

# Machine Learning: Foundations

## (NGCM Introductory Course)

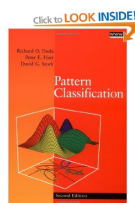
Mahezan Niranjan

School of Electronics and Computer Science  
University of Southampton

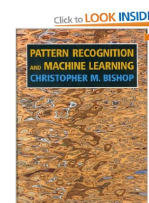
July 2019

## Good Books

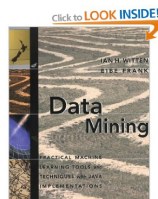
Well worth investing in books if you want to do more in this subject



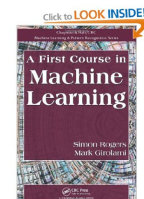
R.O.Duda, P.E.Hart & D.G.Stork  
Pattern Classification



C.M. Bishop  
Pattern Recognition and Machine Learning



I.H. Witten & E. Frank  
Data Mining



S. Rogers & M. Girolami  
A First Course in Machine Learning

*"There is nothing to be learnt from a professor, which is not to be met with in books"*  
- David Hume (1711-1776)

(**Wikipedia:** "Hume had little respect for the professors of his time [...] He did not graduate")

# Machine Learning as Data-driven Modelling

Data	$\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$	$\{\mathbf{x}_n\}_{n=1}^N$
Function Approximator	$\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + v$	
Parameter Estimation	$E_0 = \sum_{n=1}^N \{\ \mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\ \}^2$	
Prediction	$\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$	
Regularization	$E_1 = \sum_{n=1}^N \{\ \mathbf{y}_n - f(\mathbf{x}_n)\ \}^2 + g(\ \boldsymbol{\theta}\ )$	
Modelling Uncertainty	$p(\boldsymbol{\theta}   \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N)$	
Probabilistic Inference	$\mathbf{E}[g(\boldsymbol{\theta})] = \int g(\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{N_s} \sum_{n=1}^{N_s} g(\boldsymbol{\theta}^{(n)})$	
Sequential Estimation	$\boldsymbol{\theta}(n-1 n-1) \longrightarrow \boldsymbol{\theta}(n n-1) \longrightarrow \boldsymbol{\theta}(n n)$ Kalman & Particle Filters; Reinforcement Learning	

## Gaussian Densities: Univariate and Multivariate

- Univariate Gaussian

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{(x - m)^2}{\sigma^2}\right\}$$

- Multivariate Gaussian

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{p/2}(\det \mathbf{C})^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mathbf{m})^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m})\right\}$$

Mean  $\mathbf{m}$  is a vector

Covariance,  $\mathbf{C}$ , matrix: symmetric, positive definite!

$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{C}), \mathbf{y} = \mathbf{A}\mathbf{x} \implies \mathbf{y} \sim \mathcal{N}(\mathbf{A}\mathbf{m}, \mathbf{A}\mathbf{C}\mathbf{A}^T)$

---

Univariate Mean  $\hat{m} = \frac{1}{N} \sum_{n=1}^N x_n$

Univariate Covariance  $\hat{\sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{m})^2$

Multivariate Mean  $\hat{\mathbf{m}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$

Covariance Matrix  $\hat{\mathbf{C}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\mathbf{m}})(\mathbf{x}_n - \hat{\mathbf{m}})^T$

---

- These are known as maximum likelihood estimates (see later).

## Bayesian Decision Theory

- Classes:  $\omega_i, i = 1, \dots, K$
- Prior Probabilities:  $P[\omega_1], \dots, P[\omega_K];$   
 $P[\omega_i] \geq 0, \sum_{i=1}^K P[\omega_i] = 1$
- Likelihoods (class conditional probabilities):  $p(\mathbf{x}|\omega_i), i = 1, \dots, K$
- Posterior Probability:  $P[\omega_j | \mathbf{x}]$

$$P[\omega_j | \mathbf{x}] = \frac{p(\mathbf{x} | \omega_j) P[\omega_j]}{\sum_{i=1}^K p(\mathbf{x} | \omega_i) P[\omega_i]}$$

- From prior knowledge:  $P[\omega_i];$  From training data:  $p(\mathbf{x}|\omega_i)$
- Decision rule: Assign  $\mathbf{x}$  to the class that maximizes posterior probability.
- The denominator is a constant; i.e. does not depend on  $\omega_j$
- Hence the decision rule becomes:

$$\mathbf{x} \in \max_j p(\mathbf{x} | \omega_j) P[\omega_j]$$

# Bayes' Classifier for Gaussian Densities

Make assumptions, cancel common terms when making comparisons...

- Decision rule from:  $p(\mathbf{x} | \omega_j) P[\omega_j]$
- Assume the two classes are Gaussian distributed with distinct means and identical covariance matrices  
 $p(\mathbf{x} | \omega_j) = \mathcal{N}(\mathbf{m}_j, \mathbf{C})$
- Substitute into Bayes' classifier decision rule

$$\begin{aligned} P[\omega_1 | \mathbf{x}] &\leq P[\omega_2 | \mathbf{x}] \\ p(\mathbf{x} | \omega_1) P[\omega_1] &\leq p(\mathbf{x} | \omega_2) P[\omega_2] \end{aligned}$$

$$\begin{aligned} \frac{1}{(2\pi)^{p/2}(\det(\mathbf{C}))^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_1) \right\} P[\omega_1] &\leq \\ \frac{1}{(2\pi)^{p/2}(\det(\mathbf{C}))^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_2) \right\} P[\omega_2] \end{aligned}$$

## Bayes' classifier for simple densities (cont'd)

Distinct Means; Equal, isotropic covariance matrix

- Suppose the densities are isotropic and priors are equal  
*i.e.*  $\mathbf{C} = \sigma^2 \mathbf{I}$  and  $P[\omega_1] = P[\omega_2]$
- The comparison simplifies to (see algebra on board):

$$\begin{aligned} (\mathbf{x} - \mathbf{m}_1)^t (\mathbf{x} - \mathbf{m}_1) &\leq (\mathbf{x} - \mathbf{m}_2)^t (\mathbf{x} - \mathbf{m}_2) \\ |\mathbf{x} - \mathbf{m}_1| &\leq |\mathbf{x} - \mathbf{m}_2| \end{aligned}$$

- The above is a simple *distance to mean* classifier
- Under the above simplistic assumptions, we only need to store one template per class (the means)!

# Bayes' classifier for simple densities (cont'd)

Distinct Means; Common covariance matrix (but not isotropic)

- Cancel common terms and take log

$$(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_1) \leq (\mathbf{x} - \mathbf{m}_2)^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_2) - \log \left\{ \frac{P[\omega_1]}{P[\omega_2]} \right\}$$

- Also simplifies to a linear classifier!

$$\mathbf{w}^t \mathbf{x} + b \leq 0$$

$$\mathbf{w} = 2\mathbf{C}^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

$$b = (\mathbf{m}_1^t \mathbf{C}^{-1} \mathbf{m}_1 - \mathbf{m}_2^t \mathbf{C}^{-1} \mathbf{m}_2) - \log \left\{ \frac{P[\omega_1]}{P[\omega_2]} \right\}$$

- Also a distance to template classifier, where the distance is

$$(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_1)$$

Known as Mahalanobis distance

## Posterior probabilities for simple Gaussian cases

- Bayes classifier:

$$P[\omega_1|\mathbf{x}] = \frac{p(\mathbf{x}|\omega_1) P[\omega_1]}{p(\mathbf{x}|\omega_1) P[\omega_1] + p(\mathbf{x}|\omega_2) P[\omega_2]}$$

- Restrictive assumptions:

- Gaussian  $p(\mathbf{x}|\omega_j) = \mathcal{N}(\mathbf{m}_j, \mathbf{C}_j)$
- Equal covariance matrices:  $\mathbf{C}_1 = \mathbf{C}_2 = \mathbf{C}$

- Substitute, divide through by numerator term and cancel common terms to get

$$P[\omega_1|\mathbf{x}] = \frac{1}{1 + \exp\{-\mathbf{w}^t \mathbf{x} + w_0\}}$$

- The functional form  $1/(1 + \exp(-\alpha))$  is known as sigmoid / logistic.

## Summary:

Bayesian decision theory is fundamental to machine learning; we started from a probabilistic setting in which we described the posterior probability of membership and derived several results starting from this premise:

- Under certain assumptions...
  - The class boundary is linear
  - The posterior probability has a sigmoidal shaped default
  - The optimal classifier reduces to a distance to template classifier
- ...and under certain other assumptions...
  - The best classifier is still a distance to template classifier, but instead of Euclidean distance we need to use Mahalanobis distance.
- ...and under certain other assumptions...
  - The optimal classifier is a quadratic classifier

## Fisher Linear Discriminant Analysis

- In the  $p$ -dimensional ( $\mathcal{R}^p$ ) input space, find a direction on which projected data is maximally separable:
  - Projected means should be far apart
  - Projected scatter of each class should be small
- Projection of  $\mathbf{x}_n$  onto direction  $\mathbf{w}$  is  $\mathbf{w}^t \mathbf{x}_n$ ;
  - Projected mean for class  $j$  will be at  $\mathbf{w}^t \mathbf{m}_j$
  - Variance of projections is  $\mathbf{w}^t \mathbf{C}_j \mathbf{w}$ , where  $\mathbf{C}_j$  is the covariance matrix of data in class  $j$ .

- Fisher Ratio:

$$J_F = \frac{(\mathbf{w}^t \mathbf{m}_1 - \mathbf{w}^t \mathbf{m}_2)^2}{\mathbf{w}^t \mathbf{C}_1 \mathbf{w} + \mathbf{w}^t \mathbf{C}_2 \mathbf{w}}$$

- We can write the numerator as  $\mathbf{w}^t \mathbf{C}_B \mathbf{w}$ , where  $\mathbf{C}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t$ , the between-class scatter matrix.
- $\mathbf{C}_W = \mathbf{C}_1 + \mathbf{C}_2$ , the within class scatter matrix.

# Fisher Linear Discriminant Analysis

- Fisher criterion to maximize

$$J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{C}_B \mathbf{w}}{\mathbf{w}^t \mathbf{C}_W \mathbf{w}}$$

- Set gradient to zero

$$\nabla_{\mathbf{w}} = \frac{2\mathbf{C}_B \mathbf{w} \times (\mathbf{w}^t \mathbf{C}_W \mathbf{w}) - 2\mathbf{C}_W \mathbf{w} \times (\mathbf{w}^t \mathbf{C}_B \mathbf{w})}{(\mathbf{w}^t \mathbf{C}_W \mathbf{w})^2}$$

- Equate this to zero and observe
  - $\mathbf{w}^t \mathbf{C}_W \mathbf{w}$  and  $\mathbf{w}^t \mathbf{C}_B \mathbf{w}$  are scalars
  - $\mathbf{C}_B \mathbf{w}$  points in the same direction as  $\mathbf{m}_1 - \mathbf{m}_2$
  - We are only interested in the direction of  $\mathbf{w}$

$$\mathbf{w}_F = \alpha \mathbf{C}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

## Linear Regression & Perceptron

Formulation → derivation → solution

- Data:  $\{\mathbf{x}_n, f_n\}_{n=1}^N$   
Input:  $\mathbf{x}_n \in \mathcal{R}^p$ ; target / output  $f_n$  real valued
- Model:  $f = \mathbf{w}^t \mathbf{x} + w_0$   
Output linear function of input (including a constant  $w_0$ )
- Work in  $(p + 1)$  dimensional space to avoid treating  $w_0$  separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data:  $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model:  $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$  unknowns held in vector  $\mathbf{a}$

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left( \sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best  $\mathbf{a}$  we minimize  $E$  – differentiate with respect to each of the unknowns in  $\mathbf{a}$  and set to zero.

•

$$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left( \sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are  $(p + 1)$  derivatives (with respect to each  $a_i$ )
- Equating them to zero gives  $(p + 1)$  equations in  $(p + 1)$  unknowns

## Solution to Regression

- $(p + 1)$  simultaneous equations to solve:  
 $i^{\text{th}}$  row,  $j^{\text{th}}$  column shown

$$\begin{pmatrix} \dots & \dots & \dots \\ \dots & \sum_{n=1}^N y_{ni} y_{nj} & \dots \\ \vdots & \vdots & \vdots \\ \dots & \dots & \dots \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_{(p+1)} \end{pmatrix} = \begin{pmatrix} \vdots \\ \sum_{n=1}^N f_n y_{ni} \\ \vdots \end{pmatrix}$$



# Derivation in vector/matrix form

- $\mathbf{Y}$ :  $N \times (p + 1)$  matrix  $n^{\text{th}}$  row is  $\mathbf{y}_n^t$
- $\mathbf{f}$ :  $N \times 1$  vector of outputs
- Error  $E = \|\mathbf{Y} \mathbf{a} - \mathbf{f}\|^2$
- Gradient

$$\nabla_{\mathbf{a}} E = 2 \mathbf{Y}^t (\mathbf{Y} \mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned} \mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f} \end{aligned}$$

## Solution by Gradient Descent

- Gradient vector:  $\nabla_{\mathbf{a}} E = 2 \mathbf{Y}^t (\mathbf{Y} \mathbf{a} - \mathbf{f})$
- Steepest descent algorithm:

Initialize  $\mathbf{a}$  at random

Update  $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta \nabla_{\mathbf{a}} E$

Until Convergence

- Second order (Newton's) method

Initialize  $\mathbf{a}$  at random

Update  $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta \mathbf{H}^{-1} \nabla_{\mathbf{a}} E$

Until Convergence

- Rapid convergence with second order method, but cost of computing and inverting  $\mathbf{H}$  can be high (more on this under Neural Networks)

- Error  $E = \sum_{n=1}^N e_n^2$
- True gradient:

$$\nabla_{\mathbf{a}} E = 2 \sum_{n=1}^N \{ \mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n \} (\mathbf{y}_n)$$

- Gradient computed on  $n^{\text{th}}$  data:

$$\nabla_{\mathbf{a}} e_n = 2 \{ \mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n \} (\mathbf{y}_n)$$

## Regularization

- Pseudo inverse solution:  $\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}$
- This can be ill conditioned, so we could *regularize* by

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y} + \gamma \mathbf{I})^{-1} \mathbf{Y}^t \mathbf{f}$$

where  $\gamma$  is a small constant.

- We achieve precisely this by minimizing an error of the form

$$\| \mathbf{Y} \mathbf{a} - \mathbf{f} \|^2 + \gamma \| \mathbf{a} \|^2$$

Here a quadratic penalty term has been included

- Homework: Differentiate this error and derive the regularized solution
- Sparse solutions are obtained by regularizing with an  $l_1$  norm (sum of absolute values of  $\mathbf{a}$ , i.e.  $\sum_{j=1}^p |a_j|$ ): Known as *Lasso*

- Number of misclassified examples as measure of error  
Piecewise constant (cannot differentiate)
- Suitable error measure:

$$E_P = - \sum y_n^t \mathbf{a}$$

- Summation taken over misclassified examples
- We started with  $y_n^t \mathbf{a} > 0$  for positive class and  $y_n^t \mathbf{a} < 0$  for the negative class; we then switch the signs of negative class examples and required  $y_n^t \mathbf{a} > 0$  for all the training data; so for the misclassified examples  $-\sum y_n^t \mathbf{a}$  should be as small as possible.

## Perceptron

### Learning rule

- Gradient:

$$\frac{\partial E}{\partial \mathbf{a}} = - \sum y_n$$

- Gradient algorithm:  $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \sum y_n$
- Stochastic gradient algorithm:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + y_n$$

- Note what  $y_n$  is. It is an item of data that is taken at random and happens to be misclassified by the current value of  $\mathbf{a}$  at iteration  $k$ .

# Perceptron

## Convergence of the learning rule

- Learning Rule:  $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}(k)$   
where  $\mathbf{y}(k)$  is a misclassified input.
- Training criterion
  - We start with requiring  $\mathbf{a}^t \mathbf{y}(k) \leq 0$ , depending on the example belonging to class 1 or class 2.
  - If we switch the signs of examples of class 2, we require  $\mathbf{a}^t \mathbf{y}(k) > 0$  for all  $k$ .
- On misclassified data  $\mathbf{a}^t \mathbf{y}(k) < 0$
- If  $\hat{\mathbf{a}}$  is a solution (separable data), for all  $k$ ,  $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$
- We prove convergence by showing:  
 $\|\mathbf{a}^{(k+1)} - \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \hat{\mathbf{a}}\|^2$  for this update rule. i.e. the learning rule brings the guess closer to a valid solution.

# Perceptron

## Convergence of the learning rule (cont'd)

- For perceptron criterion, the magnitude of  $\mathbf{a}$  is not relevant (only the direction is). Hence for some scalar  $\alpha$ , we wish to show

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2$$

- From the update formula

$$\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}} = \mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}} + \mathbf{y}(k)$$

- Taking magnitudes

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 = \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 + 2(\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}})^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

- If we drop the negative term  $\mathbf{a}^{(k)t} \mathbf{y}(k)$  from RHS, the equality becomes an inequality

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha \hat{\mathbf{a}}^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

# Perceptron

## Convergence of the learning rule (cont'd)

- Of the three terms on the right hand side, we know  $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$ , because  $\hat{\mathbf{a}}$  is assumed to be a solution.

- If we select

$$\beta^2 = \max_i \|\mathbf{y}_i\|^2 \quad \gamma = \min_i \hat{\mathbf{a}}^t \mathbf{y}_i$$

i.e. largest of the positive term and smallest of the negative term, then for  $\alpha = \beta^2 / \gamma$ ,

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - \beta^2$$

- (Note the inequality remains true when the right hand side is replaced by a quantity larger than what it previously was.)
- Every correction takes the guess closer to a true solution.
- From an initialization  $\mathbf{a}^{(1)}$ , we will find a solution in *at most*  $k_0 = \frac{\|\mathbf{a}^{(1)} - \alpha \hat{\mathbf{a}}\|^2}{\beta^2}$  updates.

## Estimation

### Maximum likelihood on simple examples

- We have data  $\mathbf{x}_k$
- We have a model: e.g. the data came from a Gaussian density
- We have parameters relating to the model: e.g. mean of the Gaussian
- Our task is to estimate the parameters given the data
- Frequentist thought
  - The given data is a particular realization of the underlying system
  - Repeated experiments will give different estimates
  - If each experiment uses a lot of data, the variation may be small
  - We define a probabilistic model and maximize likelihood
  - Bias and Variance in estimation
- Bayesian thought
  - We are interested in the uncertainty in parameters
  - We have a prior uncertainty
  - There is some information in the data
  - We combine these to get a posterior uncertainty

# Likelihood & Log likelihood

- $p(\mathbf{x} | \omega_j)$
- Parametric form  $p(\mathbf{x} | \omega_j, \theta_j)$   
For example

$$p(\mathbf{x} | \omega_j, \theta_j) = \mathcal{N}(\mathbf{m}_j, \mathbf{C}_j)$$

- Dataset  $\mathcal{D} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  of identical and independently distributed samples (iid)
  - All samples were drawn from this distribution
  - Independent draws (previous value does not affect the next draw)
- Likelihood of an item of data (function of the parameter!)

$$p(\mathbf{x}_k | \theta)$$

- Likelihood of the set of data

$$p(\mathcal{D} | \theta) = \prod_{i=1}^n p(\mathbf{x}_i | \theta)$$

- Log likelihood

$$l(\theta) = \ln p(\mathcal{D} | \theta)$$

## Maximum Likelihood

- Maximum likelihood

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

- Maximize by taking derivative

$$\nabla_{\theta} = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_n} \end{bmatrix}$$

$$\nabla_{\theta} l = \sum_{k=1}^n \nabla_{\theta} p(\mathcal{D} | \theta)$$

... and equating to zero

$$\nabla_{\theta} l = \mathbf{0}.$$

... and solve for the unknown parameter values.

## Example: Multivariate Gaussian $\mathcal{N}(\mathbf{m}, C)$

Mean unknown, (Covariance known)

- ... product of Gaussians; taking log removes exp and turns  $\prod$  into  $\sum$
- ... write it out for a single data point

$$\ln p(\mathbf{x}_k | \mathbf{m}) = \frac{1}{2} \ln(2\pi)^d \det C - \frac{1}{2} (\mathbf{x}_k - \mathbf{m})^T C^{-1} (\mathbf{x}_k - \mathbf{m})$$

- ... the derivative

$$\nabla_{\mathbf{m}} \ln p(\mathbf{x}_k | \mathbf{m}) = C^{-1} (\mathbf{x}_k - \mathbf{m})$$

- Given  $n$  data  $\mathbf{x}_1, \dots, \mathbf{x}_n$

$$\sum_{k=1}^n C^{-1} (\mathbf{x}_k - \hat{\mathbf{m}}) = \mathbf{0}$$

- and the solution is...

$$\hat{\mathbf{m}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

## Example: Univariate Gaussian $\mathcal{N}(m, \sigma^2)$

(Mean and variance unknown)

- Two parameters:  $\theta_1 = m$  and  $\theta_2 = \sigma^2$
- Log likelihood of a single data:

$$\ln p(x_k | \boldsymbol{\theta}) = \frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

- Derivative of the log likelihood

$$\nabla_{\boldsymbol{\theta}} \ln p(x_k | \boldsymbol{\theta}) = \begin{bmatrix} \frac{1}{\theta_2} (x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{bmatrix}$$

- Given  $n$  data  $x_1, x_2, \dots, x_n$ , and considering the full log likelihood

$$\begin{aligned} \sum_{k=1}^n \frac{1}{\theta_2} (x_k - \hat{\theta}_1) &= 0 \\ -\sum_{k=1}^n \frac{1}{\theta_2} + \sum_{k=1}^n \frac{1}{\theta_2^2} (x_k - \hat{\theta}_1)^2 &= 0 \end{aligned}$$

# Example

## Univariate Gaussian, unknown mean and variance (cont'd)

- ... after some algebra

$$\begin{aligned}\hat{m} &= \frac{1}{n} \sum_{k=1}^n x_k \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{k=1}^n (x_k - \hat{m})^2\end{aligned}$$

- If we did this estimation from several datasets...

$$E \left[ \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2 \right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

expected value of estimate is not the same as the true value!

For the multivariate Gaussian, we have seen the results:

$$\begin{aligned}\hat{\mathbf{m}} &= \\ \hat{\mathbf{C}} &= \end{aligned}$$

# Bayesian Estimation

- Data:  $\mathcal{D} : x_1, \dots, x_n$
- Likelihood (as seen before)  $p(x|m) \sim \mathcal{N}(m, \sigma^2)$
- Prior uncertainty over parameters:  $p(m) \sim \mathcal{N}(m_0, \sigma_0^2)$   
 $m_0$  and  $\sigma_0^2$  are known.
- Posterior via Bayes' formula

$$p(m|\mathcal{D}) = \frac{p(\mathcal{D}|m)p(m)}{\int p(\mathcal{D}|m)p(m) dm}$$

Denominator is a constant, so we deal with

$$p(m|\mathcal{D}) = \alpha \prod_{k=1}^n p(x_k|m) p(m)$$

- Two ways forward from here
  - Maximum *a posteriori* estimation
  - Inference by integrating out parameters



# Bayesian Estimation: Univariate Gaussian

(Only the mean is unknown)

- Data:  $\mathcal{D} : x_1, \dots, x_n$
- Likelihood (as seen before)  $p(x|m) \sim \mathcal{N}(m, \sigma^2)$
- Prior uncertainty over parameters:  $p(m) \sim \mathcal{N}(m_0, \sigma_0^2)$   
 $m_0$  and  $\sigma_0^2$  are known.
- Substituting gives the posterior as a product of Gaussians

$$p(m|\mathcal{D}) = \alpha \prod_{k=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{x_k - m}{\sigma} \right)^2 \right] \times \frac{1}{\sigma_0 \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{m - m_0}{\sigma_0} \right)^2 \right]$$

- Which can be reduced to...

$$p(m|\mathcal{D}) = \alpha_2 \exp \left\{ -\frac{1}{2} \left\{ \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) m^2 - 2 \left( \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{m_0}{\sigma_0^2} \right) m \right\} \right\}$$

- But then...

$$p(m|\mathcal{D}) = \frac{1}{\sigma_n} \exp \left\{ -\frac{1}{2} \left( \frac{m - m_n}{\sigma_n} \right)^2 \right\}$$

- Matching terms...

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}$$

## Bayesian Estimation: Univariate Gaussian (cont'd)

- Finally...

$$m_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{m}_n + \left( \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right) m_0$$

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

- We now have an estimate that combines *prior* information about the parameter ( $p(m) = m_0$ ) with data ( $x_1, \dots, x_k$ ) to quantify uncertainty about the parameter:
  - Before seeing any data, we have a belief
  - As we see more and more data, our belief is taken over by what the data tells us.

# Unsupervised Learning

- Given:  $\{\mathbf{x}_n\}_{n=1}^N$  (as opposed to  $\{\mathbf{x}_n, f_n\}_{n=1}^N$ )
- We might extract cluster structures
  - Notion of distance between points of data
  - Criterion to determine how many clusters (often from prior knowledge)
  - Underlying probabilistic model
- We might project data onto a subspace  
 $\mathbf{x}_n \in \mathcal{R}^d \longrightarrow \mathbf{y}_n \in \mathcal{R}^q$ 
  - $q = 2$  helps visualization
  - Subspace representation useful for
    - Data compression
    - Sometimes used to reduce features

Semi Supervised Learning:

$$\{\mathbf{x}_n, f_n\}_{n=1}^{N_1} \text{ and } \{\mathbf{x}_n\}_{n=N_1+1}^{N_2}$$

## Principal Component Analysis:

- $N$  data  $\mathbf{x}_n \in \mathcal{R}^d$  distributed with mean  $\mathbf{m}$  and covariance matrix  $\mathbf{C}$ .
- Project onto direction  $\mathbf{u}$ ; find the direction that maximizes projected variance.
- Projected variance is  $\mathbf{u}^t \mathbf{C} \mathbf{u}$
- We are only interested in the direction; not in increasing the projected variance by choosing  $\mathbf{u}$  with large magnitude.
- Set up a constrained optimization problem

$$\max_{\mathbf{u}} \mathbf{u}^t \mathbf{C} \mathbf{u} \quad \text{subject to } \mathbf{u}^t \mathbf{u} = 1$$

- Lagrangian

$$\mathcal{L} = \mathbf{u}^t \mathbf{C} \mathbf{u} - \lambda [\mathbf{u}^t \mathbf{u} - 1]$$

- $\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = 0 \implies \mathbf{C} \mathbf{u} = \lambda \mathbf{u}$ ; i.e. principal directions are eigenvectors of covariance

# Mixture Model

We write a mixture of Gaussian densities:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- If the mixing proportions  $\pi_k$  satisfy
  - $\pi_k \geq 0$
  - $\sum_{k=1}^K \pi_k = 1$ $p(\mathbf{x})$  is a proper probability density.
- More powerful model – useful when data is multi-modal
- Parameters are: proportions, means and covariance matrices
- Parameter estimation ( $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ ) is not easy.
- $z_{nk}$ : association of  $n^{\text{th}}$  data to  $k^{\text{th}}$  mode unknown (latent)

## Algorithm:

$$\begin{aligned}\pi_k &= \frac{1}{N} \sum_{n=1}^N q_{nk} \\ \boldsymbol{\mu}_k &= \frac{\sum_{n=1}^N q_{nk} \mathbf{x}_n}{\sum_{n=1}^N q_{nk}} \\ \boldsymbol{\Sigma}_k &= \frac{\sum_{n=1}^N q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^t}{\sum_{n=1}^N q_{nk}}\end{aligned}$$

$$q_{nk} = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

# Expectation Maximization

## Auxilliary Variable as Posteriors

Interpret:

- Mixture model as a Gaussian classifier with  $K$  classes
- $\pi_k$  as prior probabilities
- Each of the  $\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  as class conditional densities / likelihoods.

$$\begin{aligned} p(z_{nk} = 1 | \mathbf{x}_n, \boldsymbol{\pi}, \boldsymbol{\Delta}) &= \frac{p(z_{nk} = 1 | \pi_k) p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K p(z_{nj} = 1 | \pi_j) p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \\ &= q_{nk} \end{aligned}$$

- Each data item has a weighted contribution to the estimation of parameters.
- Unknown assignment  $z_{nk}$ ; **E**xpected value of this unknown assignment is  $q_{nk}$ , the posterior probability
- **M**aximize (the lower bound) to re-estimate parameters.

## K-Means Clustering Algorithm

**Input:**  $\mathbf{X} = \{\mathbf{x}_n^t\}_{n=1}^N, K$

**Output:**  $\mathbf{C}, \text{Idx}$

initialize:  $\mathbf{C} = \{\mathbf{c}_j^t\}_{j=1}^K$

repeat

. assign  $n^{\text{th}}$  sample to nearest  $\mathbf{c}_j$   
.  $\text{Idx}(n) = \min_j \|\mathbf{x}_n - \mathbf{c}_j\|^2$

. recompute  $\mathbf{c}_j = \frac{1}{N_j} \sum_{n=j} \mathbf{x}_n$

until no change in  $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma_k)$$

- Set  $\Sigma_k = \sigma_k^2 I$
- At every iteration, set largest  $q_{nk}$  (largest over  $k$ ) to one and others to zero. Winner take all at each datapoint.
- Computation of  $q_{nk}$  is expectation of latent variable  $z_{nk}$  – **E** step
- Re-estimation of  $\mu_k$  and  $\Sigma_k$  become maximum likelihood estimates from data assigned to each cluster (because  $q_{nk}$  is either one or zero) – **M** step

## Radial Basis Functions

- We fix nonlinear part in a data-dependent way
- Estimate only a linear part (by linear regression)

The Model:

$$g(\mathbf{x}) = \sum_{j=1}^M \lambda_j \phi(|\mathbf{x} - \mathbf{m}_j| / \sigma)$$

- We think of *centers*  $\mathbf{m}_j$  in the input space
- ‘Radial’ comes from distances to these centers (scaled by  $\sigma$ )
- We refer to  $\phi(\cdot)$  as basis functions
- Note similarity to Fourier transform – signal expressed as sum of sine waves (basis functions)
- $\mathbf{m}_j, j = 1, \dots, M, \sigma$  and the functional form of  $\phi(\cdot)$  are fixed by some sensible (*ad hoc*) way
- $\lambda_j, j = 1, \dots, M$  are the unknowns

## Radial Basis Functions (cont'd)

- Data is  $\{\mathbf{x}_n, f_n\}_{n=1}^N$
- Similar to linear regression, we will hope to achieve

$$\begin{aligned} f_n &= g(\mathbf{x}_n) \\ &= \sum_{j=1}^M \lambda_j \phi(|\mathbf{x}_n - \mathbf{m}_j| / \sigma) \end{aligned}$$

- *i.e.* at each input  $\mathbf{x}_n$ , the function should output our target  $f_n$
- Usually,  $M < N$  *i.e.* we have fewer basis functions than there are items of data

## Radial basis functions (cont'd)

What basis functions?

- Linear:  $\phi(\alpha) = \alpha$
- Gaussian:  $\phi(\alpha) = \exp(-\alpha^2/\sigma^2)$
- Multi-quadric:  $\phi(\alpha) = \sqrt{1 + \alpha^2}$
- Thin plate splines:  $\phi(\alpha) = \alpha^2 \log(\alpha)$
- Gaussian very popular in practice
  - Local basis functions
  - $\mathbf{m}_j$  obtained by clustering

## Radial basis functions (cont'd): Estimating $\lambda_j$

- $N \times M$  matrix  $Y$

$$\begin{bmatrix} \phi(|\mathbf{x}_1 - \mathbf{m}_1|/\sigma) & \phi(|\mathbf{x}_1 - \mathbf{m}_2|/\sigma) & \dots & \phi(|\mathbf{x}_1 - \mathbf{m}_M|/\sigma) \\ \phi(|\mathbf{x}_2 - \mathbf{m}_1|/\sigma) & \phi(|\mathbf{x}_2 - \mathbf{m}_2|/\sigma) & \dots & \phi(|\mathbf{x}_2 - \mathbf{m}_M|/\sigma) \\ \vdots & \vdots & \dots & \vdots \\ \phi(|\mathbf{x}_N - \mathbf{m}_1|/\sigma) & \phi(|\mathbf{x}_N - \mathbf{m}_2|/\sigma) & \dots & \phi(|\mathbf{x}_N - \mathbf{m}_M|/\sigma) \end{bmatrix}$$

- Targets  $N \times 1$  vector

$$\begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_N \end{bmatrix}$$

- Minimize error

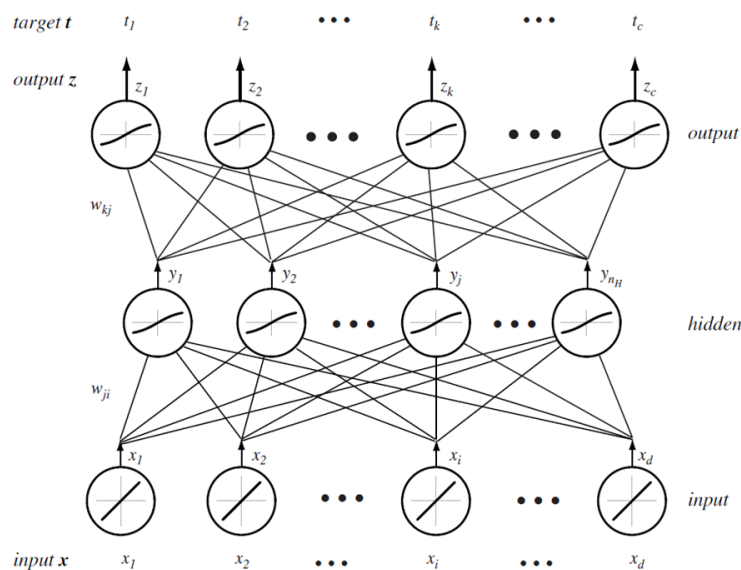
$$\hat{\lambda} = \min_{\lambda} \|Y\lambda - \mathbf{f}\|^2$$

- Solution similar to linear regression; e.g. pseudo inverse

$$\hat{\lambda} = (Y^T Y)^{-1} Y^T \mathbf{f}$$

- Clustering and other rules to set  $\mathbf{m}_j$  and  $\sigma$

## Multi-layer perceptron



$$g_k(\mathbf{x}) = f \left( \sum_{j=1}^{n_H} w_{jk} f \left( \sum_{i=1}^d w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$

# Each unit is a logistic

- Weighted sum of inputs – net activation

$$\text{net}_j = \sum_{i=1}^d x_i w_{ji} + w_{j0}$$

- Squashed by a nonlinearity

$$y_j = f(\text{net}_j)$$

- Simple threshold

$$f(\text{net}) = \begin{cases} 1 & \text{if } \text{net} \geq 0 \\ -1 & \text{if } \text{net} < 0 \end{cases}$$

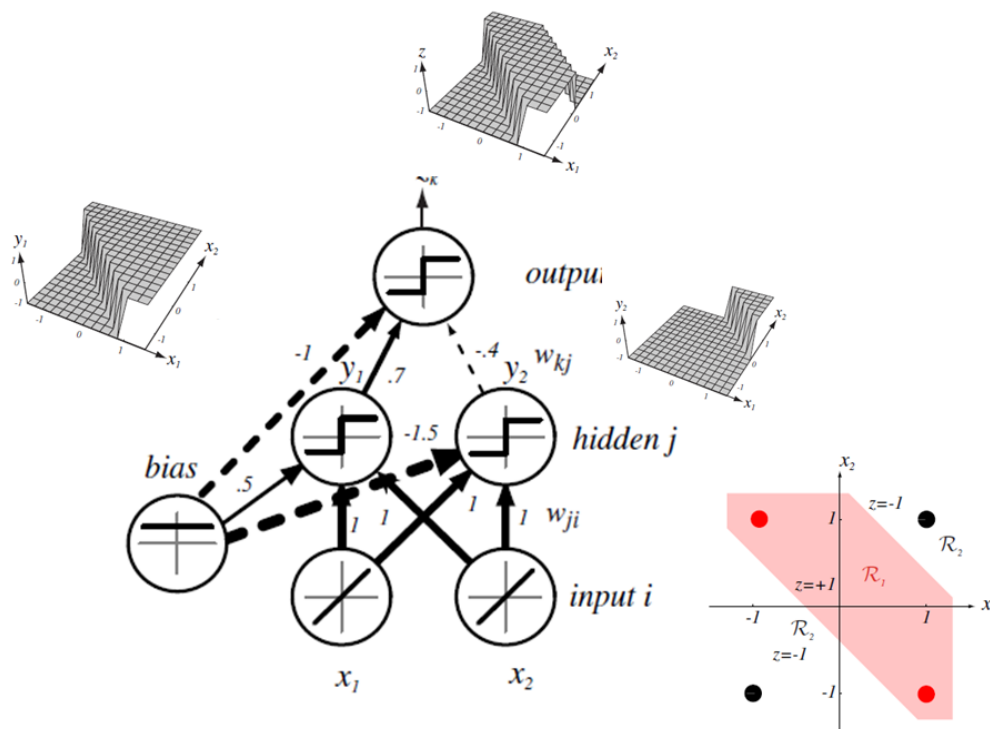
- Sigmoid (logistic)

$$f(\text{net}) = \frac{1}{1 + \exp(-\text{net})}$$

- Hyperbolic tangent

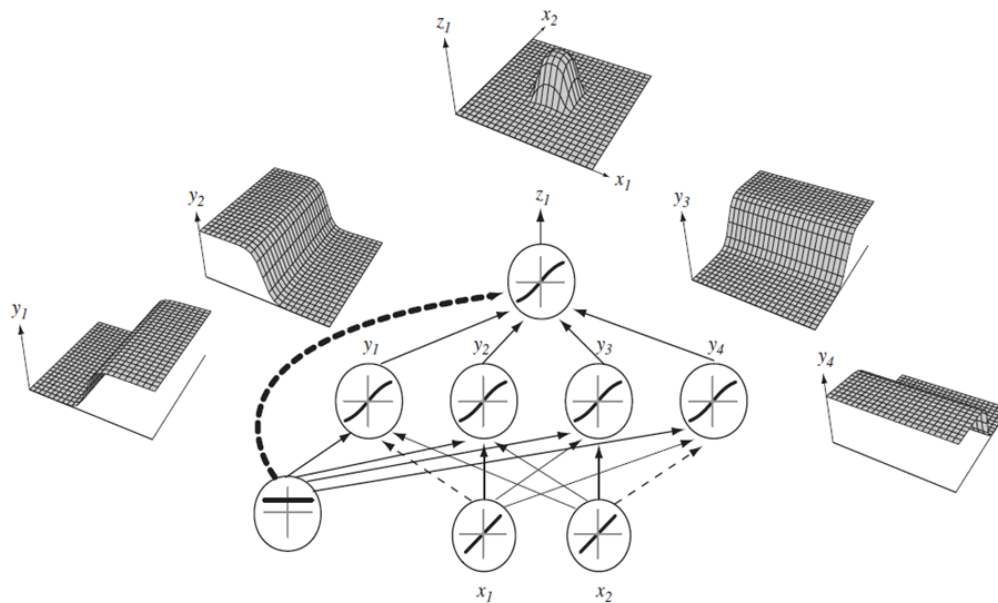
$$\begin{aligned} f(\text{net}) &= \tanh(b \text{ net}) \\ &= \frac{\exp(b \text{ net}) - \exp(-b \text{ net})}{\exp(b \text{ net}) + \exp(-b \text{ net})} \end{aligned}$$

## XOR Using a MLP

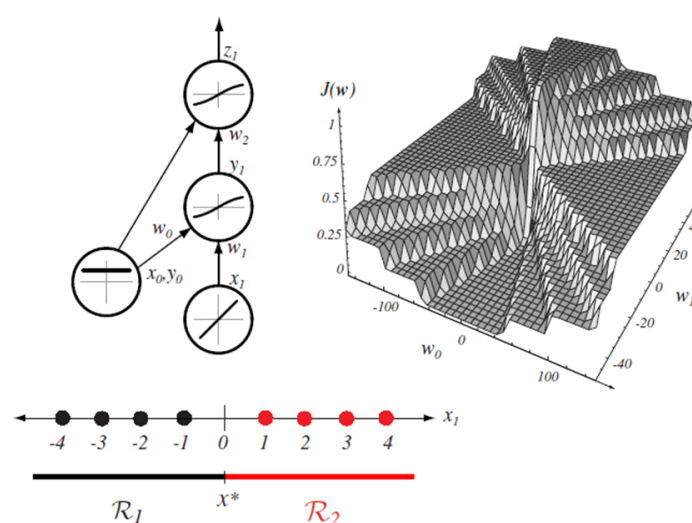




# MLP Constructing Complex Functions



## Error Function for MLP



- Not smooth and quadratic like for the linear model
- Large variation in gradient
- Local minima
- Gradient descent training with clever tricks

# Training Multi-Layer Neural Networks

- Error

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 \\ &= \frac{1}{2} \|\mathbf{t} - \mathbf{z}\|^2 \end{aligned}$$

- Change in weight for gradient descent

$$\Delta \mathbf{w} = -\eta \frac{\partial J}{\partial \mathbf{w}}$$

Or in component form

$$\Delta w_{pq} = -\eta \frac{\partial J}{\partial w_{pq}}$$

- Gradient descent update

$$\mathbf{w}(m+1) = \mathbf{w}(m) + \Delta \mathbf{w}(m)$$

## Error Back Propagation Algorithm

- Hidden to output layer weights

$$\begin{aligned} \frac{\partial J}{\partial w_{kj}} &= \frac{\partial J}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial w_{kj}} \\ &= -\delta_k \frac{\partial \text{net}_k}{\partial w_{kj}} \end{aligned}$$

where, change in error with respect to the activation of the unit,

$$\delta_k = -\frac{\partial J}{\partial \text{net}_k}$$

- Easy form for  $\delta_k$

$$\delta_k = -\frac{\partial J}{\partial \text{net}_k} = -\frac{\partial J}{\partial z_k} \frac{\partial z_k}{\partial \text{net}_k} = (t_k - z_k) f'(\text{net}_k)$$

Error at output  $\times$  slope of nonlinearity

- Also with respect to weights net is differentiated easily

$$\frac{\partial \text{net}_k}{\partial w_{kj}} = y_j$$

## Error back propagation (cont'd)

- Units internal to the network have no explicit error signal
- Chain rule of differentiation helps!

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \frac{\partial y_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

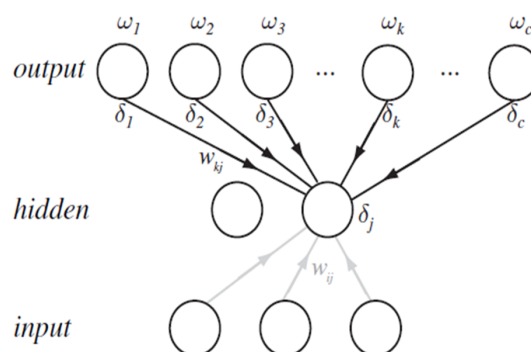
$$\begin{aligned} \frac{\partial J}{\partial y_j} &= \frac{\partial}{\partial y_j} \left[ \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 \right] \\ &= - \sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial y_j} \\ &= - \sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial y_j} \\ &= - \sum_{k=1}^c (t_k - z_k) f'(\text{net}_k) w_{kj} \end{aligned}$$

## Error propagation (cont'd)

$$\delta_j = f'(\text{net}_j) \sum_{k=1}^c w_{kj} \delta_k$$

and the update rule

$$\delta w_{ji} = \eta x_i \delta_j = \eta \left[ \sum_{k=1}^c w_{kj} \delta_k \right] f'(\text{net}_j) x_i$$



- Signal  $y_j$  propagates forward
- $\delta_j$ 's propagate backwards

- Having gradient enables us to update weights
- Can sum over all data and update with true gradient
- Or use sample-by-sample stochastic update
- Gradient descent can be slow (large flat regions)
- Can get stuck in local minima

## Speeding-up Training: Newton's Method

Taylor expansion

$$J(\mathbf{w} + \Delta \mathbf{w}) = J(\mathbf{w}) + \left( \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right)^T \Delta \mathbf{w} + \frac{1}{2} \Delta \mathbf{w}^T H \Delta \mathbf{w} + \dots$$

- We can consider the change in objective function

$$\Delta J(\mathbf{w}) = J(\mathbf{w} + \Delta \mathbf{w}) - J(\mathbf{w})$$

and ask for what  $\Delta \mathbf{w}$  is this minimized?

$$\left( \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right) + H \Delta \mathbf{w} = 0$$

$$\Delta \mathbf{w} = -H^{-1} \left( \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right)$$

... giving us the update

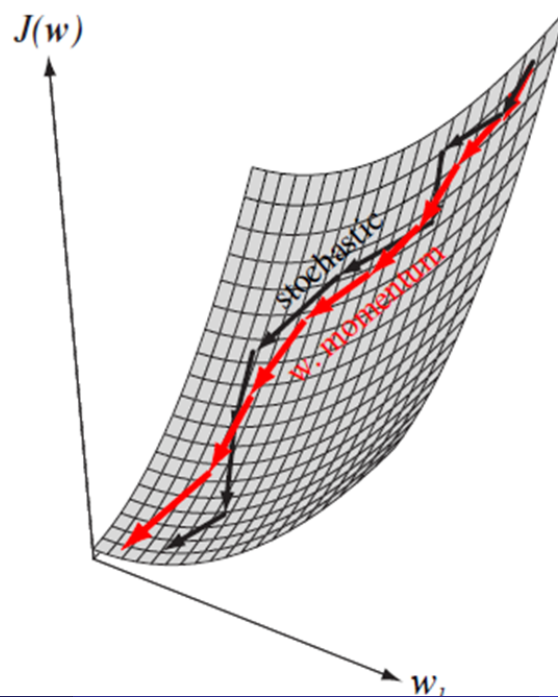
$$\begin{aligned} \mathbf{w}(m+1) &= \mathbf{w}(m) + \Delta \mathbf{w} \\ &= \mathbf{w}(m) - H^{-1}(m) \left( \frac{\partial J(\mathbf{w}(m))}{\partial \mathbf{w}} \right) \end{aligned}$$

- However
  - $N$  weights  $\rightarrow N^2$  storage for Hessian
  - Matrix inversion  $\mathcal{O}(N^3)$  complexity

# Momentum

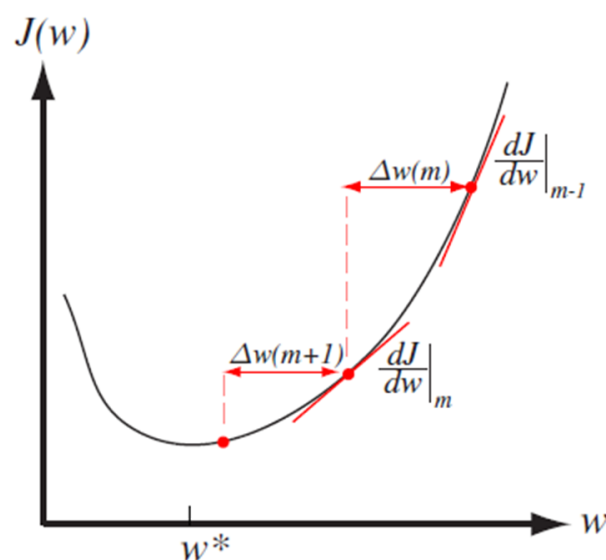
- Slow learning in regions in which gradient is small

$$\mathbf{w}(m+1) = \mathbf{w}(m) + (1 - \alpha) \Delta \mathbf{w}_{\text{BP}}(m) + \alpha \Delta \mathbf{w}(m-1)$$



## QuickProp

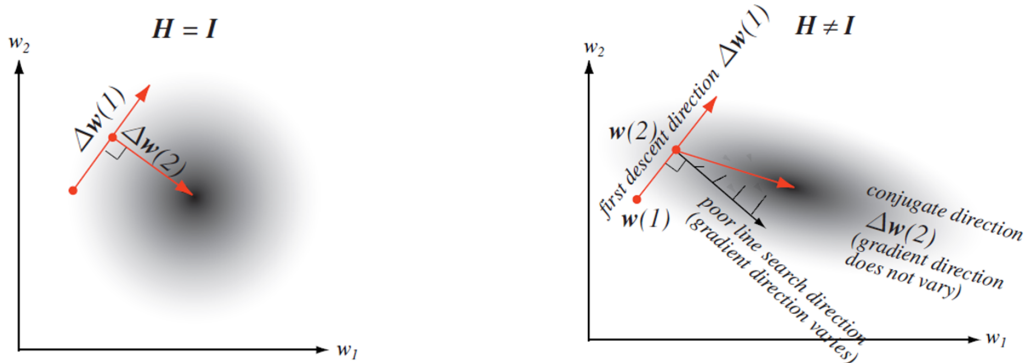
Two successive evaluations to approximate local curvature



$$\Delta w(m+1) = \frac{\partial J / \partial w|_m}{\partial J / \partial w|_{m-1} - \partial J / \partial w|_m} \Delta w(m)$$

# Conjugate Gradients

- Sequence of line searches (not just gradient update)



## Conjugate Gradients (cont'd)

- One dimensional line searches in multiple dimensions
- Careful choice of successive directions to search, informed by local curvature
- $\Delta \mathbf{w}(m-1)$  direction of line search at step  $m-1$
- For new direction to search, find  $\Delta \mathbf{w}$

$$\Delta \mathbf{w}^T(m-1) H \Delta \mathbf{w} = 0$$

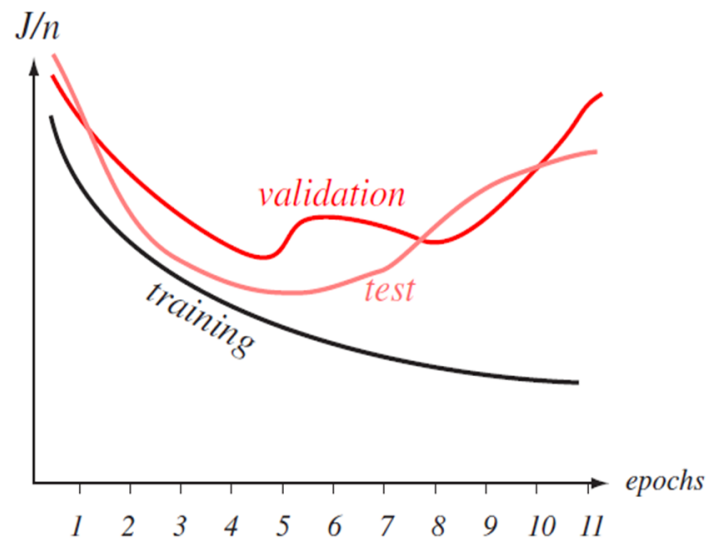
- Can be expressed as the gradient plus a component of previous search direction

$$\Delta \mathbf{w}(m) = \frac{\partial J(\mathbf{w}(m))}{\partial \mathbf{w}} + \beta_m \Delta \mathbf{w}(m-1)$$

- Need clever (approximate) way to set  $\beta_m$ : Fletcher-Reeves method

$$\beta_m = \frac{\nabla J^T(\mathbf{w}(m)) \nabla J(\mathbf{w}(m))}{\nabla J^T(\mathbf{w}(m-1)) \nabla J(\mathbf{w}(m-1))}$$

## Performance on Out-of-Sample data



- Validation set to monitor error
- Stop (early) when validation error begins to increase

## MLP As Posterior Probability Estimator

- Recall Bayes' decision

$$P[\omega_k|\mathbf{x}] = \frac{p(\mathbf{x}|\omega_k) p[\omega_k]}{\sum_{i=1}^c p(\mathbf{x}|\omega_i) p[\omega_i]}$$

- Suppose we train a neural net with  $c$  outputs

$$t_k(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \omega_k \\ 0 & \text{otherwise} \end{cases}$$

- Minimising the error

$$\sum_{\mathbf{x}} [g_k(\mathbf{x}, \mathbf{w}) - t_k]^2$$

when we have infinite data, can be shown to be the same as minimizing

$$\sum_{k=1}^c \int [g_k(\mathbf{x}, \mathbf{w}) - P(\omega_k|\mathbf{x})] p(\mathbf{x}) d\mathbf{x}$$