continuous for any x in X .

Reproducing Kernel Hilbert Space



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

The next theorem is essential to show that every Hilbert space uniquely defines a kernel (from Hilbert space to kernel):

Theorem 1 (Riesz representation theorem)

Let H^ denote H 's dual space, consisting of all continuous linear functionals from H into complex numbers \mathbb{C} . If x is an element of H , then the function ϕ_x defined by*

$$\phi_x(y) = \langle y, x \rangle \quad \forall y \in H,$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product of the Hilbert space, is an element of H^ . The Riesz representation theorem states that every element of H^* can be written uniquely in this form, that is the mapping*

$$\Phi : H \rightarrow H^*, \quad \Phi(x) = \phi_x$$

is an isometric (anti-) isomorphism.

Reproducing Kernel Hilbert Space



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Riesz theorem for kernel k : for every x in X there exists an element k_x of H , the „point-evaluation functional“ at the point x , with the property:

$$f(x) = \langle f, k_x \rangle \quad \forall f \in H$$

We define the kernel as a function $k : X \times X \rightarrow \mathbb{C}$

$$k(x, y) \stackrel{\text{def}}{=} k_x(y)$$

k is called the **reproducing kernel** for the Hilbert space H .

k is completely determined by H because of Riesz theorem.

H is called „**reproducing kernel Hilbert space**“ (RKHS).

Reproducing Kernel Hilbert Space



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex

5 Statistical Learning Theory

5.1 Error Bound Example

5.2 Empirical Risk Minimization

5.2.1 Complexity: Finite Number of Functions

5.2.2 Complexity: VC-Dimension

5.3 Error Bounds

5.4 Structural Risk Minimization

5.5 Margin

5.6 Average Bounds

6 Kernels and Dot Products

6.1 Mercer's Theorem

6.2 Reproducing Kernel Hilbert Space

The next theorem states that every symmetric, positive definite kernel uniquely defines an RKHS (from kernel to Hilbert space):

Theorem 1 (Moore-Aronszajn)

Suppose k is a symmetric, positive definite kernel on a set X . Then there is a unique Hilbert space of functions on X for which k is a reproducing kernel.

Proof.

Define, for all x in X , $k_x = k(x, \cdot)$. Let H_0 be the linear span of $\{k_x : x \in X\}$. Define an inner product on H_0 by

$$\left\langle \sum_{j=1}^n b_j k_{y_j}, \sum_{i=1}^m a_i k_{x_i} \right\rangle = \sum_{i=1}^m \sum_{j=1}^n \overline{a_i} b_j k(y_j, x_i).$$

The symmetry of this inner product follows from the symmetry of k and the non-degeneracy follows from the fact that k is positive definite.

Let H be the completion of H_0 with respect to this inner product. Then H consists of functions of the form

$$f(x) = \sum_{i=1}^{\infty} a_i k_{x_i}(x),$$

where $\sum_{i=1}^{\infty} a_i^2 k(x_i, x_i) < \infty$. The fact that the above sum converges for every x follows from the Cauchy-Schwartz inequality.

We confirm the RKHS property:

$$\langle f, k_x \rangle = \left\langle \sum_{i=1}^{\infty} a_i k_{x_i}, k_x \right\rangle = \sum_{i=1}^{\infty} a_i k(x_i, x) = f(x).$$

Therefore we have

$$f(x) = \sum_{i=1}^{\infty} a_i k(x_i, x).$$

which is for example the model class of support vector machines.

To prove uniqueness, let G be another Hilbert space of functions for which k is a reproducing kernel. For any x and y in X , The RKHS property implies that

$$\langle k_x, k_y \rangle_H = k(x, y) = \langle k_x, k_y \rangle_G.$$

By linearity,

$$\langle \cdot, \cdot \rangle_H = \langle \cdot, \cdot \rangle_G$$

on the span of $\{k_x : x \in X\}$. Then $G = H$ by the uniqueness of the completion. **End Proof.**

Reproducing Kernel Hilbert Space



- 4 Noise Models
 - 4.1 Gaussian Noise
 - 4.2 Laplace Noise
 - 4.3 Binary Models
 - 4.3.1 Cross-Entropy
 - 4.3.2 Logistic Regression
 - 4.3.3 Log. Regression Convex
 - 4.3.4 Softmax
 - 4.3.5 Softmax Convex
- 5 Statistical Learning Theory
 - 5.1 Error Bound Example
 - 5.2 Empirical Risk Minimization
 - 5.2.1 Complexity: Finite Number of Functions
 - 5.2.2 Complexity: VC-Dimension
 - 5.3 Error Bounds
 - 5.4 Structural Risk Minimization
 - 5.5 Margin
 - 5.6 Average Bounds
- 6 Kernels and Dot Products
 - 6.1 Mercer's Theorem
 - 6.2 Reproducing Kernel Hilbert Space

Properties of the reproducing kernel and the RKHS:

- **reproducing property:** $k(x, y) = \langle k(x, \cdot), k(y, \cdot) \rangle$
- **orthonormal sequences, kernel expansion:** If $\{\phi_k\}_{k=1}^{\infty}$ is an orthonormal sequence and the closure of its span is H , then

$$k(x, y) = \sum_{k=1}^{\infty} \phi_k(x) \phi_k(y)$$

- **Moore-Aronszajn Theorem:** every symmetric, positive definite kernel defines a unique reproducing kernel Hilbert space.

For machine learning, the model class

$$f(x) = \sum_{i=1}^{\infty} a_i k_{x_i}(x) = \sum_{i=1}^{\infty} a_i k(x_i, x)$$

is of importance and can be represented by the RKHS.