

# Theoretical Concepts of Machine Learning Part 2

Sepp Hochreiter

Institute of Bioinformatics  
Johannes Kepler University, Linz, Austria

---

1 Introduction

2 Generalization Error

3 Maximum Likelihood

4 Noise Models

5 Statistical Learning Theory

7 Optimization Techniques

8 Bayes Techniques

9 Linear Models

## **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

## 9 Linear Models

- 9.1 Linear Regression
- 9.2 Analysis of Variance
- 9.3 Analysis of Covariance
- 9.4 Mixed Effects Models
- 9.5 Generalized Linear Models
- 9.6 Regularization

# Chapter 7

## Optimization Techniques

---

# Error Minimization and Model Selection



## 7 Optimization Techniques

### 7.1 Parameter Optimization

#### 7.1.1 Search and Evolutionary Methods

#### 7.1.2 Gradient Descent

#### 7.1.3 Step-size

#### 7.1.4 Update Direction

#### 7.1.5 Levenberg-Marquardt Algorithm

#### 7.1.6 Predictor

#### Corrector Methods

#### 7.1.7 Convergence

### 7.2 On-line

### 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:

### Evidence Framework

### 8.6 Hyper-parameter:

### Integrate Out

### 8.7 Model Comparison

### 8.8 Posterior Sampling

Model class: **parameterized** models with parameter vector  $w$

Goal: to find the optimal model in parameter space

optimal model: optimizes the objective function

objective function: is defined by the problem to solve

- empirical error
- complexity

to find the optimal model: search in the parameter space

# Search Methods and Evolutionary Approaches



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

In principle any search algorithm can be used

**random search:** randomly a parameter vector is selected and then evaluated -- best solution kept

**exhaustive search:** parameter space is searched systematically

will find the optimal solution for a finite set of possible parameters  
do not use any dependency between objective function and parameter  
→ singular solution

In general there are dependencies between objective function and parameters which should be utilized

# Search Methods and Evolutionary Approaches



7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

first dependency: good solutions are near good solutions  
(objective is not continuous)

**stochastic gradient:** locally optimize the objective function

**genetic algorithms:** a stochastic gradient corresponds to mutations which are small

another dependency: good solutions share properties (independent!)

*genetic algorithms:* - “crossover mutations” combines parts of different parameter vectors  
- solutions have independent *building blocks*

# Search Methods and Evolutionary Approaches



7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

other evolutionary strategies to optimize models:

- genetic programming
- swarm algorithms
- ant algorithms
- self-modifying policies
- Optimal Ordered Problem Solver  
(class is predefined by the search algorithm,  
model class is not parameterized)  
the latter modify their own search strategy

parallel search: different locations in the parameter space



# Search Methods and Evolutionary Approaches



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

**simulated annealing** (SA): overcome the problem of local optima

- find the global solution with probability one if the annealing process is slow enough
- jumps from the current state into another state where the probability is given by the objective function
- objective function = energy function
- energy function follows the Maxwell-Boltzmann distribution
- sampling is similar to the Metropolis-Hastings algorithm
- **temperature**: which jumps are possible
  - beginning: high temperature → large jumps into energetically worse regions
  - end: low temperature → small jumps into energetically favorable regions

# Search Methods and Evolutionary Approaches



7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

## Advantages:

- discrete problems and non-differentiable objective functions
- very easy to implement

## Disadvantages:

- computationally expensive for large parameter spaces
- depend on the representation of the model

parameter vector of length  $W$ :  $2^W$  decisions to make to decide whether to in- or decrease components

# Search Methods and Evolutionary Approaches



7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

in the following: the objective function is differentiable with respect to the parameters

→ use of gradient information

Already treated: convex optimization (only one minimum)

# Gradient Descent



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

**objective function:**  $R(g(\cdot; \mathbf{w}))$  is a differentiable function

**parameter vector**  $\mathbf{w}$  is  $W$ -dimensional

**model:**  $g$

for simplicity:  $R(\mathbf{w}) = R(g(\cdot; \mathbf{w}))$

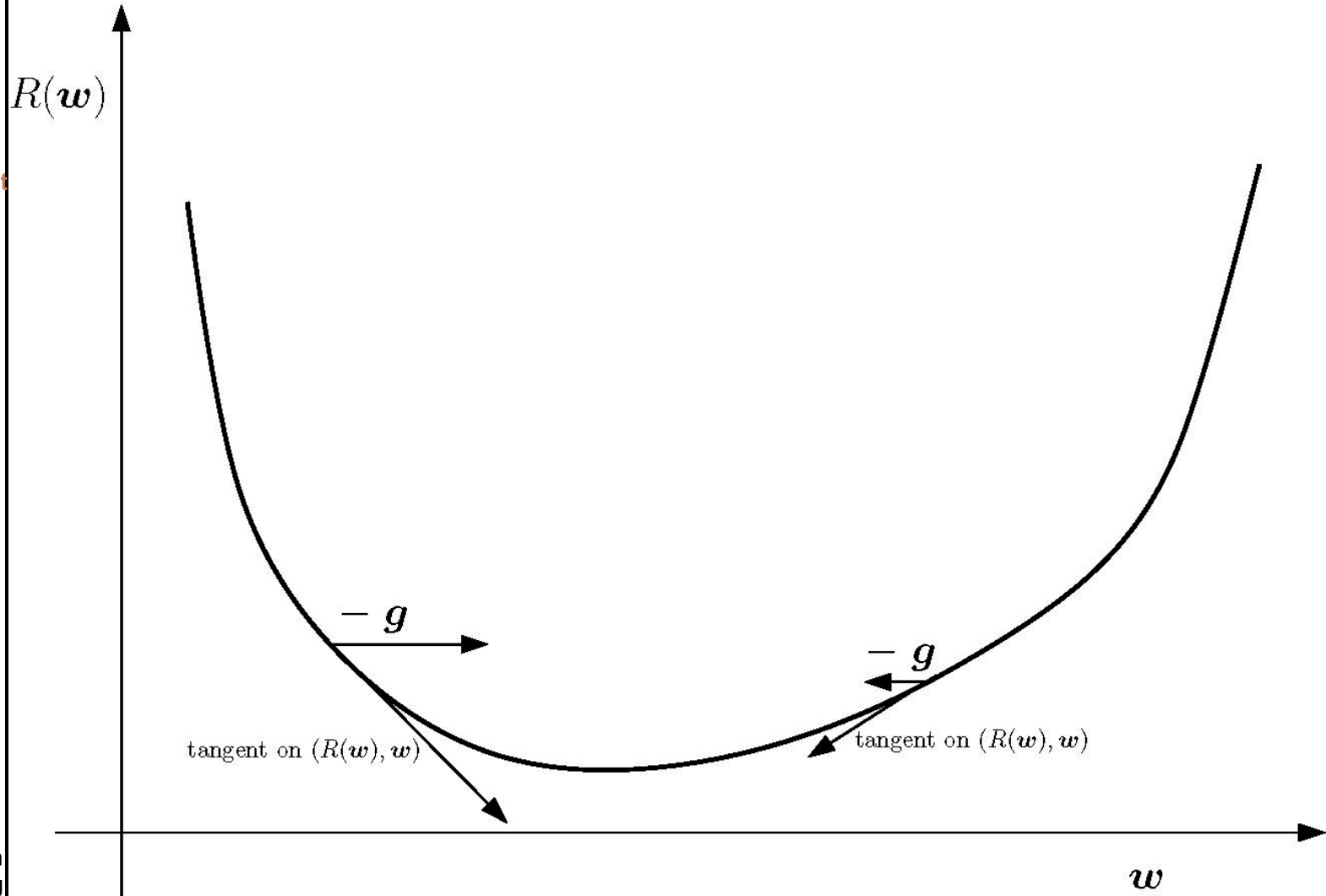
**gradient:**  $\frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} = \nabla_{\mathbf{w}} R(\mathbf{w}) = \left( \frac{\partial R(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial R(\mathbf{w})}{\partial w_W} \right)^T$

short:  $\mathbf{g} = \nabla_{\mathbf{w}} R(\mathbf{w})$

negative gradient: direction of the steepest descent of the objective

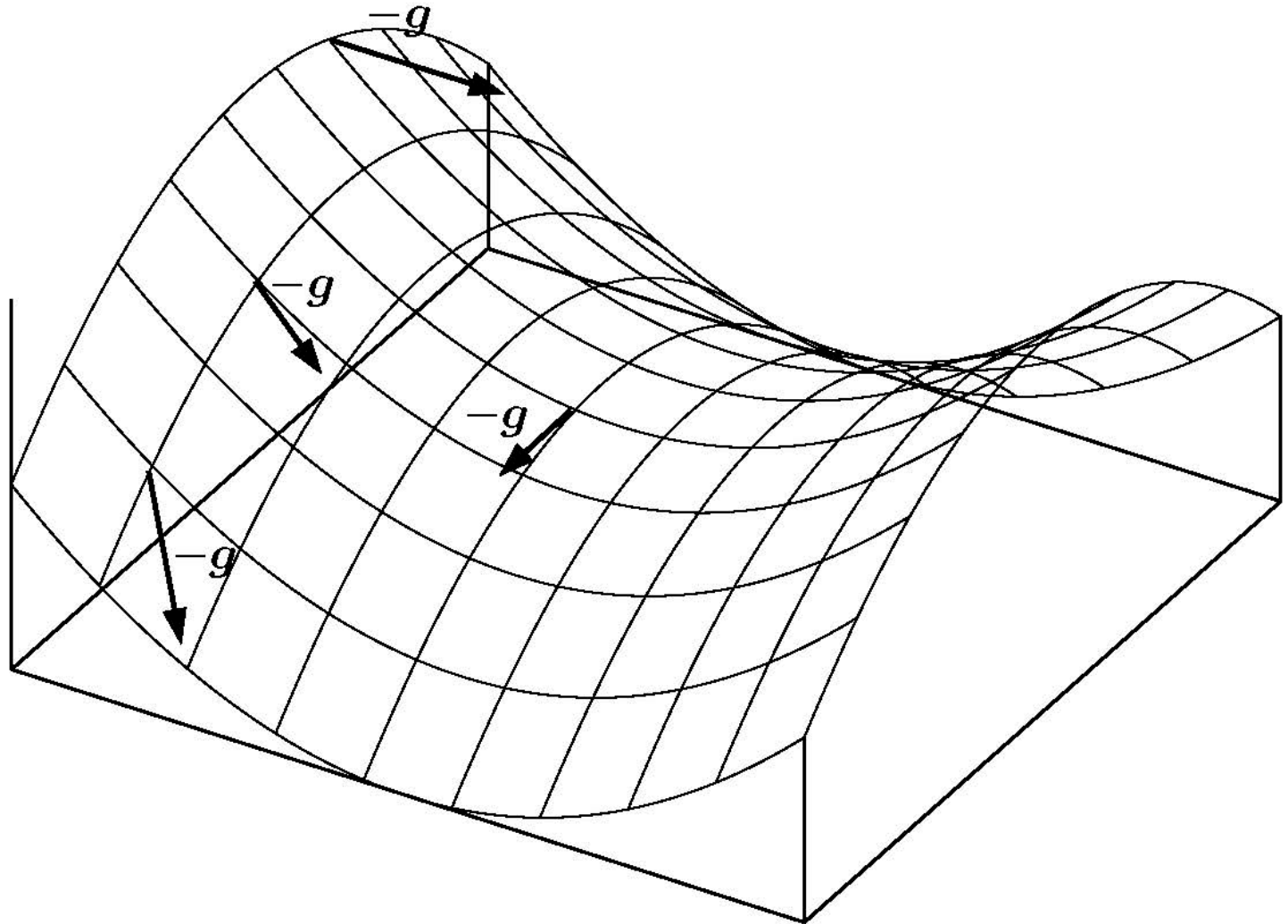
# Gradient Descent

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling



# Gradient Descent

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling



# Gradient Descent



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

gradient is locally  $\rightarrow$  small step in the negative gradient direction

**learning rate**  $0 < \eta$  : controls step-size

gradient descent update:

$$\Delta \mathbf{w} = -\eta \nabla_{\mathbf{w}} R(\mathbf{w})$$
$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \Delta \mathbf{w}$$

**momentum term:** - gradient descent oscillates  
- flat plateau  
 $\rightarrow$  learning slows down

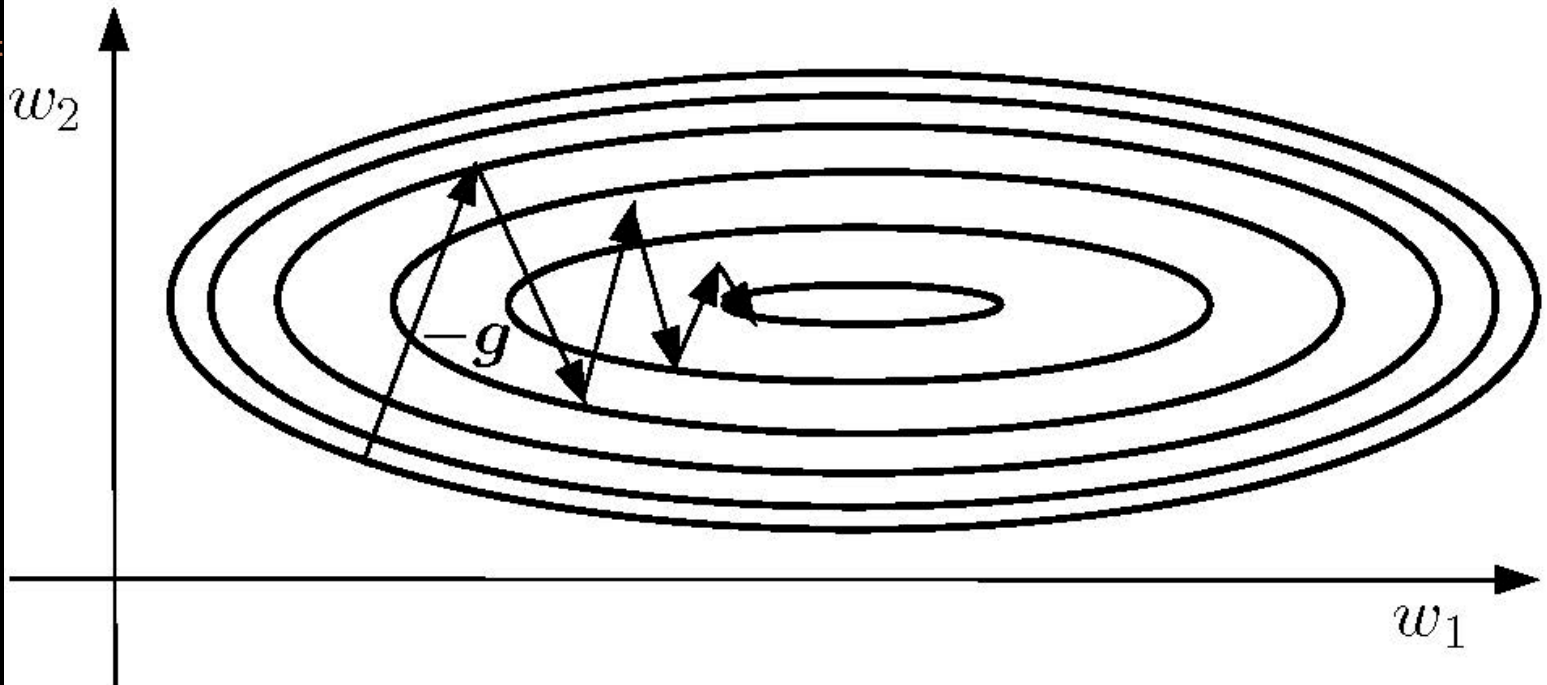
$$\Delta^{\text{new}} \mathbf{w} = -\eta \nabla_{\mathbf{w}} R(\mathbf{w})$$
$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \Delta^{\text{new}} \mathbf{w} + \mu \Delta^{\text{old}}$$
$$\Delta^{\text{old}} = \Delta^{\text{new}}$$

**momentum factor** or **momentum parameter:**  $0 \leq \mu \leq 1$

# Gradient Descent

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 7.3 Convex Optimization

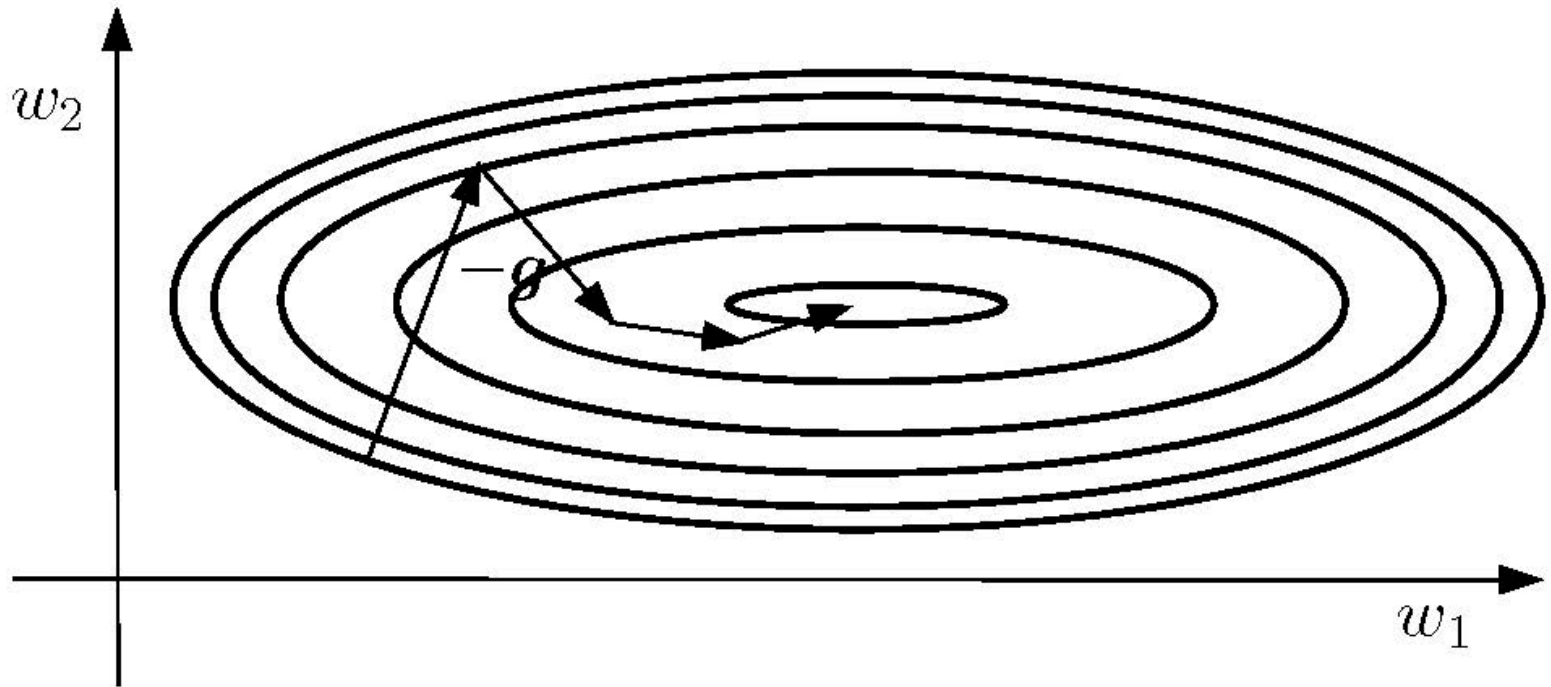
- 8 Bayes Techniques
  - 8.1 Likelihood, Prior, Posterior, Evidence
  - 8.2 Maximum A Posteriori Approach
  - 8.3 Posterior Approximation
  - 8.4 Error Bars and Confidence Intervals
  - 8.5 Hyper-parameter: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling





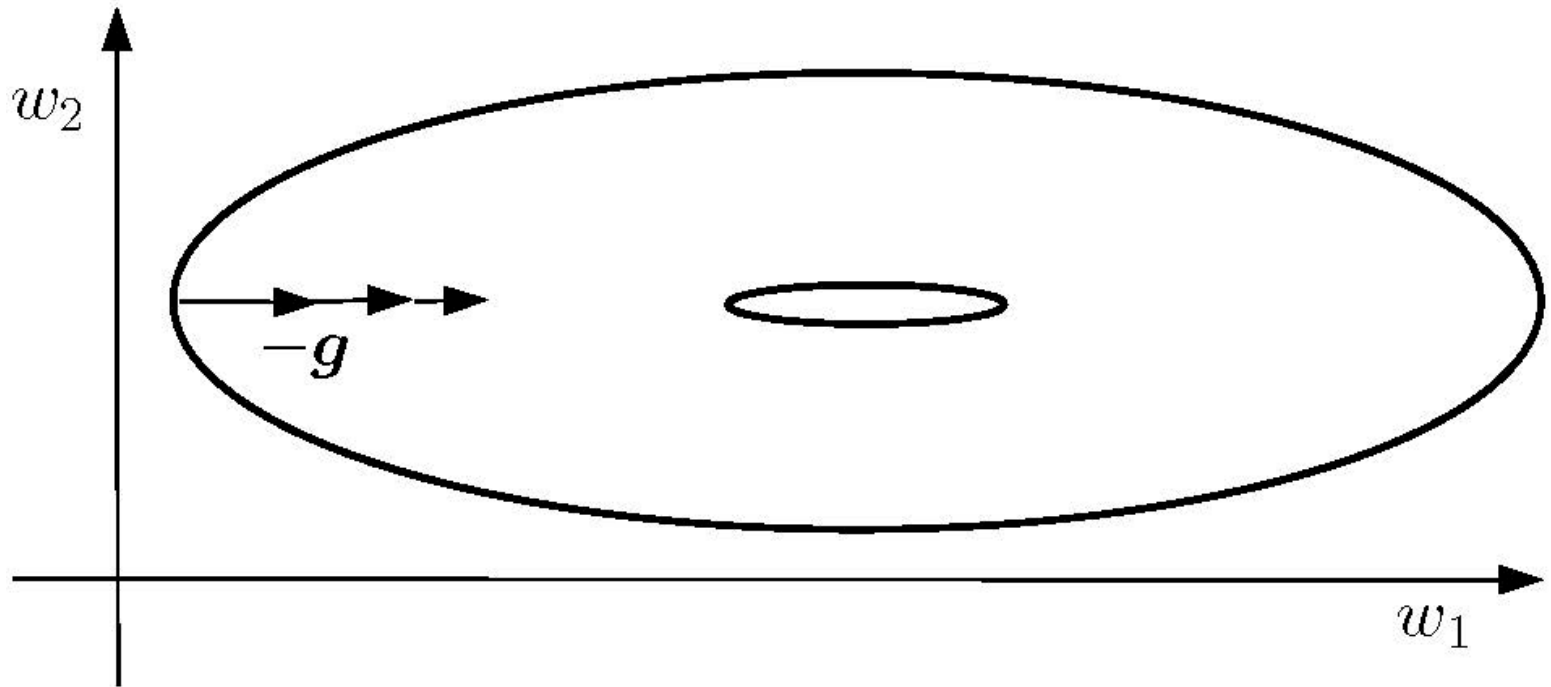
# Gradient Descent

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling



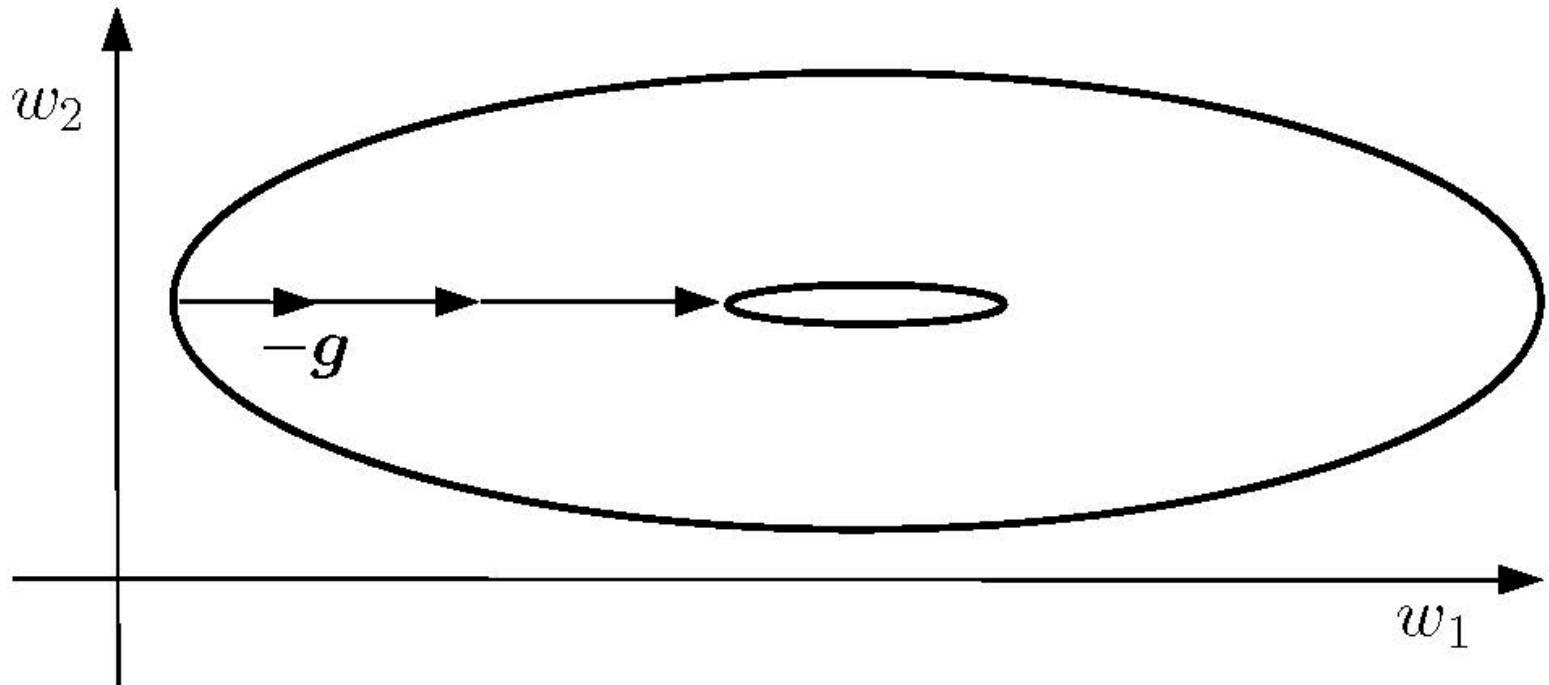
# Gradient Descent

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling



# Gradient Descent

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling



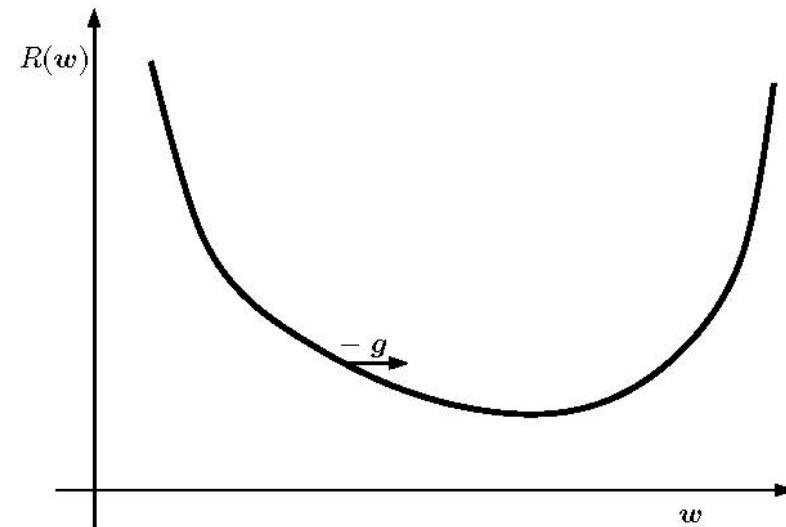
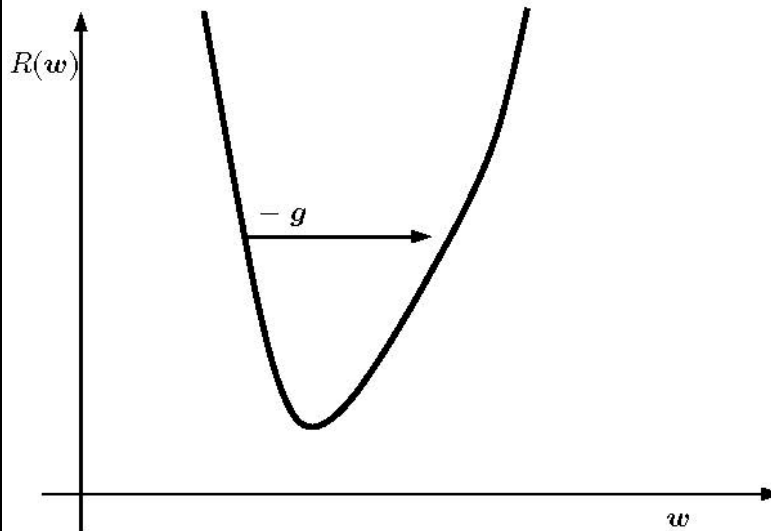
# Step-size Optimization

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

step-size should be adjusted to the curvature of the error surface

flat curve: large step-size

steep minima: small step-size



# Step-size Optimization: Heuristics



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

**Learning rate adjustments:** resilient backpropagation (Rprop)

change of the risk  $\Delta R = R(\mathbf{w} + \Delta \mathbf{w}) - R(\mathbf{w})$

learning rate adjustment: 
$$\eta^{\text{new}} = \begin{cases} \rho \eta^{\text{old}} & \text{if } \Delta R \leq 0 \\ \sigma \eta^{\text{old}} & \text{if } \Delta R > 0 \end{cases}$$

$\rho > 1$     $\sigma < 1$    typical:  $\rho = 1.1$ ,  $\sigma = 0.5$

# Step-size Optimization: Heuristics



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Largest eigenvalue of the Hessian

$$\text{Hessian } \mathbf{H}: H_{ij} = \frac{\partial^2 R(\mathbf{w})}{\partial w_i \partial w_j}$$

Largest eigenvalue  $e_{\max}$  bounds the learning rate:

$$\eta \leq \frac{2}{\lambda_{\max}}$$

maximal eigenvalue: matrix iteration  $\mathbf{a} = \sum_{i=1}^W \alpha_i \mathbf{e}_i$

$$\mathbf{H}^s \mathbf{a} = \sum_{i=1}^W \lambda_i^s \alpha_i \mathbf{e}_i \approx \lambda_{\max}^s \alpha_{\max} \mathbf{e}_{\max}$$

Important question:

how time intensive is it to obtain the Hessian or the product of the Hessian with a vector?

# Step-size Optimization: Heuristics



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Individual learning rate for each parameter

$$\Delta w_i = -\eta_i g_i \quad g_i = \left[ \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} \right]_i$$

parameters not independent from each other  
if  $\eta_i > 0$  then the error decreases

**delta-delta rule:**  $\Delta \eta_i = \gamma g_i^{\text{new}} g_i^{\text{old}}$

**delta-bar-delta rule:**  $\Delta \eta_i = \begin{cases} \kappa & \text{if } \bar{g}_i^{\text{old}} g_i^{\text{new}} > 0 \\ -\phi g_i^{\text{new}} & \text{if } \bar{g}_i^{\text{old}} g_i^{\text{new}} \leq 0 \end{cases}$

where  $\bar{g}_i^{\text{new}} = (1 - \theta) g_i^{\text{new}} + \theta \bar{g}_i^{\text{old}}$

$\bar{g}_i$  is an exponentially weighted average

disadvantage: many hyper-parameters

# Step-size Optimization: Heuristics



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Quickprop

developed in the context of neural networks (“backpropagation”)

$$\Delta^{\text{new}} w_i = \frac{g_i}{g_i^{\text{old}} - g_i^{\text{new}}} \Delta^{\text{old}} w_i$$

Taylor expansion

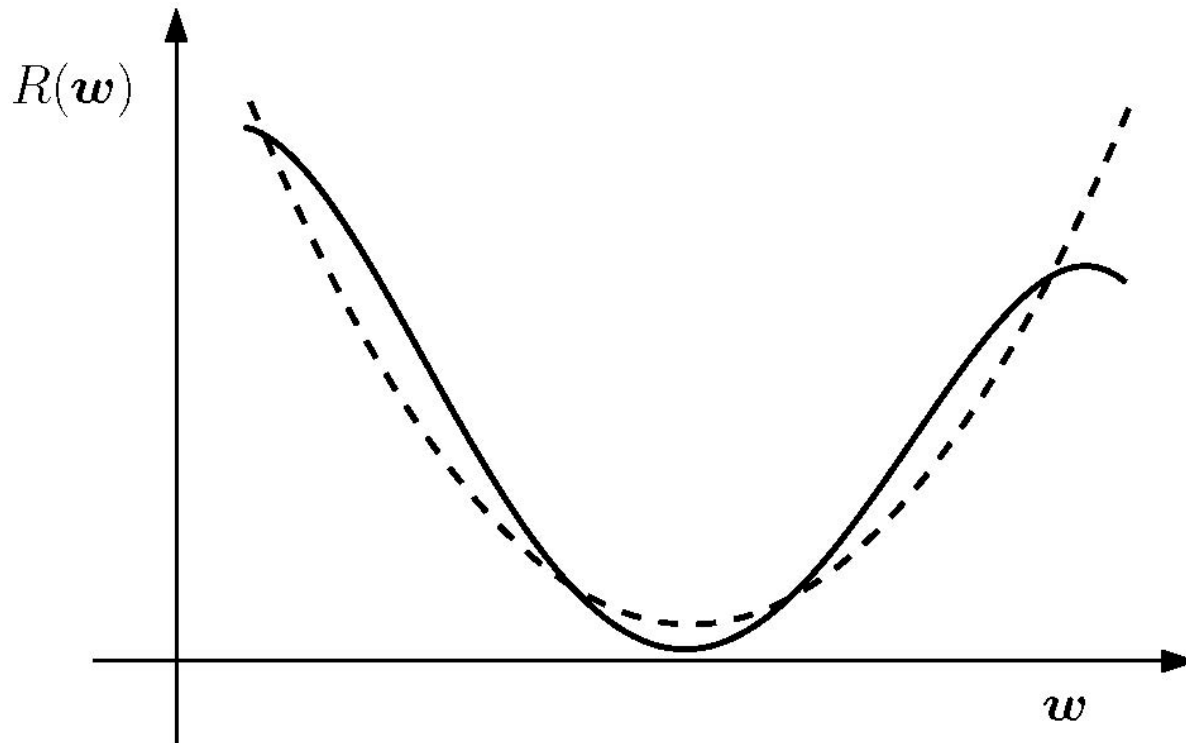
$$\begin{aligned} R(w_i + \Delta w_i) &= R(w_i) + \frac{\partial R(\mathbf{w})}{\partial w_i} \Delta w_i + \\ &\frac{1}{2} \frac{\partial^2 R(\mathbf{w})}{(\partial w_i)^2} (\Delta w_i)^2 + O((\Delta w_i)^3) = \\ R(w_i) + g_i \Delta w_i + \frac{1}{2} g'_i (\Delta w_i)^2 + O((\Delta w_i)^3) \end{aligned}$$



# Step-size Optimization: Heuristics

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

quadratic approximation



# Step-size Optimization: Heuristics

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

minimize  $R(w_i + \Delta w_i) - R(w_i)$  with respect to  $\Delta w_i$

zero derivative:  $\Delta w_i = -\frac{g_i}{g'_i}$

Now approximate  $g'_i = g'_i(w_i)$  by  $g'_i(w_i^{\text{old}})$

difference quotient:  $g'_i = \frac{g_i^{\text{new}} - g_i^{\text{old}}}{\Delta^{\text{old}} w_i}$

$$g_i^{\text{old}} = g_i(w_i - \Delta^{\text{old}} w_i)$$

- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

update direction  $\mathbf{d}$  is given  $\Delta \mathbf{w} = \eta \mathbf{d}$

minimize  $R(\mathbf{w} + \eta \mathbf{d})$  with respect to  $\eta$

quadratic functions with minimum  $\mathbf{w}^*$  : 
$$\eta = \frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T \mathbf{H}(\mathbf{w}^*) \mathbf{g}}$$

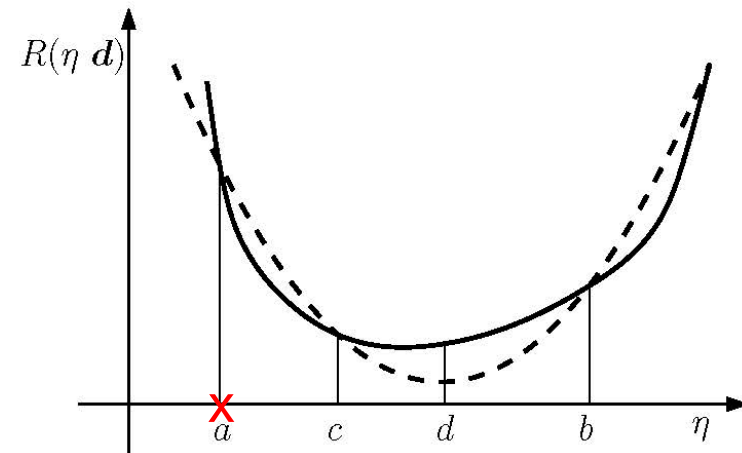
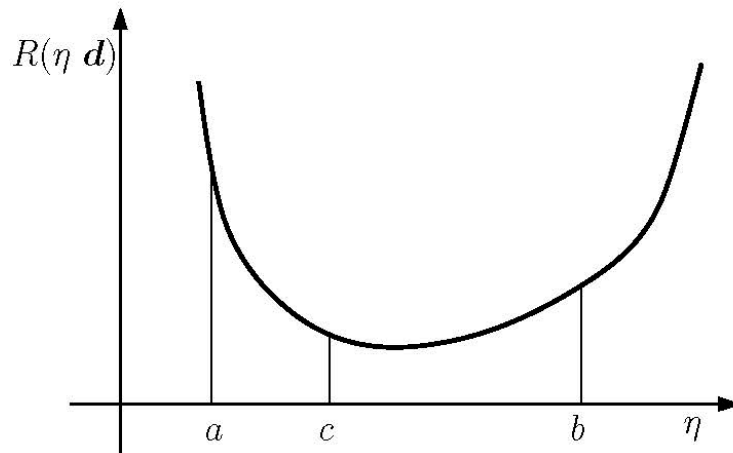
However: 1) valid only near the minimum and  
2)  $\mathbf{H}(\mathbf{w}^*)$  unknown

# Line Search

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Line search:

1. Fits a parabola through three points and determines its minimum
2. Adds minimum point to the points
3. Point with largest  $R$  is discharged



## Algorithm: Line-Search

```
BEGIN initialization  $a_0, b_0, c_0; R(a_0) > R(c_0);$   
     $R(b_0) > R(c_0), \text{Stop}=\text{false}, i = 0$   
END initialization  
BEGIN line search  
    while STOP=false do  
        fit quadratic function through  $a_i, b_i, c_i$   
        determine minimum  $d_i$  of quadratic function  
        if stop criterion fulfilled, e.g.  $|a_i - b_i| < \epsilon$  or  $|R(b_0) - R(c_0)| < \epsilon$  then  
            Stop=true  
        else  
             $c_{i+1} = d_i$   
             $b_{i+1} = c_i$   
             $a_{i+1} = \begin{cases} a_i & \text{if } R(a_i) \leq R(b_i) \\ b_i & \text{if } R(a_i) > R(b_i) \end{cases}$   
        end if  
         $i = i + 1$   
    end while  
END line search
```

7 Optimization Techniques

7.1 Parameter Optimization

7.1.1 Search and Evolutionary Methods

7.1.2 Gradient Descent

7.1.3 Step-size

7.1.4 Update Direction

7.1.5 Levenberg-Marquardt Algorithm

7.1.6 Predictor

Corrector Methods

7.1.7 Convergence

7.2 On-line

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: Evidence Framework

8.6 Hyper-parameter: Integrate Out

8.7 Model Comparison

8.8 Posterior Sampling

# Optimization of the Update Direction



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

The default direction is the negative gradient:  $-g$

However there are better approaches

# Newton and Quasi-Newton Method



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

The gradient vanishes at the global minimum  $\mathbf{w}^*$

$$\nabla_{\mathbf{w}} R(\mathbf{w}^*) = \mathbf{g}(\mathbf{w}^*) = \mathbf{0}$$

Taylor series:

$$R(\mathbf{w}) = R(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w}^*) (\mathbf{w} - \mathbf{w}^*) + O\left(\|\mathbf{w} - \mathbf{w}^*\|^3\right)$$

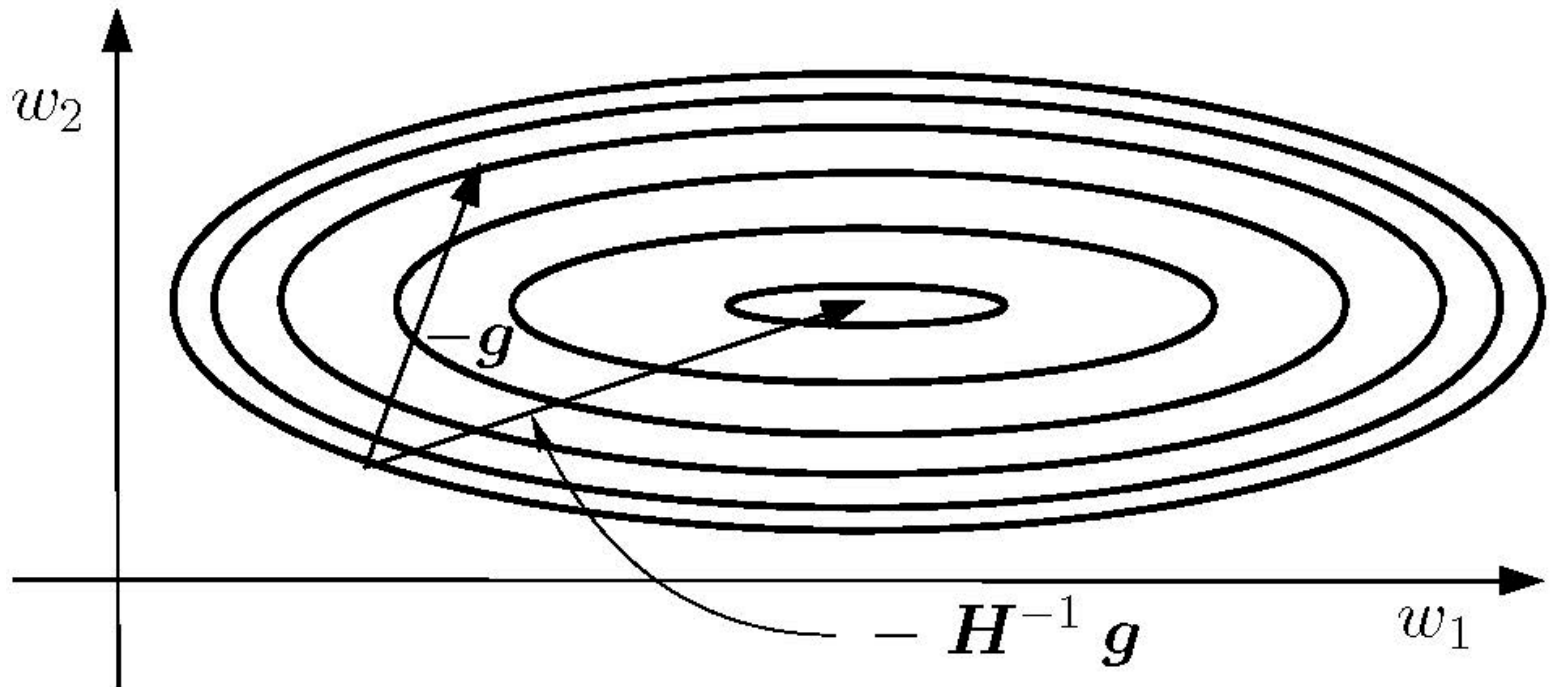
quadratic approximation:  $\mathbf{g} = \frac{\partial R}{\partial \mathbf{w}} = \mathbf{H}(\mathbf{w}^*) (\mathbf{w} - \mathbf{w}^*)$

Thus:  $\mathbf{w}^* = \mathbf{w} - \mathbf{H}^{-1} \mathbf{g}$

**Newton direction:**  $\mathbf{H}^{-1} \mathbf{g}$

# Newton and Quasi-Newton Method

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling





# Newton and Quasi-Newton Method



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

## Disadvantages

1. computationally expensive: inverse Hessian
2. only valid near the minimum (Hessian is positive definite)

*model trust region* approach: model is trusted up to a certain value

$H + \lambda I$  compromise between gradient and Newton direction

Hessian can be approximated by a diagonal matrix

# Newton and Quasi-Newton Method



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Quasi-Newton Method

$$\text{Newton equation: } \mathbf{w}^{\text{new}} - \mathbf{w}^{\text{old}} = -\mathbf{H}^{-1} (\mathbf{g}^{\text{new}} - \mathbf{g}^{\text{old}})$$

quasi-Newton condition

## Broyden-Fletcher-Goldfarb-Shanno (BFGS) method

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta \mathbf{G}^{\text{old}} \mathbf{g}^{\text{old}} \quad \text{where } \eta \text{ is found by line search and}$$

$$\mathbf{G}^{\text{new}} = \mathbf{G}^{\text{old}} + \frac{\mathbf{p} \mathbf{p}^T}{\mathbf{p}^T \mathbf{v}} - \frac{(\mathbf{G}^{\text{old}} \mathbf{v}) \mathbf{v}^T \mathbf{G}^{\text{old}}}{\mathbf{v}^T \mathbf{G}^{\text{old}} \mathbf{v}} + (\mathbf{v}^T \mathbf{G}^{\text{old}} \mathbf{v}) \mathbf{u} \mathbf{u}^T$$

$\mathbf{G}$  approximates the inverse Hessian, therefore the update is more complicated

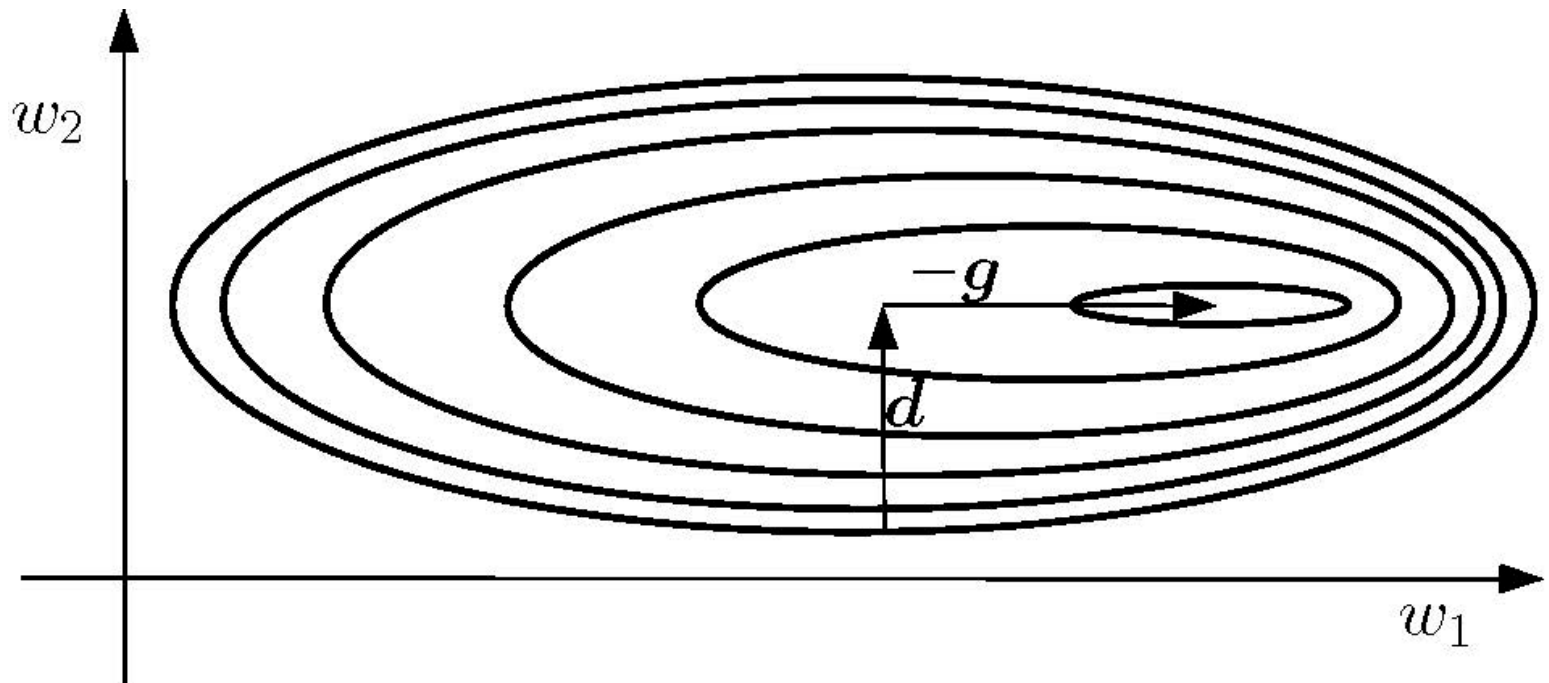
$$\begin{aligned} \mathbf{p} &= \mathbf{w}^{\text{new}} - \mathbf{w}^{\text{old}} \\ \mathbf{v} &= \mathbf{g}^{\text{new}} - \mathbf{g}^{\text{old}} \\ \mathbf{u} &= \frac{\mathbf{p}}{\mathbf{p}^T \mathbf{v}} - \frac{\mathbf{G}^{\text{old}} \mathbf{v}}{\mathbf{v}^T \mathbf{G}^{\text{old}} \mathbf{v}} \end{aligned}$$

# Conjugate Gradient

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

$R(\mathbf{w} + \eta \mathbf{d})$  is minimized with respect to  $\eta$

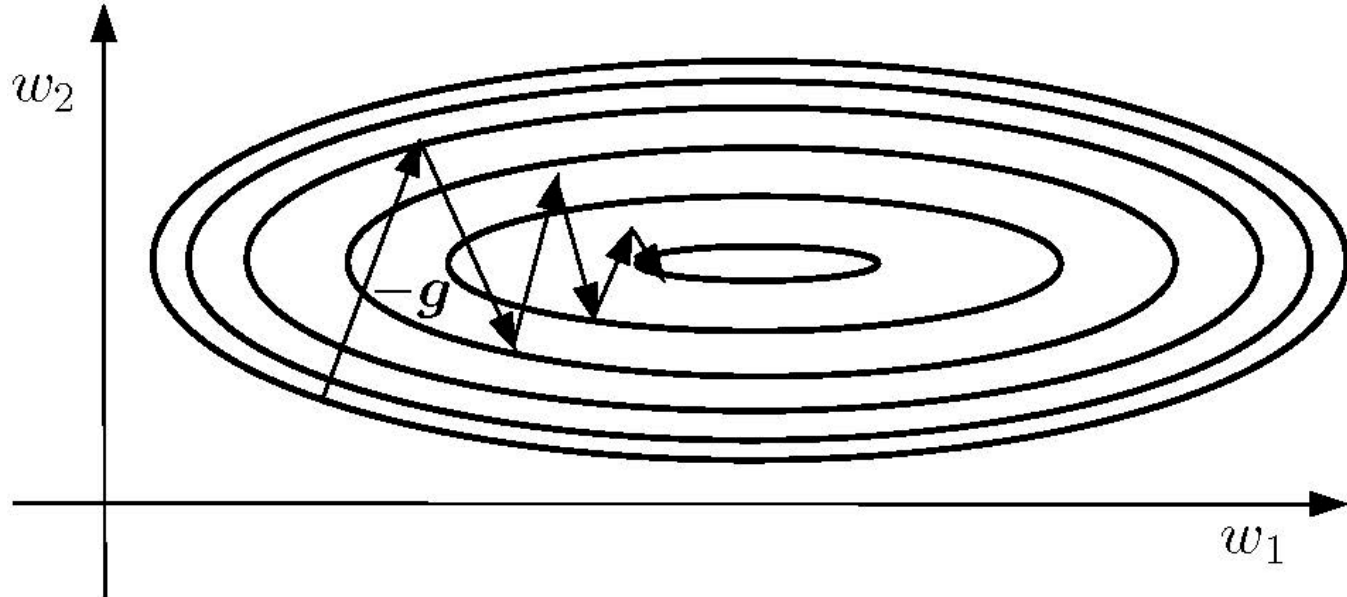
$$\frac{\partial}{\partial \eta} R(\mathbf{w} + \eta \mathbf{d}) = 0 \quad \square \quad (\mathbf{g}^{\text{new}})^T \mathbf{d}^{\text{old}} = 0$$



# Conjugate Gradient

- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling



still oscillations: different two-dimensional subspaces which alternate

avoid oscillations: new search directions are orthogonal to **all** previous

orthogonal is defined via the Hessian

in the parameter space the quadratic volume element is regarded

# Conjugate Gradient

- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

We require new gradient is orthogonal to old direction:

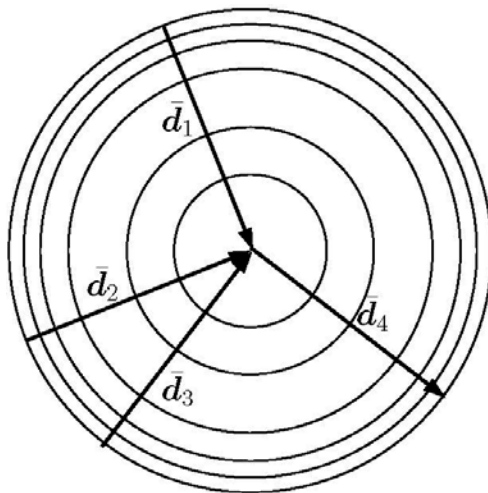
$$g(w^{\text{new}} + \eta d^{\text{new}})^T d^{\text{old}} = 0$$

Taylor expansion:

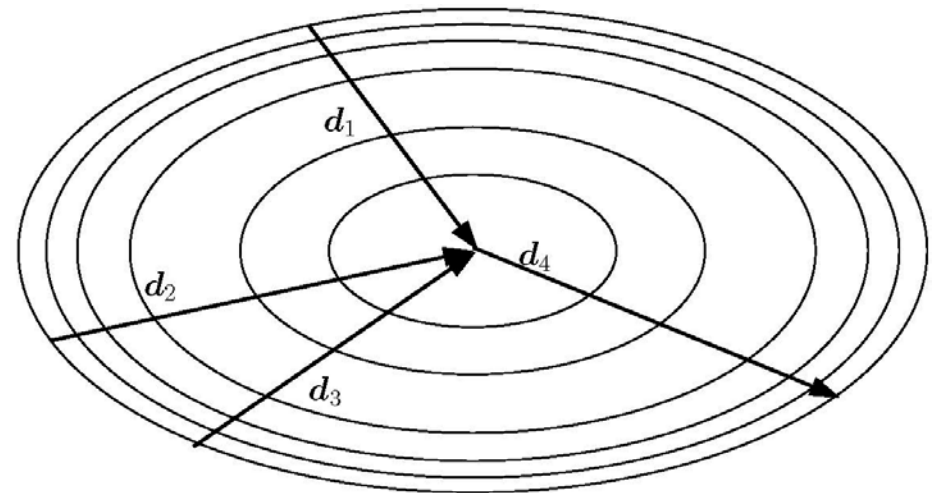
$$g(w^{\text{new}} + \eta d^{\text{new}})^T = g(w^{\text{new}})^T + \eta H(w^{\text{new}}) d^{\text{new}} + O(\eta^2)$$

and  $g(w^{\text{new}})^T d^{\text{old}} = 0$  follows:

$$(d^{\text{new}})^T H(w^{\text{new}}) d^{\text{old}} = 0 \text{ conjugate directions}$$



$$d = H^{-1/2} \bar{d}$$



$$d_1^T H d_2 = 0 \quad d_3^T H d_4 = 0$$

Conjugate directions: orthogonal directions in other coordinate system

# Conjugate Gradient



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

quadratic problem  $R(w) = \frac{1}{2} w^T H w + c^T w + k$

$0 = H w^* + c$       $g(w^*) = H w^* + c$       $g(w_j) = g_j = H w_j + c$

conjugate directions are linearly independent, therefore

$$w^* - w_1 = \sum_{i=1}^W \eta_i d_i \qquad w_j - w_1 = \sum_{i=1}^{j-1} \eta_i d_i$$

multiplied by  $d_j^T H$

multiplied by  $d_j^T H$

$H w^* = -c$

$d_j^T H d_i = 0$

$d_j^T H w_j = d_j^T H w_1$

$$-d_j^T (c + H w_1) = \sum_{i=1}^W \eta_i d_j^T H d_i = \eta_j d_j^T H d_j$$

$g_j = c + H w_j$

Gives:

$$\eta_j = - \frac{d_j^T g_j}{d_j^T H d_j}$$

# Conjugate Gradient



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

$$\mathbf{w}_{j+1} = \mathbf{w}_j + \eta_j \mathbf{d}_j$$

search directions  $\mathbf{d}_j$ ? We set

$$\mathbf{d}_{j+1} = -\mathbf{g}_{j+1} + \beta_j \mathbf{d}_j$$

Multiplying by  $\mathbf{d}_j^T \mathbf{H}$  gives  $\beta_j = \frac{\mathbf{g}_{j+1}^T \mathbf{H} \mathbf{d}_j}{\mathbf{d}_j^T \mathbf{H} \mathbf{d}_j}$

$$\mathbf{g}_{j+1} - \mathbf{g}_j = \mathbf{H} (\mathbf{w}_{j+1} - \mathbf{w}_j) = \eta_j \mathbf{H} \mathbf{d}_j$$

**Hestenes – Stiefel :**

$$\beta_j = \frac{\mathbf{g}_{j+1}^T (\mathbf{g}_{j+1} - \mathbf{g}_j)}{\mathbf{d}_j^T (\mathbf{g}_{j+1} - \mathbf{g}_j)}$$

# Conjugate Gradient

7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

$$d_k^T g_j = 0 \text{ for } k < j$$

gives

$$d_j^T g_j = -g_j^T g_j$$

**Polak – Ribiere :**

$$\beta_j = \frac{g_{j+1}^T (g_{j+1} - g_j)}{g_j^T g_j}$$

**Fletcher – Reeves :**

$$\beta_j = \frac{g_{j+1}^T g_{j+1}}{g_j^T g_j}$$



# Conjugate Gradient



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Updates are mathematically equivalent but numerical different  
The Polak-Ribiere has an edge over the other update rules

## Conjugate Gradient (Polak-Ribiere)

```
BEGIN initialization  $\mathbf{g}_0 = \nabla_{\mathbf{w}}R(\mathbf{w}_0)$ ,  $i = 0$ ,  $\mathbf{d}_0 = -\mathbf{g}_0$ , Stop=false
END initialization
BEGIN Conjugate Gradient
  while STOP=false do
    determine  $\eta_i$  by line search
     $\mathbf{w}_{j+1} = \mathbf{w}_j + \eta_j \mathbf{d}_j$ 
     $\mathbf{g}_{i+1} = \nabla_{\mathbf{w}}R(\mathbf{w}_i)$ 
     $\beta_j = \frac{\mathbf{g}_{j+1}^T(\mathbf{g}_{j+1} - \mathbf{g}_j)}{\mathbf{g}_j^T \mathbf{g}_j}$ 
     $\mathbf{d}_{j+1} = -\mathbf{g}_{j+1} + \beta_j \mathbf{d}_j$ 
    if stop criterion fulfilled, e.g.  $\|\mathbf{g}_{i+1}\| < \epsilon$  or  $|R(\mathbf{w}_{j+1}) - R(\mathbf{w}_j)| < \epsilon$  then
      STOP=true
    end if
     $i = i + 1$ 
  end while
END Conjugate Gradient
```

# Conjugate Gradient



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

$\eta$  are in most implementations found by line search

disadvantage of conjugate gradient compared to quasi-Newton:  
the line search must be done precisely (orthogonal gradients)

advantage of conjugate gradient compared to quasi-Newton:  
the storage is  $O(W)$  compared to  $O(W^2)$  for quasi-Newton

# Levenberg-Marquardt Algorithm



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

This algorithm is designed for **quadratic loss**  $R(\mathbf{w}) = \sum_{i=1}^l (e^i(\mathbf{w}))^2$

combine the errors  $e^i$  into a vector  $\mathbf{e}$

Jacobian of this vector is  $Z_{ij} = \frac{\partial e^i}{\partial w_j}$

linear approximation  $e(\mathbf{w}^{\text{new}}) = e(\mathbf{w}^{\text{old}}) + \mathbf{Z} (\mathbf{w}^{\text{new}} - \mathbf{w}^{\text{old}})$

Hessian of the loss function:

$$H_{jk} = \frac{\partial^2 R}{\partial w_j \partial w_k} = \sum_{i=1}^l \left( \frac{\partial e^i}{\partial w_j} \frac{\partial e^i}{\partial w_k} + e^i \frac{\partial^2 e^i}{\partial w_j \partial w_k} \right)$$

Small  $e^i$  or  $e^i \frac{\partial^2 e^i}{\partial w_j \partial w_k}$  averages out (negative and positive  $e^i$ )

approximate the Hessian by  $\mathbf{H} = \mathbf{Z}^T \mathbf{Z}$  (outer product)

Note that this is only valid for quadratic loss functions

# Levenberg-Marquardt Algorithm



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

unconstrained minimization problem

$$\frac{1}{2} \|e(\mathbf{w}^{\text{old}}) + \mathbf{Z}(\mathbf{w}^{\text{new}} - \mathbf{w}^{\text{old}})\|^2 + \lambda \|\mathbf{w}^{\text{new}} - \mathbf{w}^{\text{old}}\|^2$$

first term: minimizing the error

second term: minimizing the step size (for valid linear approximations)

**Levenberg-Marquardt** update rule

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^T e(\mathbf{w}^{\text{old}})$$

Small  $\lambda$  gives the Newton formula while large  $\lambda$  gives gradient descent

The Levenberg-Marquardt algorithm is a model trust region approach

# Predictor Corrector Methods for $R(\mathbf{w}) = 0$



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Problem: solve  $R(\mathbf{w}) = 0$

Idea: simple update (predictor step) by neglecting higher order terms  
then: correction by the higher order terms (corrector step)

Taylor series

$$R(\mathbf{w}^{\text{new}}) = R(\mathbf{w}^{\text{old}}) + S(\mathbf{w}^{\text{old}}, \Delta\mathbf{w}) + TS(\mathbf{w}^{\text{old}}, \Delta\mathbf{w})$$
$$S((\mathbf{w}^{\text{old}}, \mathbf{0}) = 0 \quad T((\mathbf{w}^{\text{old}}, \mathbf{0}) = 0$$

predictor step

$$R(\mathbf{w}^{\text{old}}) + S(\mathbf{w}^{\text{old}}, \Delta\mathbf{w}) = 0 \rightarrow \Delta\mathbf{w}$$

corrector step

$$R(\mathbf{w}^{\text{old}}) + S(\mathbf{w}^{\text{old}}, \Delta\mathbf{w}) + TS(\mathbf{w}^{\text{old}}, \Delta_{\text{pred}}\mathbf{w}) = 0$$

Iterative algorithm

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Gradient Descent

$$R(\mathbf{w}^{\text{old}}) - R(\mathbf{w}^{\text{new}}) = R(\mathbf{w}^{\text{old}}) \left( \frac{(\mathbf{g}^T \mathbf{g})^2}{(\mathbf{g}^T \mathbf{H}(\mathbf{w}^*) \mathbf{g}) (\mathbf{g}^T \mathbf{H}^{-1}(\mathbf{w}^*) \mathbf{g})} \right)$$

Kantorovich inequality:

$$\frac{(\mathbf{g}^T \mathbf{g})^2}{(\mathbf{g}^T \mathbf{H} \mathbf{g}) (\mathbf{g}^T \mathbf{H}^{-1} \mathbf{g})} \geq \frac{4 \lambda_{\min} \lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2} \geq \frac{1}{\text{cond}(\mathbf{H})}$$

condition of a matrix

$$\text{cond}(\mathbf{H}) = \frac{\lambda_{\max}}{\lambda_{\min}}$$

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Newton Method

the Newton method is quadratic convergent in  $\|w^{\text{old}} - w^*\|$  :

$$w^{\text{new}} - w^* = O\left(\|w^{\text{old}} - w^*\|^2\right)$$

assumed that

$$\left\| \frac{1}{2} H^{-1} \left( \xi_1^T H_1(w^{\text{old}}) \xi_1, \dots, \xi_W^T H_W(w^{\text{old}}) \xi_W \right)^T \right\| < 1$$

where

$$\xi_i = \lambda (w^* - w^{\text{old}}), \quad 0 \leq \lambda \leq 1,$$

and  $H_i$  is the Hessian of  $g_i$

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Until now fixed training set where update uses all examples  
→ batch update

Incremental update, where a single example is used for update  
→ on-line update

A) overfitting

Cheap training examples → overfitting is avoided if the training size is large enough (empirical risk minimization)

large training sets → computationally too expensive

→ on-line methods

B) non-stationary problems

non-stationary problems (dynamics change)

→ on-line methods can track data dependencies



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

goal: find  $\mathbf{w}^*$ , for which  $g(\mathbf{w}^*) = 0$

$$g(\mathbf{w}) = \mathbb{E}(f(\mathbf{w}) | \mathbf{w})$$

with finite variance  $\mathbb{E}((g - f)^2 | \mathbf{w}) < \infty$

## Robbins-Monro procedure

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta_i f(\mathbf{w}^i)$$

where  $g(\mathbf{w}^i)$  is a random variable and the learning rate sequence  $\eta_i$  satisfies

$$\lim_{i \rightarrow \infty} \eta_i = 0 \quad \text{convergence}$$

$$\sum_{i=1}^{\infty} \eta_i = \infty \quad \text{changes are sufficient large to find the root}$$

$$\sum_{i=1}^{\infty} \eta_i^2 < \infty \quad \text{noise variance is limited}$$

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Theorem 1 (Robbins-Monro)

*Under above conditions the above sequence converges to the root  $\mathbf{w}^*$  of  $\mathbf{g}$  with probability 1.*

Robbins-Monro for maximum likelihood:

$$\frac{1}{l} \frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^l \ln p(\mathbf{x}^i | \mathbf{w}) = 0$$

$$\lim_{l \rightarrow \infty} \frac{1}{l} \sum_{i=1}^l \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{x}^i | \mathbf{w}) = \mathbb{E} \left( \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{x}^i | \mathbf{w}) \right)$$

maximum likelihood solution is asymptotically equivalent to

$$\mathbb{E} \left( \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{x}^i | \mathbf{w}) \right) = 0$$

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Robbins-Monro procedure

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta_i \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{x}^{i+1} | \mathbf{w}) |_{\mathbf{w}^i}$$

online update formula for maximum likelihood

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Convex Problems

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \forall_i : c_i(\mathbf{x}) \leq 0 \\ & \forall_j : e_j(\mathbf{x}) = 0 \end{aligned}$$

$f, c_i$  and  $e_j$  are convex functions

solution is a convex set

$f$  is strictly convex: solution is unique

constraint convex minimization

SVM optimization problems

- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

## Lagrangian

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_i \alpha_i c_i(\mathbf{x}) + \sum_j \mu_j e_j(\mathbf{x}) \quad \alpha_i \geq 0$$

where  $\alpha$  and  $\mu$  are called “Lagrange multipliers”

Lagrangian also for non-convex functions

## 7 Optimization Techniques

### 7.1 Parameter Optimization

#### 7.1.1 Search and Evolutionary Methods

#### 7.1.2 Gradient Descent

#### 7.1.3 Step-size

#### 7.1.4 Update Direction

#### 7.1.5 Levenberg-Marquardt Algorithm

#### 7.1.6 Predictor Corrector Methods

#### 7.1.7 Convergence

### 7.2 On-line

### 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: Evidence Framework

### 8.6 Hyper-parameter: Integrate Out

### 8.7 Model Comparison

### 8.8 Posterior Sampling

## KKT conditions

feasible solution exists then the following statements are equivalent:

a) an  $\mathbf{x}$  exists with  $c_i(\mathbf{x}) < 0$  for all  $i$  (Slater's condition)

b) an  $\mathbf{x}$  and  $\alpha_i \geq 0$  exist such that  $\sum_i \alpha_i c_i(\mathbf{x}) \leq 0$

(Karlin's condition)

Above statements a) or b) follow from the following statement

c) there exist at least two feasible solutions and a feasible  $\mathbf{x}$  such that all  $c_i$  are strictly convex at  $\mathbf{x}$  wrt. the feasible set (strict constraint qualification).

The saddle point condition of Kuhn-Tucker: if one of a) – c) holds then

$$L(\hat{\mathbf{x}}, \boldsymbol{\alpha}, \boldsymbol{\mu}) \leq L(\hat{\mathbf{x}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\mu}}) \leq L(\mathbf{x}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\mu}})$$

is necessary and sufficient for  $(\hat{\mathbf{x}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\mu}})$  being a solution to the optimization problem.

Note, that "sufficient" also holds for non-convex functions

# Convex Optimization



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

From  $L(\hat{\mathbf{x}}, \boldsymbol{\alpha}, \boldsymbol{\mu}) \leq L(\hat{\mathbf{x}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\mu}})$  follows that

$$\sum_i (\alpha_i - \hat{\alpha}_i) c_i(\hat{\mathbf{x}}) + \sum_j (\mu_j - \hat{\mu}_j) e_j(\hat{\mathbf{x}}) \leq 0$$

$$\mu_j = \hat{\mu}_j, \alpha_i = \hat{\alpha}_i, \alpha_k = \hat{\alpha}_k + 1 \rightarrow c_k(\hat{\mathbf{x}}) \leq 0$$

Analog for  $e_k(\hat{\mathbf{x}}) = 0 \rightarrow \hat{\mathbf{x}}$  fulfills the constraints

$$\mu_j = \hat{\mu}_j, \alpha_i = \hat{\alpha}_i, \alpha_k = 0 \rightarrow \hat{\alpha}_k c_k(\hat{\mathbf{x}}) \geq 0$$

From above  $c_k(\hat{\mathbf{x}}) \leq 0$  and  $\hat{\alpha}_k \geq 0 \rightarrow \hat{\alpha}_k c_k(\hat{\mathbf{x}}) \leq 0$

$$\hat{\alpha}_i c_i(\hat{\mathbf{x}}) = 0$$

analog

$$\hat{\mu}_j e_j(\hat{\mathbf{x}}) = 0$$

“Karush-Kuhn-Tucker” (KKT) conditions

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

differentiable problems: minima and maxima can be determined

## Theorem 1 (KKT and Differentiable Convex Problems)

*A solution to the problem with convex, differentiable  $f$ ,  $c_i$ , and  $e_j$  is given by  $\hat{x}$  if  $\hat{\alpha}_i \geq 0$  and  $\hat{\mu}_j$  exist which satisfy:*

$$\frac{\partial L(\hat{x}, \hat{\alpha}, \hat{\mu})}{\partial x} = \frac{\partial f(\hat{x})}{\partial x} + \sum_i \hat{\alpha}_i \frac{\partial c_i(\hat{x})}{\partial x} + \sum_j \hat{\mu}_j \frac{\partial e_j(\hat{x})}{\partial x} = 0$$

$$\frac{\partial L(\hat{x}, \hat{\alpha}, \hat{\mu})}{\partial \alpha_i} = c_i(\hat{x}) \leq 0$$

$$\frac{\partial L(\hat{x}, \hat{\alpha}, \hat{\mu})}{\partial \mu_j} = e_j(\hat{x}) = 0$$

$$\forall_i : \hat{\alpha}_i c_i(\hat{x}) = 0$$

$$\forall_j : \hat{\mu}_j e_j(\hat{x}) = 0$$



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Constraints fulfilled by  $x$  and  $\alpha$ :

$$f(\mathbf{x}) \geq f(\hat{\mathbf{x}}) \geq f(\mathbf{x}) + \sum_i \alpha_i c_i(\mathbf{x}) \quad e_j(\mathbf{x}) = 0$$

$$\text{KKT gap: } \sum_i \alpha_i c_i(\mathbf{x})$$

Wolfe's dual to the optimization problem

$$\max_{\mathbf{x}, \alpha, \mu} \quad f(\mathbf{x}) + \sum_i \alpha_i c_i(\mathbf{x}) + \sum_j \mu_j e_j(\mathbf{x})$$

$$\text{s.t.} \quad \forall_i : \alpha_i \geq 0$$
$$\frac{\partial L(\mathbf{x}, \alpha, \mu)}{\partial \mathbf{x}} = 0$$

The solutions of the dual are the solutions of the primal  
If  $\frac{\partial L(\mathbf{x}, \alpha, \mu)}{\partial \mathbf{x}} = 0$  can be solved for  $\mathbf{x}$  and inserted into the dual,  
then we obtain a maximization problem in  $\alpha$  and  $\mu$

# Convex Optimization



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Linear Programs

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} + \mathbf{d} \leq \mathbf{0} \end{aligned}$$

where  $\mathbf{A} \mathbf{x} + \mathbf{d} \leq \mathbf{0}$  means  $\forall_i : \sum_{j=1}^l A_{ij} x^j + d_j \leq 0$

$$\text{Lagrangian } L = \underbrace{\mathbf{c}^T \mathbf{x} + \boldsymbol{\alpha}^T (\mathbf{A} \mathbf{x} + \mathbf{d})}_{\text{optimality conditions}}$$

optimality conditions

$$\frac{\partial L}{\partial \mathbf{x}} = \mathbf{A}^T \boldsymbol{\alpha} + \mathbf{c} = \mathbf{0}$$

$$\frac{\partial L}{\partial \boldsymbol{\alpha}} = \mathbf{A} \mathbf{x} + \mathbf{d} \leq \mathbf{0}$$

$$\boldsymbol{\alpha}^T (\mathbf{A} \mathbf{x} + \mathbf{d}) = 0$$

$$\boldsymbol{\alpha} \geq \mathbf{0}$$

dual formulation

$$\max_{\boldsymbol{\alpha}} \quad \mathbf{d}^T \boldsymbol{\alpha}$$

$$\text{s.t.} \quad \mathbf{A}^T \boldsymbol{\alpha} + \mathbf{c} = \mathbf{0}$$

$$\boldsymbol{\alpha} \geq \mathbf{0}$$

- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

dual of the dual

$$\begin{aligned} \min_{\alpha, \mu} \quad & c^T x' \\ \text{s.t.} \quad & A x' + d + \mu = 0 \\ & \mu \geq 0 \end{aligned}$$

$\mu$  can be chosen free if we ensure  $A x' + d \leq 0$

we obtain again the primal

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Quadratic Programs

$$\begin{aligned} \min_x \quad & \frac{1}{2} \mathbf{x}^T \mathbf{K} \mathbf{x} + \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} + \mathbf{d} \leq \mathbf{0} \end{aligned}$$

## Lagrangian

$$L(\mathbf{x}, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{x}^T \mathbf{K} \mathbf{x} + \mathbf{c}^T \mathbf{x} + \boldsymbol{\alpha}^T (\mathbf{A} \mathbf{x} + \mathbf{d})$$

## optimality conditions

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{x}} &= \mathbf{K} \mathbf{x} + \mathbf{A}^T \boldsymbol{\alpha} + \mathbf{c} = \mathbf{0} \\ \frac{\partial L}{\partial \boldsymbol{\alpha}} &= \mathbf{A} \mathbf{x} + \mathbf{d} \leq \mathbf{0} \\ \boldsymbol{\alpha}^T (\mathbf{A} \mathbf{x} + \mathbf{d}) &= \mathbf{0} \\ \boldsymbol{\alpha} &\geq \mathbf{0} \end{aligned}$$

## dual

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{A} \mathbf{K}^{-1} \mathbf{A}^T \boldsymbol{\alpha} \\ & - (\mathbf{d}^T - \mathbf{c}^T \mathbf{K}^{-1} \mathbf{A}^T) \boldsymbol{\alpha} \\ \text{s.t.} \quad & \mathbf{0} \leq \boldsymbol{\alpha} \end{aligned}$$

- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

## Optimization of Convex Problems

constraint gradient descent

interior point methods (interior point is a pair  $(\mathbf{x}, \alpha)$  which satisfies both the primal and dual constraints)

predictor-corrector methods: anneal KKT conditions (which are the only quadratic conditions in the variables)

# Chapter 8

## Bayes Techniques

---

7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

probabilistic framework for the empirical error and regularization

Focus: neural networks but valid for other models

Bayes techniques allow

- to introduce a probabilistic framework
- to deal with hyper-parameters
- to supply error bars and confidence intervals for the model output
- to compare different models
- to select relevant features
- to make averages and committees

# Likelihood, Prior, Posterior, Evidence



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

training data  $\{z^1, \dots, z^l\}$  ( $z^i = (x^i, y^i)$ )       $\{z\} = \{z^1, \dots, z^l\}$

matrix of feature vectors  $\mathbf{X} = (x^1, \dots, x^l)^T$

vector of labels  $\mathbf{y} = (y^1, \dots, y^l)^T$

training data matrix  $\mathbf{Z} = (z^1, \dots, z^l)$

Likelihood  $\mathcal{L}$ :  $\mathcal{L}(\{z\}; \mathbf{w}) = p(\{z\}; \mathbf{w})$

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

iid data:  $\mathcal{L}(\{z\}; \mathbf{w}) = p(\{z\}; \mathbf{w}) = \prod_{i=1}^l p(z^i; \mathbf{w})$



# Likelihood, Prior, Posterior, Evidence



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

supervised learning:  $p(\mathbf{z}; \mathbf{w}) = p(\mathbf{x}) p(y | \mathbf{x}; \mathbf{w})$

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

$\prod_{i=1}^l p(\mathbf{x}^i)$  is independent of the parameters

sufficient to maximize:  $\mathcal{L}(\{y\} | \{\mathbf{x}\}; \mathbf{w}) = \prod_{i=1}^l p(y^i | \mathbf{x}^i; \mathbf{w})$

likelihood is conditioned on  $\mathbf{w}$

**likelihood**  $p(\{\mathbf{z}\} | \mathbf{w})$

# Likelihood, Prior, Posterior, Evidence



7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

Overfitting: training examples with likelihood  
other data with probability zero  
sum of Dirac delta-distributions

avoid overfitting: some models are more likely in the world

**prior distribution:**  $p(\mathbf{w})$

from prior problem knowledge without seeing the data

# Likelihood, Prior, Posterior, Evidence



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Bayes formula 
$$p(\mathbf{w} | \{\mathbf{z}\}) = \frac{p(\{\mathbf{z}\} | \mathbf{w}) p(\mathbf{w})}{p_{\mathbf{w}}(\{\mathbf{z}\})}$$

posterior distribution:  $p(\mathbf{w} | \{\mathbf{z}\})$

evidence: 
$$p_{\mathbf{w}}(\{\mathbf{z}\}) = \int_{\mathcal{W}} p(\{\mathbf{z}\} | \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

normalization constant:

accessible volume of the configuration space (statistical mechanics)  
error moment generating function

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

# Likelihood, Prior, Posterior, Evidence



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

if  $p(\mathbf{w})$  is indeed the distribution of model parameters in real world

$$p_{\mathbf{w}}(\{\mathbf{z}\}) = p(\{\mathbf{z}\})$$

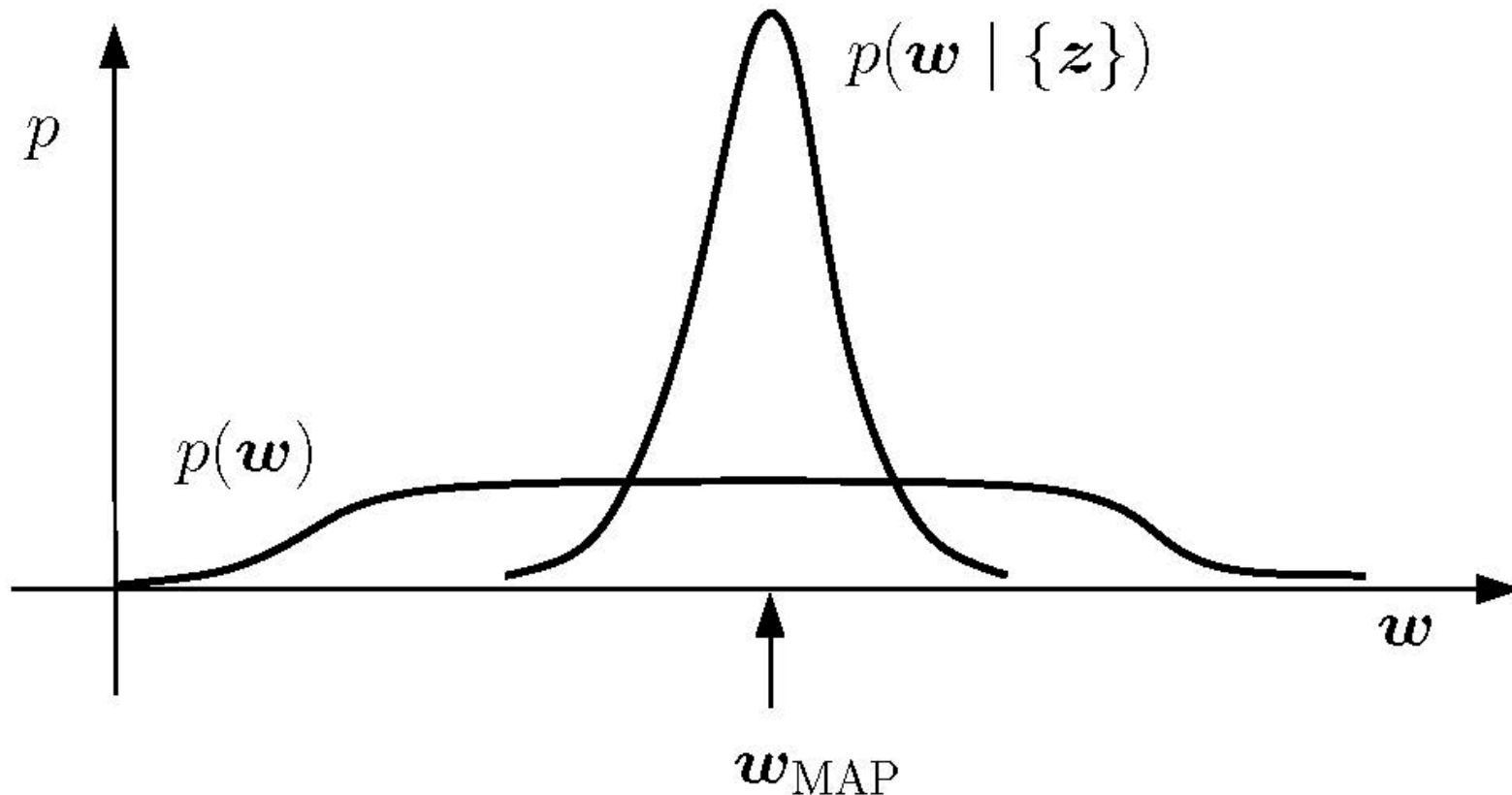
The real data is indeed produced by first choosing  $\mathbf{w}$  according to  $p(\mathbf{w})$  and then generating  $\{\mathbf{z}\}$  through  $p(\{\mathbf{z}\} | \mathbf{w})$

data in real world is not produced according to  $p(\mathbf{w})$   
therefore  $p_{\mathbf{w}}(\{\mathbf{z}\})$  is not the distribution of occurrence of data  
but gives the probability of observing data with the model class

# Maximum A Posteriori Approach

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Maximum A Posteriori Approach (MAP):  
maximal posterior  $p(w | \{z\})$



# Maximum A Posteriori Approach



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Prior: weight decay  $\Omega(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} = \sum_{ij} w_{ij}^2$

Gaussian weight prior

$$p(\mathbf{w}) = \frac{1}{Z_w(\alpha)} \exp\left(-\frac{1}{2} \alpha \|\mathbf{w}\|^2\right)$$

$$Z_w(\alpha) = \int_W \exp\left(-\frac{1}{2} \alpha \|\mathbf{w}\|^2\right) d\mathbf{w} = \left(\frac{2\pi}{\alpha}\right)^{W/2}$$

$\alpha$  hyper-parameter: trades error term against the complexity term

# Maximum A Posteriori Approach



7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

other weight decay terms:

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_1 \quad \text{Laplace distribution}$$

$$p(\mathbf{w}) = \frac{1}{Z_{\mathbf{w}}(\alpha)} \exp\left(-\frac{1}{2} \alpha \|\mathbf{w}\|_1\right)$$

$$Z_{\mathbf{w}}(\alpha) = \int_{\mathcal{W}} \exp\left(-\frac{1}{2} \alpha \|\mathbf{w}\|_1\right) d\mathbf{w}$$

$$\Omega(\mathbf{w}) = \log(1 + \|\mathbf{w}\|^2) \quad \text{Cauchy distribution}$$

$$p(\mathbf{w}) = \frac{1}{Z_{\mathbf{w}}(\alpha)} \exp\left(-\frac{1}{2} \alpha\right) (1 + \|\mathbf{w}\|^2)$$

$$Z_{\mathbf{w}}(\alpha) = \int_{\mathcal{W}} \exp\left(-\frac{1}{2} \alpha\right) (1 + \|\mathbf{w}\|^2) d\mathbf{w}$$

# Maximum A Posteriori Approach



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Gaussian noise models

$$p(\{z\} | w) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (y - Xw)^T \Sigma^{-1} (y - Xw)\right) p(\{x\})$$

and for  $\Sigma = \sigma^2 I$

$$p(\{z\} | w) = \frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)\right) p(\{x\})$$

$$R_{\text{emp}} = (y - Xw)^T (y - Xw) \text{ mean squared error}$$



# Maximum A Posteriori Approach



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

negative log-posterior

$$-\log p(\mathbf{w} | \{\mathbf{z}\}) = -\log p(\{\mathbf{z}\} | \mathbf{w}) - \log p(\mathbf{w}) + \log p_{\mathbf{w}}(\{\mathbf{z}\})$$

where  $p_{\mathbf{w}}(\{\mathbf{z}\})$  does not depend on  $\mathbf{w}$

maximum a posteriori

$$\tilde{R}(\mathbf{w}) = \frac{1}{2\sigma^2} R_{\text{emp}} + \frac{1}{2}\alpha \Omega(\mathbf{w}) = \frac{1}{2}\beta R_{\text{emp}} + \frac{1}{2}\alpha \Omega(\mathbf{w})$$

where  $\beta^{-1} = \sigma^2$

$$R(\mathbf{w}) = \tilde{R}(\mathbf{w}) 2\sigma^2 \qquad \lambda = \sigma^2 \alpha$$

$$R(\mathbf{w}) = R_{\text{emp}} + \lambda \Omega(\mathbf{w}) \qquad \text{this is exactly weight decay}$$

empirical error plus complexity: maximum a posteriori

# Maximum A Posteriori Approach



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Likelihood: exponential function of empirical error

$$p(\{z\} | \mathbf{w}) = \frac{1}{Z_R(\beta)} \exp\left(-\frac{1}{2} \beta R_{\text{emp}}\right)$$

Prior: exponential function of the complexity

$$p(\mathbf{w}) = \frac{1}{Z_w(\alpha)} \exp\left(-\frac{1}{2} \alpha \Omega(\mathbf{w})\right)$$

Posterior:

$$p(\mathbf{w} | \{z\}) = \frac{1}{Z(\alpha, \beta)} \exp\left(-\frac{1}{2} (\alpha \Omega(\mathbf{w}) + \beta R_{\text{emp}})\right)$$

where

$$Z(\alpha, \beta) = \int_{\mathcal{W}} \exp\left(-\frac{1}{2} (\alpha \Omega(\mathbf{w}) + \beta R_{\text{emp}})\right) d\mathbf{w}$$

# Posterior Approximation



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

approximate the posterior: Gaussian assumption

$$\tilde{R}(\mathbf{w}) = \frac{1}{2\sigma^2} R_{\text{emp}} + \frac{1}{2}\alpha \Omega(\mathbf{w}) = \frac{1}{2}\beta R_{\text{emp}} + \frac{1}{2}\alpha \Omega(\mathbf{w})$$

Taylor expansion of  $R(\mathbf{w})$  around its minimum  $\mathbf{w}_{\text{MAP}}$

$$\tilde{R}(\mathbf{w}) = \tilde{R}(\mathbf{w}_{\text{MAP}}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{H} (\mathbf{w} - \mathbf{w}_{\text{MAP}})$$

where the first order derivatives vanish at the minimum  
 $\mathbf{H}$  is the Hessian

posterior is now a Gaussian

$$p(\mathbf{w} | \{z\}) = \frac{1}{Z} \exp(-\tilde{R}(\mathbf{w})) = \frac{1}{Z} \exp\left(-\tilde{R}(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{H} (\mathbf{w} - \mathbf{w}_{\text{MAP}})\right)$$

where  $Z$  is normalization constant

# Posterior Approximation



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Hessian for weight decay

$$\mathbf{H} = \frac{1}{2 \sigma^2} \mathbf{H}_{\text{emp}} + \alpha \mathbf{I} = \frac{\beta}{2} \mathbf{H}_{\text{emp}} + \alpha \mathbf{I}$$

where  $\mathbf{H}_{\text{emp}}$  is the Hessian of the empirical error (the factor  $\frac{1}{2}$  vanishes for the weight decay term)

normalization constant:

$$Z(\alpha, \beta) = \exp\left(-\tilde{R}(\mathbf{w}_{\text{MAP}})\right) (2 \pi)^{W/2} |\mathbf{H}|^{1/2}$$

# Error Bars and Confidence Intervals



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

confidence intervals for model outputs  
how reliable is the prediction

$$\text{Output distribution } p(y | \mathbf{x}, \{z\}) = \int_{\mathcal{W}} p(y | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \{z\}) d\mathbf{w}$$

where we used the posterior distribution  $p(\mathbf{w} | \{z\})$   
and a noise model  $p(y | \mathbf{x}, \mathbf{w})$

Gaussian noise model (one dimension)

$$p(y | \mathbf{x}, \mathbf{w}) = \frac{1}{Z_R(\beta)} \exp\left(-\frac{\beta}{2} (g(\mathbf{x}; \mathbf{w}) - y)^2\right)$$

$$Z_R(\beta) = \left(\frac{2\pi}{\beta}\right)^{1/2} \quad \beta = \frac{1}{\sigma^2}$$

# Error Bars and Confidence Intervals



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Linear approximation

$$g(\mathbf{x}; \mathbf{w}) = g(\mathbf{x}; \mathbf{w}_{\text{MAP}}) + \mathbf{g}^T (\mathbf{w} - \mathbf{w}_{\text{MAP}})$$

$$p(y | \mathbf{x}, \{\mathbf{z}\}) \propto$$

$$\int_{\mathbf{w}} \exp \left( - \underbrace{\frac{\beta}{2} (y - g(\mathbf{x}; \mathbf{w}_{\text{MAP}}) + \mathbf{g}^T (\mathbf{w} - \mathbf{w}_{\text{MAP}}))^2}_{p(y | \mathbf{x}, \mathbf{w})} - \underbrace{\frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{H} (\mathbf{w} - \mathbf{w}_{\text{MAP}})}_{p(\mathbf{w} | \{\mathbf{z}\})} \right) d\mathbf{w}$$

Using the fact that the integrand is a Gaussian in  $\mathbf{w}$ , this gives

$$p(y | \mathbf{x}, \{\mathbf{z}\}) = \frac{1}{\sqrt{2\pi}\sigma_y} \exp \left( - \frac{1}{2\sigma_y^2} (y - g(\mathbf{x}; \mathbf{w}_{\text{MAP}}))^2 \right)$$

$$\text{where } \sigma_y^2 = \frac{1}{\beta} + \mathbf{g}^T \mathbf{H}^{-1} \mathbf{g} = \sigma^2 + \mathbf{g}^T \mathbf{H}^{-1} \mathbf{g}$$

inherent data noise  $\sigma^2$  plus the approximation uncertainty  $\mathbf{g}^T \mathbf{H}^{-1} \mathbf{g}$

# Error Bars and Confidence Intervals



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

To derive the formula for this Gaussian, we use following identity:

$$\int_{\mathbf{w}} \exp\left(-\frac{1}{2}\mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{v}^T \mathbf{w}\right) d\mathbf{w} = (2\pi)^{W/2} |\mathbf{A}|^{-1/2} \exp\left(\frac{1}{2}\mathbf{v}^T \mathbf{A}^{-1} \mathbf{v}\right)$$

After collecting terms in  $(\mathbf{w} - \mathbf{w}_{\text{MAP}})$  the quadratic part is

$$\mathbf{A} = \beta \mathbf{g} \mathbf{g}^T + \mathbf{H}$$

and the linear part is

$$\mathbf{v} = \beta (y - g(\mathbf{x}; \mathbf{w}_{\text{MAP}})) \mathbf{g}$$

and the constant term is

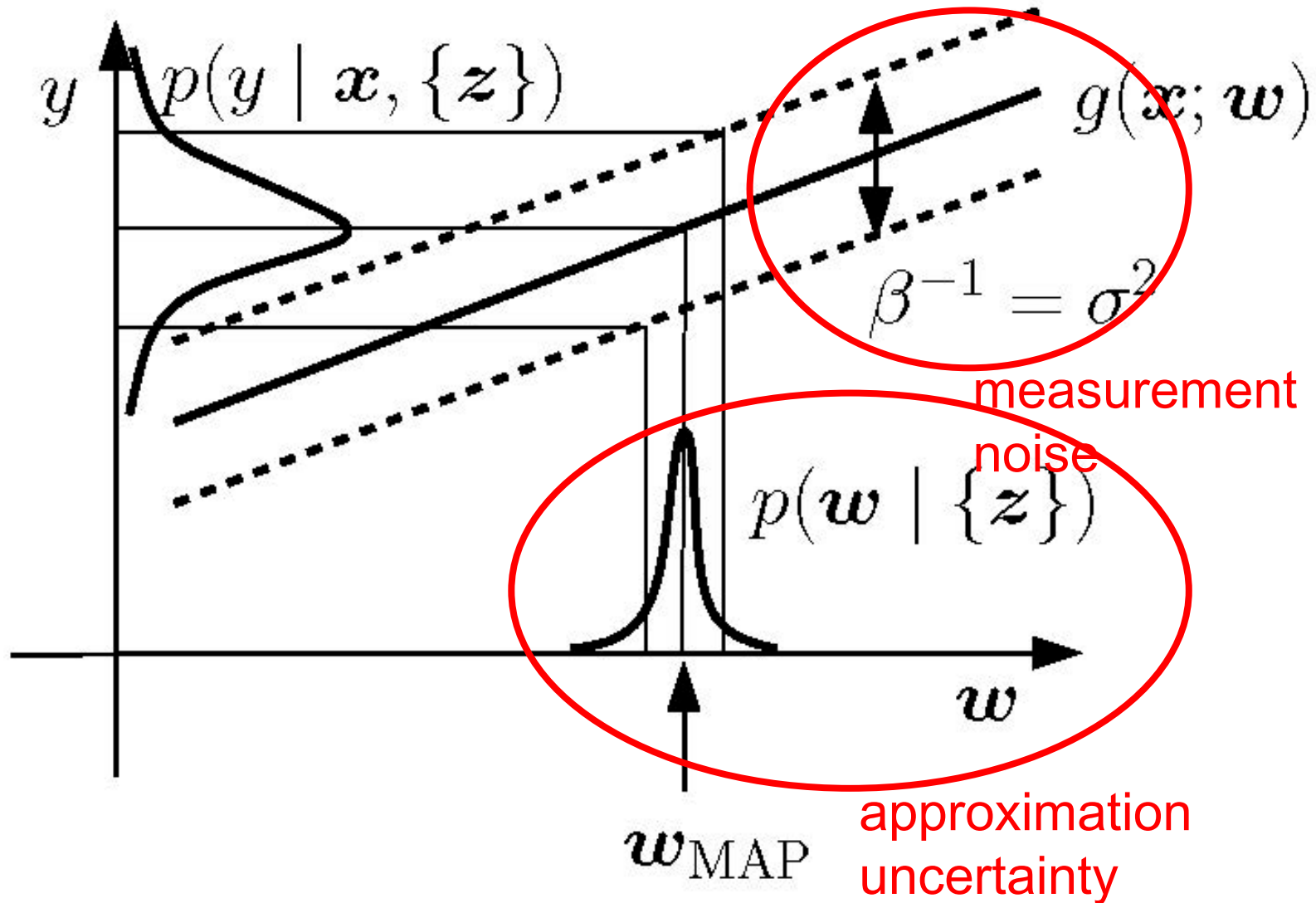
$$c = \beta (y - g(\mathbf{x}; \mathbf{w}_{\text{MAP}}))^2$$

We obtain in the exponent:

$$-\frac{1}{2}\mathbf{v}^T \mathbf{A}^{-1} \mathbf{v} + \frac{1}{2}c = \frac{1}{2} (y - g(\mathbf{x}; \mathbf{w}_{\text{MAP}}))^2 \left(\beta - \beta^2 \mathbf{g}^T (\mathbf{H} + \beta \mathbf{g} \mathbf{g}^T)^{-1} \mathbf{g}\right)^{-1} = \frac{1}{\beta} + \mathbf{g}^T \mathbf{H}^{-1} \mathbf{g}$$

# Error Bars and Confidence Intervals

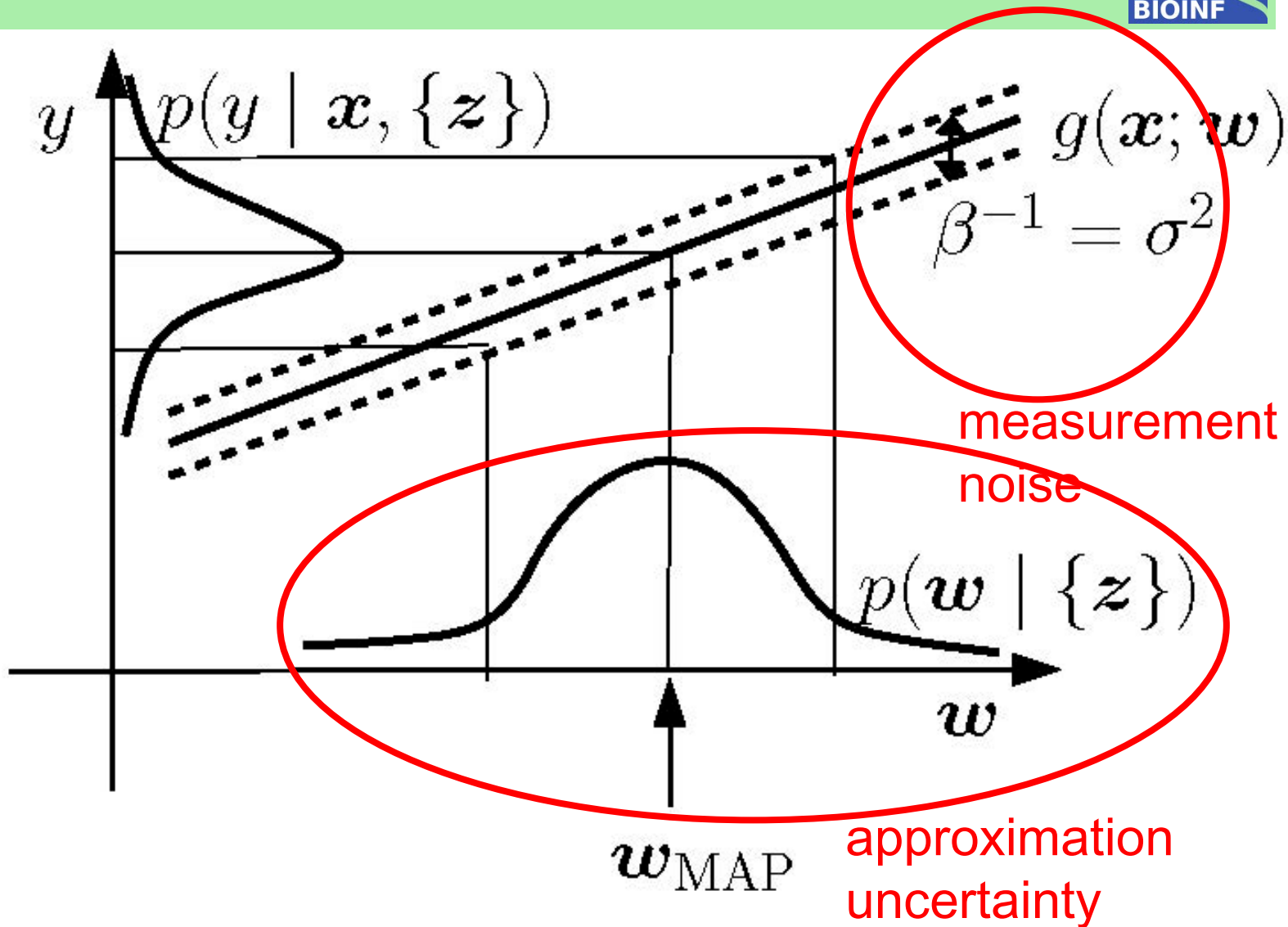
- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling





# Error Bars and Confidence Intervals

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling



# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

hyper-parameters  $\alpha$  and  $\beta$

$\beta$ : assumption on the noise in the data

$\alpha$ : assumption on the optimal complexity relative to empirical error

integrating out  $\alpha$  and  $\beta \rightarrow$  **marginalization**

$$p(\mathbf{w} \mid \{\mathbf{z}\}) =$$

$$\int_{S_\alpha} \int_{S_\beta} p(\mathbf{w}, \alpha, \beta \mid \{\mathbf{z}\}) d\alpha d\beta =$$

$$\int_{S_\alpha} \int_{S_\beta} p(\mathbf{w} \mid \alpha, \beta, \{\mathbf{z}\}) p(\alpha, \beta \mid \{\mathbf{z}\}) d\alpha d\beta$$

compute the integrals  $\rightarrow$  later in next section

# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

approximate the posterior

assume: sharply peaked around the maximal values  $\alpha_{\text{MAP}}$  and  $\beta_{\text{MAP}}$

high values of  $p(\alpha, \beta | \{z\})$  and constant  $p(\mathbf{w} | \alpha, \beta, \{z\})$  with value  $p(\mathbf{w} | \alpha_{\text{MAP}}, \beta_{\text{MAP}}, \{z\})$

$$p(\mathbf{w} | \{z\}) = p(\mathbf{w} | \alpha_{\text{MAP}}, \beta_{\text{MAP}}, \{z\})$$

$$\int_{S_\alpha} \int_{S_\beta} p(\alpha, \beta | \{z\}) d\alpha d\beta = p(\mathbf{w} | \alpha_{\text{MAP}}, \beta_{\text{MAP}}, \{z\})$$

searching for the hyper-parameters which maximize the posterior

# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

posterior of

$$p(\alpha, \beta | \{z\}) = \frac{p(\{z\} | \alpha, \beta) p(\alpha, \beta)}{p_{\alpha, \beta}(\{z\})}$$

prior for  $\alpha$  and  $\beta$ :  $p(\alpha, \beta)$

**non-informative priors**: equal probability for all values

marginalization over  $w$ :

$$p(\{z\} | \alpha, \beta) = \int_W p(\{z\} | w, \alpha, \beta) p(w | \alpha, \beta) dw =$$

$$\int_W p(\{z\} | w, \beta) p(w | \alpha) dw$$

Def. conditional probability:

$$p(\{z\} | w) p(w) = p(\{z\}) p(w | \{z\})$$

$$p(\{z\} | \alpha, \beta) = \frac{Z(\alpha, \beta)}{Z_w(\alpha) Z_R(\beta)}$$

$$p(\{z\} | w) = \frac{1}{Z_R(\beta)} \exp\left(-\frac{1}{2} \beta R_{\text{emp}}\right)$$

$$p(w) = \frac{1}{Z_w(\alpha)} \exp\left(-\frac{1}{2} \alpha \Omega(w)\right)$$

$$p(w | \{z\}) = \frac{1}{Z(\alpha, \beta)} \exp\left(-\frac{1}{2} (\alpha \Omega(w) + \beta R_{\text{emp}})\right)$$

$$p(w | \{z\}) = \frac{Z_R(\beta) Z_w(\alpha)}{Z(\alpha, \beta)} p(\{z\} | w) p(w)$$

$$p(\{z\} | w) p(w) = \frac{Z(\alpha, \beta)}{Z_R(\beta) Z_w(\alpha)} p(w | \{z\})$$

# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

## Example

mean squared error and regularization by quadratic weight decay

$$Z(\alpha, \beta) = \exp\left(-\tilde{R}(\mathbf{w}_{\text{MAP}})\right) (2\pi)^{W/2} |\mathbf{H}|^{-1/2}$$

where  $\tilde{R}(\mathbf{w}_{\text{MAP}}) = \frac{1}{2} \beta R_{\text{emp}} + \frac{1}{2} \alpha \Omega(\mathbf{w}_{\text{MAP}})$

$$Z_R(\beta) = \left(\frac{2\pi}{\beta}\right)^{l/2}$$

$$Z_w(\alpha) = \left(\frac{2\pi}{\alpha}\right)^{W/2}$$

thus

$$\ln p(\{\mathbf{z}\} | \alpha, \beta) = -\frac{1}{2} \alpha \Omega(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} \beta R_{\text{emp}} - \frac{1}{2} \ln |\mathbf{H}| +$$

$$\frac{W}{2} \ln \alpha + \frac{l}{2} \ln \beta - \frac{l}{2} \ln(2\pi)$$

$(2\pi)^{W/2}$  cancels

$$\mathbf{H} = \frac{\beta}{2} \mathbf{H}_{\text{emp}} + \alpha \mathbf{I}$$

$$\tilde{R}(\mathbf{w}) = \frac{1}{2} \beta R_{\text{emp}}(\mathbf{w}) + \frac{1}{2} \alpha \Omega(\mathbf{w})$$

$$\tilde{R}(\mathbf{w}) = \tilde{R}(\mathbf{w}_{\text{MAP}}) + (\mathbf{w} - \mathbf{w}_{\text{MAP}})^T \mathbf{H} (\mathbf{w} - \mathbf{w}_{\text{MAP}})$$

- first order derivative is zero
- $\exp(-R(\mathbf{w}))$  is a Gaussian with normalizing constant  $Z(\alpha, \beta)$

# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Assumption: eigenvalues  $\lambda_j$  of  $1/2 \mathbf{H}_{\text{emp}}$  do not depend on  $\alpha$

$$\frac{\partial}{\partial \alpha} \ln |\mathbf{H}| = \frac{\partial}{\partial \alpha} \ln \prod_{j=1}^W (\beta \lambda_j + \alpha) =$$

$$\frac{\partial}{\partial \alpha} \sum_{j=1}^W \ln (\beta \lambda_j + \alpha) = \sum_{j=1}^W \frac{1}{\beta \lambda_j + \alpha} = \text{Tr} \mathbf{H}^{-1}$$

Note, Hessian  $\mathbf{H}$  was evaluated at  $w_{\text{MAP}}$  which depends on  $\alpha$   
→ terms  $\frac{\partial \lambda_j}{\partial \alpha}$  neglected

# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

$$\frac{\partial}{\partial \alpha} \ln p(\{z\} | \alpha, \beta) = -\frac{1}{2} \Omega(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} \sum_{j=1}^W \frac{1}{\beta \lambda_j + \alpha} + \frac{1}{2} W \frac{1}{\alpha} = 0$$

which gives

$$\alpha \Omega(\mathbf{w}_{\text{MAP}}) = -\sum_{j=1}^W \frac{\alpha}{\beta \lambda_j + \alpha} + W = \sum_{j=1}^W \frac{\beta \lambda_j}{\beta \lambda_j + \alpha} = \gamma$$

$\Omega(\mathbf{w}_{\text{MAP}})$  how far are weights pushed away from prior (zero) by data term  $\frac{\beta \lambda_j}{\beta \lambda_j + \alpha}$  is in  $[0;1]$ ; 1  $\rightarrow$  data governed; 0  $\rightarrow$  prior driven

$\gamma$ : effective number of weights which are driven by the data

# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Note, Hessian is not evaluated at the minimum of  $R_{\text{emp}}$  but at  $w_{\text{MAP}}$

eigenvalues  $\lambda_j$  are not guaranteed to be positive

→ terms  $\frac{\beta \lambda_j}{\beta \lambda_j + \alpha}$  may be negative

derivative of the negative log-posterior with respect to  $\beta$

$$\frac{\partial}{\partial \beta} \ln |\mathbf{H}| = \frac{\partial}{\partial \beta} \sum_{j=1}^W \ln (\beta \lambda_j + \alpha) =$$

$$\sum_{j=1}^W \frac{\lambda_j}{\beta \lambda_j + \alpha}$$

derivative set to zero:  $\beta R_{\text{emp}} = l - \sum_{j=1}^W \frac{\lambda_j}{\beta \lambda_j + \alpha} = l - \gamma$



# Hyper-parameter Selection: Evidence Framework



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

updates for the hyper-parameters:

$$\alpha^{\text{new}} = \frac{\gamma}{\Omega(\mathbf{w}_{\text{MAP}})}$$
$$\beta^{\text{new}} = \frac{l - \gamma}{R_{\text{emp}}(\mathbf{w}_{\text{MAP}})}$$

with these new hyper-parameters the new values of  $\mathbf{w}_{\text{MAP}}$  can be estimated through gradient based methods

Iteratively again the hyper-parameters are updated and so forth

If all parameters are well defined:  $\gamma = W$

If many training examples:  $l \gg W$

approximation for update formulae

$$\alpha^{\text{new}} = \frac{W}{\Omega(\mathbf{w}_{\text{MAP}})}$$
$$\beta^{\text{new}} = \frac{l}{R_{\text{emp}}(\mathbf{w}_{\text{MAP}})}$$

# Hyper-parameter Selection: Integrate Out



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

Posterior: obtained by integrating out  $\alpha$  and  $\beta$

$$p(\mathbf{w} | \{\mathbf{z}\}) = \int_{S_\alpha} \int_{S_\beta} p(\mathbf{w}, \alpha, \beta | \{\mathbf{z}\}) d\alpha d\beta =$$
$$\int_{S_\alpha} \int_{S_\beta} p(\mathbf{w} | \alpha, \beta, \{\mathbf{z}\}) p(\alpha, \beta | \{\mathbf{z}\}) d\alpha d\beta =$$
$$\frac{1}{p_{\mathbf{w}}(\{\mathbf{z}\})} \int_{S_\alpha} \int_{S_\beta} p(\{\mathbf{z}\} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) p(\alpha) p(\beta) d\alpha d\beta$$

where we used

$$p(\alpha, \beta | \{\mathbf{z}\}) = p(\alpha) p(\beta)$$

$$p(\mathbf{w} | \{\mathbf{z}\}) = \frac{p(\{\mathbf{z}\} | \mathbf{w}) p(\mathbf{w})}{p_{\mathbf{w}}(\{\mathbf{z}\})}$$

$$p(\{\mathbf{z}\} | \mathbf{w}, \alpha, \beta) = p(\{\mathbf{z}\} | \mathbf{w}, \beta)$$

$$p(\mathbf{w} | \alpha, \beta) = p(\mathbf{w} | \alpha)$$

# Hyper-parameter Selection: Integrate Out



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

$\alpha$  and  $\beta$  are scaling parameters

- output range:  $\beta$
- weights:  $\alpha$

Non-informative (uniform) on a logarithmic scale, therefore  $p(\ln(\alpha))$  and  $p(\ln(\beta))$  are constant

Using  $p_x(\mathbf{x}) = p_g(g(\mathbf{x})) \left| \frac{\partial g}{\partial \mathbf{x}} \right|$  gives

$$p(\alpha) = \frac{1}{\alpha}$$
$$p(\beta) = \frac{1}{\beta}$$

# Hyper-parameter Selection: Integrate Out



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

prior over the weights

$$\begin{aligned} p(\mathbf{w}) &= \int_0^\infty p(\mathbf{w} | \alpha) p(\alpha) d\alpha = \\ &= \int_0^\infty \frac{1}{Z_{\mathbf{w}}(\alpha)} \exp\left(-\frac{1}{2} \alpha \Omega(\mathbf{w})\right) \frac{1}{\alpha} d\alpha = \\ &= (2\pi)^{-W/2} \int_0^\infty \exp\left(-\frac{1}{2} \alpha \Omega(\mathbf{w})\right) \alpha^{W/2-1} d\alpha = \frac{\Gamma(W/2)}{(2\pi \Omega(\mathbf{w}))^{W/2}} \end{aligned}$$

where  $\Gamma$  is the gamma function

$$\text{Analog } p(\{\mathbf{z}\} | \mathbf{w}) = \frac{\Gamma(l/2)}{(2\pi R_{\text{emp}})^{l/2}}$$

# Hyper-parameter Selection: Integrate Out



Bayes formula  $\rightarrow$  negative log-posterior:

$$-\log p(\mathbf{w} | \{\mathbf{z}\}) = -\log p(\{\mathbf{z}\} | \mathbf{w}) - \log p(\mathbf{w}) + \log p_{\mathbf{w}}(\{\mathbf{z}\})$$

$$-\ln p(\mathbf{w} | \{\mathbf{z}\}) = \frac{l}{2} \ln R_{\text{emp}} + \frac{W}{2} \ln \Omega(\mathbf{w}) + \text{const}$$

$$\text{compare with } -\ln p(\mathbf{w} | \{\mathbf{z}\}) = \frac{1}{2} \beta R_{\text{emp}} + \frac{1}{2} \alpha \Omega(\mathbf{w})$$

Compute gradients with respect to parameters!

Comparing the coefficients on front of the gradients gives:

$$\alpha^{\text{new}} = \frac{W}{\Omega(\mathbf{w}_{\text{MAP}})}$$

$$\beta^{\text{new}} = \frac{l}{R_{\text{emp}}(\mathbf{w}_{\text{MAP}})}$$

iterative algorithm

$\mathbf{w}_{\text{MAP}}$  gradient descent

7 Optimization Techniques  
7.1 Parameter Optimization  
7.1.1 Search and Evolutionary Methods  
7.1.2 Gradient Descent  
7.1.3 Step-size  
7.1.4 Update Direction  
7.1.5 Levenberg-Marquardt Algorithm  
7.1.6 Predictor Corrector Methods  
7.1.7 Convergence  
7.2 On-line  
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: Evidence Framework  
8.6 Hyper-parameter: Integrate Out  
8.7 Model Comparison  
8.8 Posterior Sampling

- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

compare model classes  $\mathcal{M}$

Bayes formula 
$$p(\mathcal{M} | \{z\}) = \frac{p(\{z\} | \mathcal{M}) p(\mathcal{M})}{p_{\mathcal{M}}(\{z\})}$$

evidence for model class

$$p(\{z\} | \mathcal{M}) = \int_{\mathcal{W}} p(\{z\} | \mathbf{w}, \mathcal{M}) p(\mathbf{w} | \mathcal{M}) d\mathbf{w}$$

Approximate posterior by a box around MAP value

$$p(\{z\} | \mathcal{M}) \approx p(\{z\} | \mathbf{w}_{\text{MAP}}, \mathcal{M}) p(\mathbf{w}_{\text{MAP}} | \mathcal{M}) \Delta \mathbf{w}_{\text{MAP}}$$

Gaussian:  $\Delta \mathbf{w}_{\text{MAP}}$  estimated from Hessian  $\mathbf{H}$

# Model Comparison



- 7 Optimization Techniques
- 7.1 Parameter Optimization
- 7.1.1 Search and Evolutionary Methods
- 7.1.2 Gradient Descent
- 7.1.3 Step-size
- 7.1.4 Update Direction
- 7.1.5 Levenberg-Marquardt Algorithm
- 7.1.6 Predictor Corrector Methods
- 7.1.7 Convergence
- 7.2 On-line
- 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: Evidence Framework
- 8.6 Hyper-parameter: Integrate Out
- 8.7 Model Comparison
- 8.8 Posterior Sampling

$$\ln p(\{\mathbf{z}\} \mid \alpha, \beta) = -\frac{\alpha}{2} \Omega(\mathbf{w}_{\text{MAP}}) - \frac{\beta}{2} R_{\text{emp}} - \frac{1}{2} \ln |\mathbf{H}| + \frac{W}{2} \ln \alpha + \frac{l}{2} \ln \beta - \frac{l}{2} \ln(2\pi)$$

A more exact term (not derived) is

$$\ln p(\{\mathbf{z}\} \mid \mathcal{M}) = -\frac{\alpha_{\text{MAP}}}{2} \Omega(\mathbf{w}_{\text{MAP}}) - \frac{\beta_{\text{MAP}}}{2} R_{\text{emp}} - \frac{1}{2} \ln |\mathbf{H}| + \frac{W}{2} \ln \alpha_{\text{MAP}} + \frac{l}{2} \ln \beta_{\text{MAP}} + \ln M! + 2 \ln M + \frac{1}{2} \ln \left( \frac{2}{\gamma} \right) + \frac{1}{2} \ln \left( \frac{2}{l - \gamma} \right)$$

$M$ : number of hidden units in the network  
→ posterior is only locally approximated, thus equivalent solutions in weights space exist  
→ permutation invariant solutions

# Posterior Sampling



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

Approximate integrals like  $A(f) = \int_{\mathcal{W}} f(\mathbf{w}) p(\mathbf{w} | \{\mathbf{z}\}) d\mathbf{w}$

sample weight vectors  $\mathbf{w}_i$  according to  $p(\mathbf{w} | \{\mathbf{z}\})$

$$A(f) \approx \frac{1}{L} \sum_{i=1}^L f(\mathbf{w}_i)$$

we cannot easily sample from  $p(\mathbf{w} | \{\mathbf{z}\})$

simpler distribution  $q(\mathbf{w})$  where we can sample from

$$A(f) = \int_{\mathcal{W}} f(\mathbf{w}) \frac{p(\mathbf{w} | \{\mathbf{z}\})}{q(\mathbf{w})} q(\mathbf{w}) d\mathbf{w}$$

$$A(f) \approx \frac{1}{L} \sum_{i=1}^L f(\mathbf{w}_i) \frac{p(\mathbf{w}_i | \{\mathbf{z}\})}{q(\mathbf{w}_i)}$$

$\mathbf{w}_i$  are sampled according to  $q(\mathbf{w})$



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

avoid the normalization of  $p(\mathbf{w} \mid \{\mathbf{z}\})$

$$A(f) \approx \frac{\sum_{i=1}^L f(\mathbf{w}_i) \tilde{p}(\mathbf{w}_i \mid \{\mathbf{z}\}) / q(\mathbf{w}_i)}{\sum_{i=1}^L \tilde{p}(\mathbf{w}_i \mid \{\mathbf{z}\}) / q(\mathbf{w}_i)}$$

$\tilde{p}(\mathbf{w} \mid \{\mathbf{z}\})$ : unnormalized posterior  $\rightarrow$  product likelihood and prior

## importance sampling

sample in regions with large probability mass  
 $\rightarrow$  **Markov Chain Monte Carlo**

# Posterior Sampling



- 7 Optimization Techniques
  - 7.1 Parameter Optimization
    - 7.1.1 Search and Evolutionary Methods
    - 7.1.2 Gradient Descent
    - 7.1.3 Step-size
    - 7.1.4 Update Direction
    - 7.1.5 Levenberg-Marquardt Algorithm
    - 7.1.6 Predictor Corrector Methods
    - 7.1.7 Convergence
  - 7.2 On-line
  - 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: Evidence Framework
  - 8.6 Hyper-parameter: Integrate Out
  - 8.7 Model Comparison
  - 8.8 Posterior Sampling

**Metropolis algorithm:** random walk in a way that regions with large probability mass are sampled

$w^{\text{new}} = w^{\text{candidate}}$  with probability

$$\begin{cases} 1 & \text{if } p(w^{\text{candidate}} | \{z\}) > p(w^{\text{old}} | \{z\}) \\ \frac{p(w^{\text{candidate}} | \{z\})}{p(w^{\text{old}} | \{z\})} & \text{if } p(w^{\text{candidate}} | \{z\}) < p(w^{\text{old}} | \{z\}) \end{cases}$$

simulated annealing: can be used to estimate the expectation