



Artificial Intelligence

Linear Models for Classification

Donghui Wang
AI Institute@ZJU
2015.4



Contents

- Basic Concepts
- Discriminant Functions
- Probabilistic Generative Models
- Probabilistic Discriminative Models
- The Laplace Approximation
- Bayesian Logistic Regression

References:

1. Bishop. *“Pattern Recognition and Machine Learning”*, Chapter 4. 2006.



浙江大学

ZheJiang University



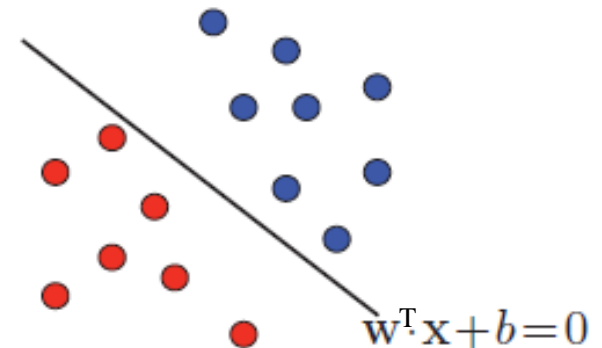
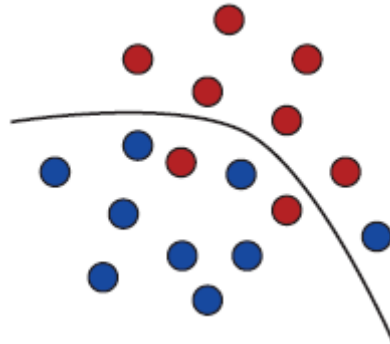
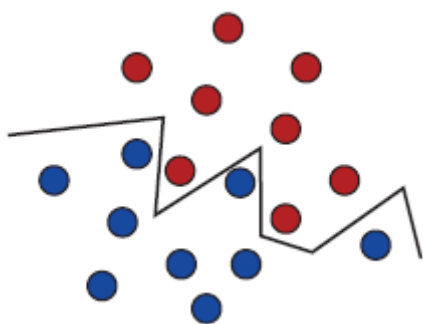
人工智能研究所

Institute of Artificial Intelligence

Basic Concepts

Linearly separable

- Decision regions:
 - Input space is divided into several regions
- Decision boundaries (surfaces):
 - Under linear models, it's a linear function of the input vector x
 - $(D-1)$ -dimensional hyper-plane within the D -dimensional input space
- Data sets whose classes can be separated exactly by linear decision surfaces are said to be *linearly separable*.





Representation of Class Labels

- Two classes ($K=2$):
 - Target variable $t \in \{0,1\}$, $t=1$ represents class C_1 , else class C_2
- K-classes ($K>2$):
 - 1-of-K coding scheme: $\mathbf{t} = (0, 1, 0, 0, 0)^T$
- Predict discrete class labels:
 - Linear model prediction (linear discriminant function): $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - Nonlinear function $f(\cdot) : \mathbb{R} \rightarrow (0, 1)$
 - Generalized linear models:

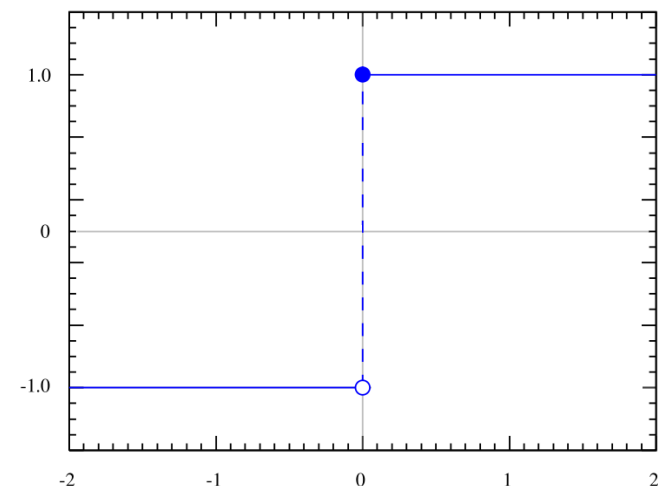
\mathbf{W} : weight vector

w_0 : bias/threshold

$f(\cdot)$: activation function
 $f^{-1}(\cdot)$: link function

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$

- Decision surface:
 - $y(\mathbf{x}) = \text{constant} \rightarrow \mathbf{w}^T \mathbf{x} + w_0 = \text{constant}$





Three classification approaches

- Discriminant function:
 - *Least-squares approach*: making the model predictions as close as possible to a set of target values
 - *Fisher's linear discriminant*: maximum class separation in the output space
 - *The perceptron algorithm of Rosenblatt*: generalized linear model
- Generative approach:
 - Model the class-conditional densities and the class priors
 - Compute posterior probabilities through Bayes's theorem
- Discriminative approach:
 - Directly training posterior probabilities.



浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Discriminant Functions (**nonprobabilistic methods**)

Two classes

- Linear discriminant function: $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - if $y(\mathbf{x}) \geq 0$, assign \mathbf{x} to class C_1 , else class C_2
 - decision surface Ω : $y(\mathbf{x}) = 0$
 - the normal distance from the origin to the decision surface: $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$

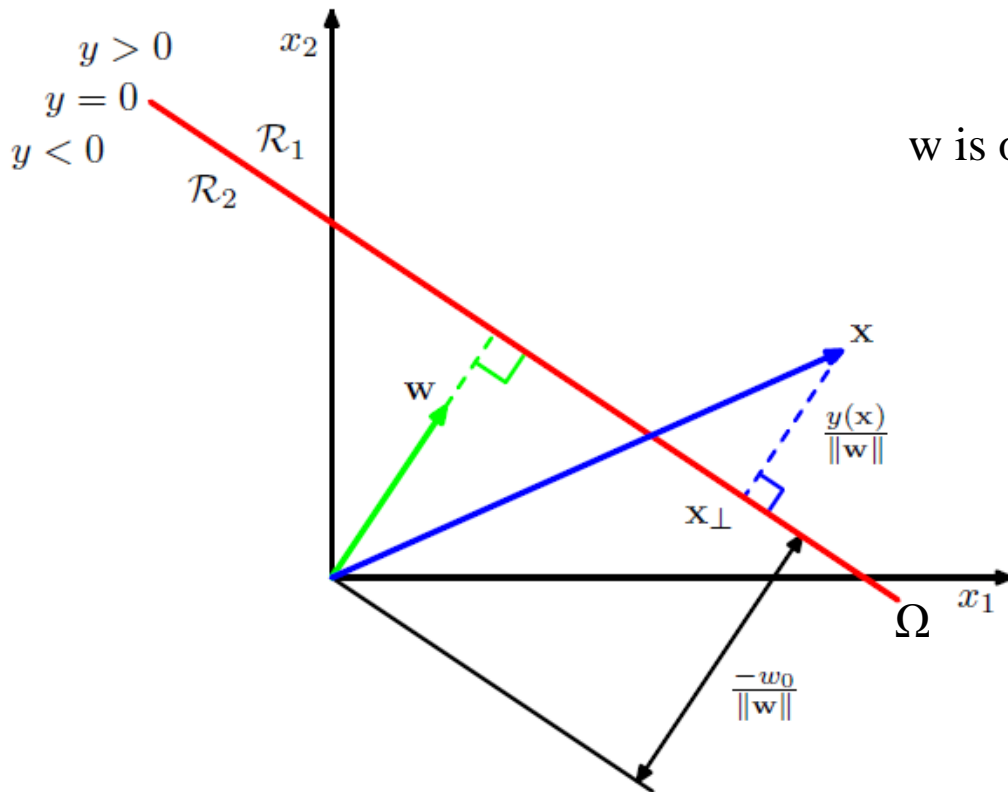
\mathbf{x}_A and \mathbf{x}_B lie on the decision surface: $y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$

$$\mathbf{w}^T (\mathbf{x}_A - \mathbf{x}_B) = 0$$

\mathbf{w} is orthogonal to every vector lying within Ω

$\frac{\mathbf{w}}{\|\mathbf{w}\|}$ is the normal vector of Ω

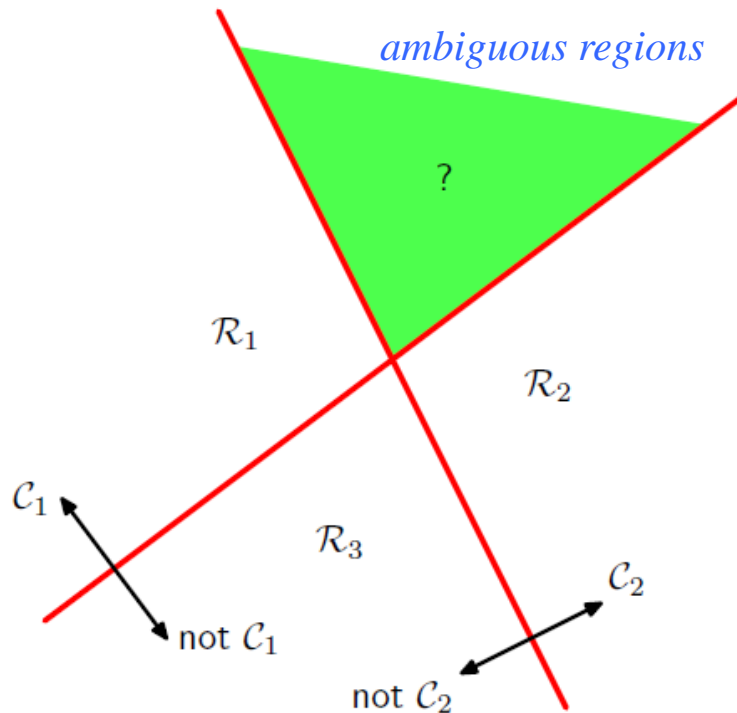
$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad \Rightarrow \quad r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$



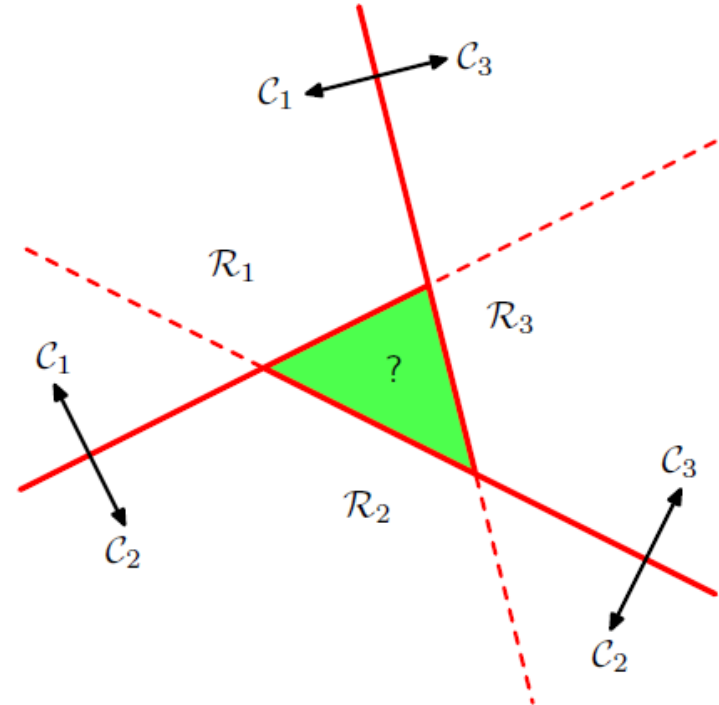
$$\begin{aligned} \tilde{\mathbf{w}} &= (w_0, \mathbf{w}) & \tilde{\mathbf{x}} &= (x_0, \mathbf{x}) \\ y(\mathbf{x}) &= \tilde{\mathbf{w}}^T \tilde{\mathbf{x}} \end{aligned}$$

Multiple classes

- How to build a K-class discriminant function?
 - *One-versus-the-rest classifier*
 - K-1 classifiers each of which solves a two-class problem
 - *One-versus-one classifier*
 - $K(K - 1)/2$ binary discriminant functions



One-versus-the-rest classifier



One-versus-one classifier

Multiple classes

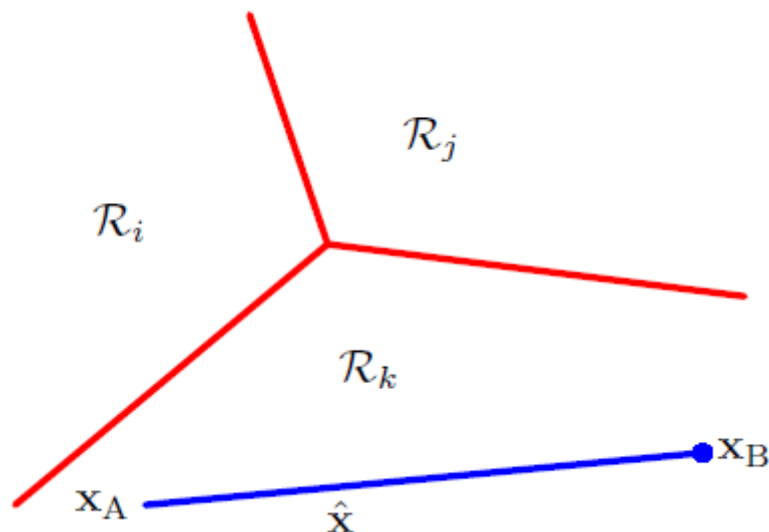
- Single K-class discriminant comprising K linear functions:

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

- assigning a point \mathbf{x} to class \mathcal{C}_k if $y_k(\mathbf{x}) > y_j(\mathbf{x})$ for all $j \neq k$.
- decision boundary between class \mathcal{C}_k and class \mathcal{C}_j is given by $y_k(\mathbf{x}) = y_j(\mathbf{x})$

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

- \mathcal{R}_k is singly connected and convex.



$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B \quad \text{where } 0 \leq \lambda \leq 1$$

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B)$$

Because both \mathbf{x}_A and \mathbf{x}_B lie inside \mathcal{R}_k , it follows that $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A)$, $y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$, for all $j \neq k$, and hence $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$, and so $\hat{\mathbf{x}}$ also lies inside \mathcal{R}_k .

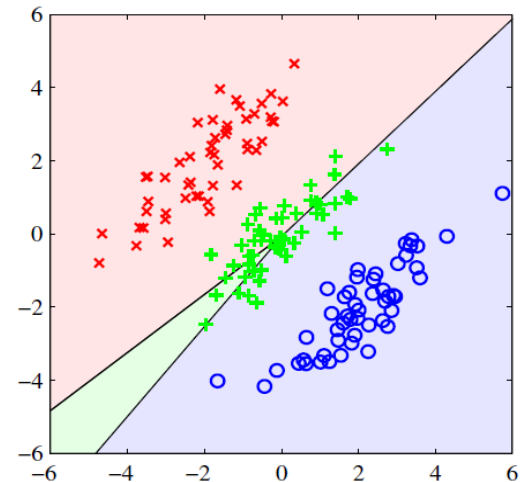
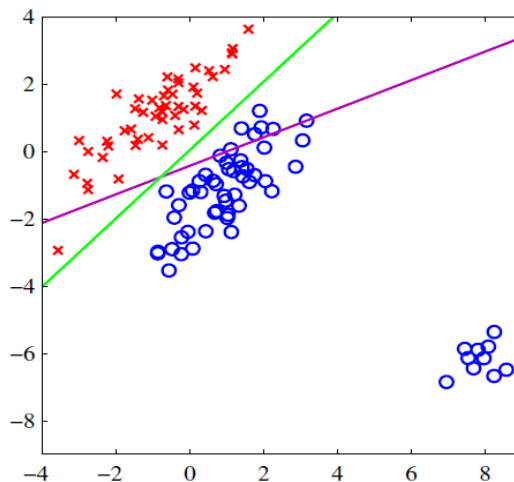
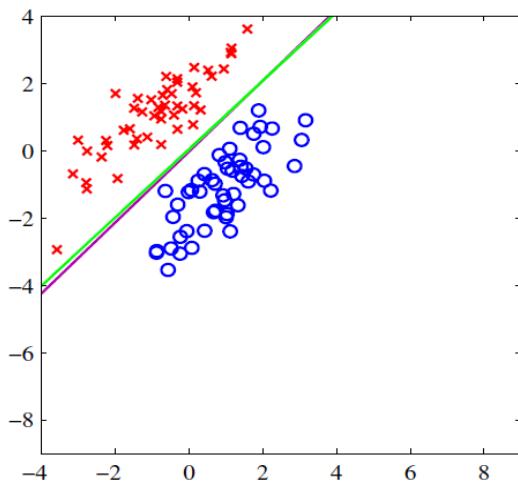


Learning the parameters of LDF

- Three approaches:
 - *Least-squares approach*:
 - making the model predictions as close as possible to a set of target values
 - *Fisher's linear discriminant*:
 - maximum class separation in the output space
 - *The perceptron algorithm of Rosenblatt*:
 - generalized linear model

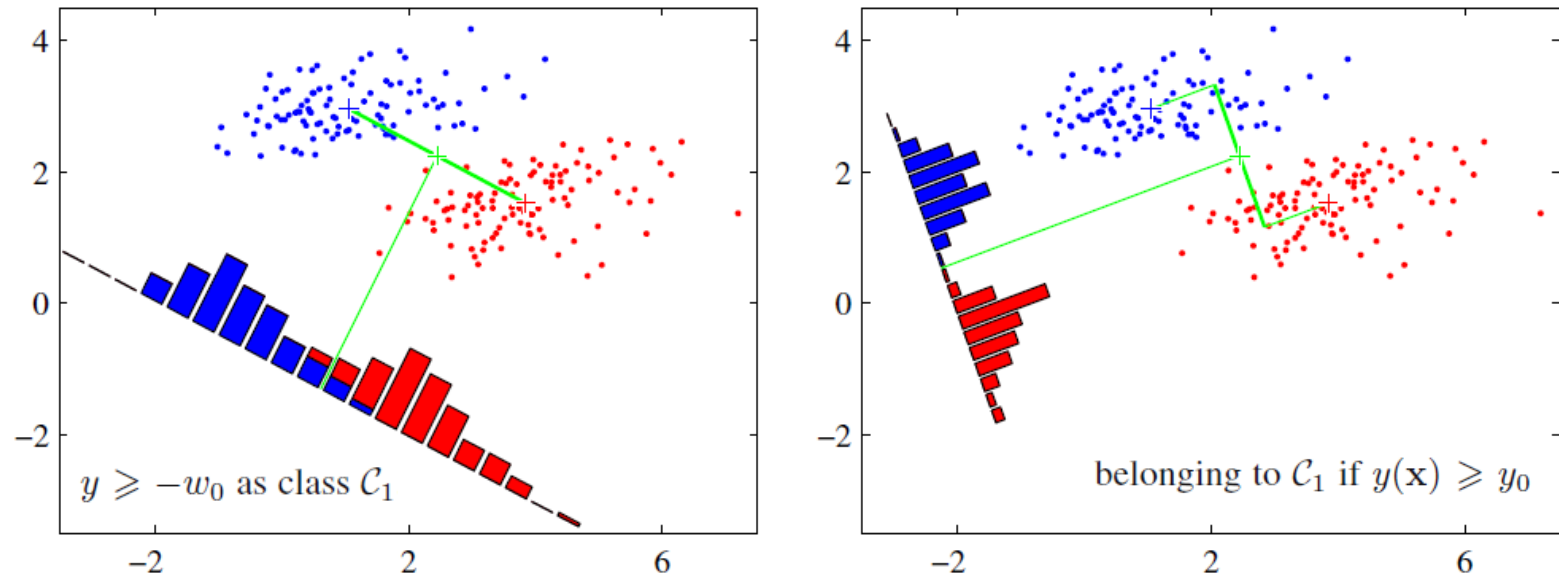
Least squares for classification

- Problem:
 - Each class \mathcal{C}_k is described by its own linear model:
$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad \text{where } k = 1, \dots, K$$
 - group together:
$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} \quad \widetilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T \quad \widetilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$$
 - new input \mathbf{x} is then assigned to the class for which the output $y_k = \widetilde{\mathbf{w}}_k^T \widetilde{\mathbf{x}}$ is largest.
- Learning $\widetilde{\mathbf{W}}$ with training data set: $\{\mathbf{x}_n, \mathbf{t}_n\}$ where $n = 1, \dots, N$
 - SSE function:
$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T}) \right\} \quad \Rightarrow \quad \widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T}$$
- Discriminant function:
$$y(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} = \mathbf{T}^T (\widetilde{\mathbf{X}}^\dagger)^T \widetilde{\mathbf{x}}$$



Fisher's linear discriminant

- From the view of dimensionality reduction:



- The simplest measure of the separation of the classes is the separation of the projected class means:

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n \quad \xrightarrow[\mathbf{m}_k = \mathbf{w}^T \mathbf{m}_k]{y = \mathbf{w}^T \mathbf{x}} \quad m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$

- Problem:** we can increase the magnitude of \mathbf{w} to make $(m_2 - m_1)$ arbitrarily large!

$$\sum_i w_i^2 = 1 \quad \longrightarrow \quad \mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$$

Fisher's linear discriminant

- *The Fisher's criterion*: maximize the separation between the projected class means as well as the inverse of the total within-class variance.

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2 \quad y = \mathbf{w}^T \mathbf{x} \quad m_k = \mathbf{w}^T \mathbf{m}_k$$

➡ $J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$ *Generalized Rayleigh quotient*

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad \text{Between-class covariance matrix}$$

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \quad \text{Within-class covariance matrix}$$

-
- *Fisher's linear discriminant*:

$$\nabla J(\mathbf{w}) = 0 \quad \Rightarrow \quad (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w} \quad \Rightarrow \quad \boxed{\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)}$$

Relation to least squares

- The Fisher criterion can be obtained as a special case of least squares if we consider following target coding scheme:
 - The target for class C_1 to be N/N_1 , for class C_2 to be $-N/N_2$
 - The sum-of-squares error function:

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2$$

$$\frac{\partial E}{\partial w_0} = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0$$

$$\sum_{n=1}^N t_n = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} = 0$$

$$\Rightarrow w_0 = -\mathbf{w}^T \mathbf{m}$$

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)$$

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0$$

$$\Rightarrow \left(\mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2)$$

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$\Rightarrow \mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Fisher's discriminant for multiple classes

- Assume input space dimensionality $D > K$ (number of classes, $K > 2$):

$$y = \mathbf{W}^T \mathbf{x} \quad y_k = \mathbf{w}_k^T \mathbf{x}$$

covariance matrices defined in the original x-space

- Total covariance matrix: $\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T \quad \mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k \quad N = \sum_k N_k$$

- The generalization of the within-class covariance matrix:

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad \mathbf{S}_k = \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \quad \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

- The generalization of the between-class covariance matrix:

➡
$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

Fisher's discriminant for multiple classes

- Assume input space dimensionality $D > K$ (number of classes, $K > 2$):

$$y = W^T x \quad y_k = w_k^T x$$

covariance matrices defined in the projected y-space

- The generalization of the within-class and between-class covariance matrix:

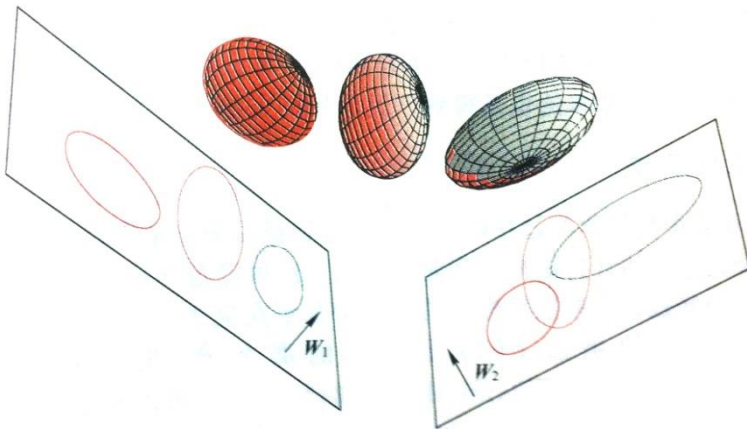
$$S_W = \sum_{k=1}^K \sum_{n \in C_k} (y_n - \mu_k)(y_n - \mu_k)^T$$

$$S_B = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T$$

$$\mu_k = \frac{1}{N_k} \sum_{n \in C_k} y_n$$

$$\mu = \frac{1}{N} \sum_{k=1}^K N_k \mu_k$$

- The Fisher's criterion:* $J(W) = \text{Tr} \{ S_W^{-1} S_B \} = \text{Tr} \{ (W^T S_W W)^{-1} (W^T S_B W) \}$



$$S_W = \sum_{k=1}^K \sum_{n \in C_k} (x_n - m_k)(x_n - m_k)^T$$

$$S_B = \sum_{k=1}^K N_k (m_k - m)(m_k - m)^T$$

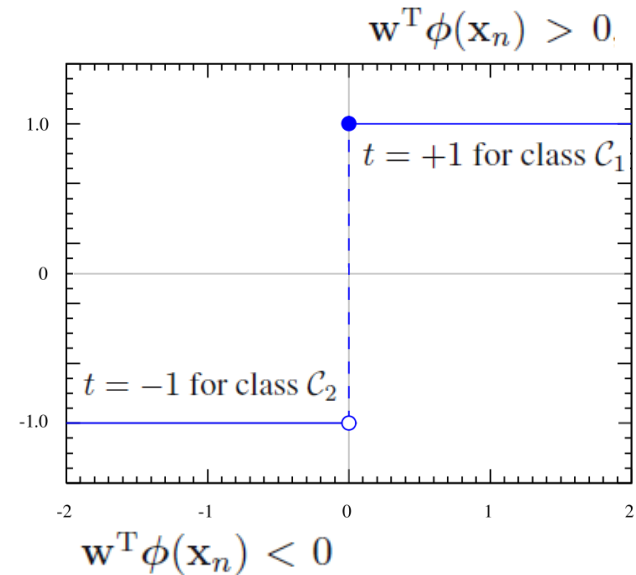
The perceptron algorithm

- Construct a generalized linear model:

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x})) \quad f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

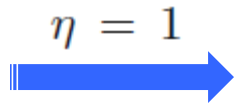
- Perceptron criterion (need to be minimized):

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n$$



- Stochastic gradient descent algorithm:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$



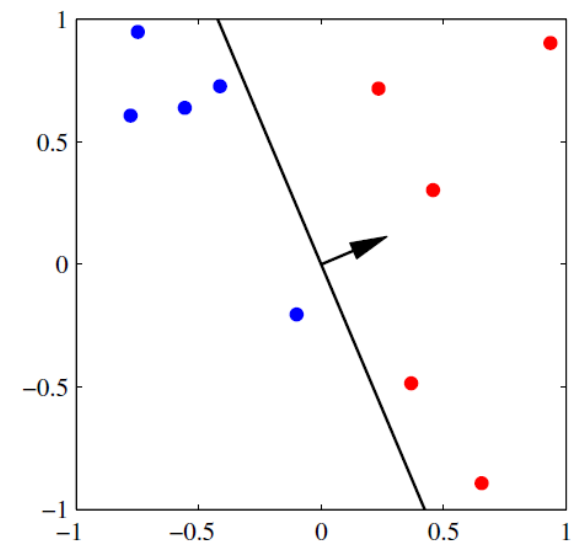
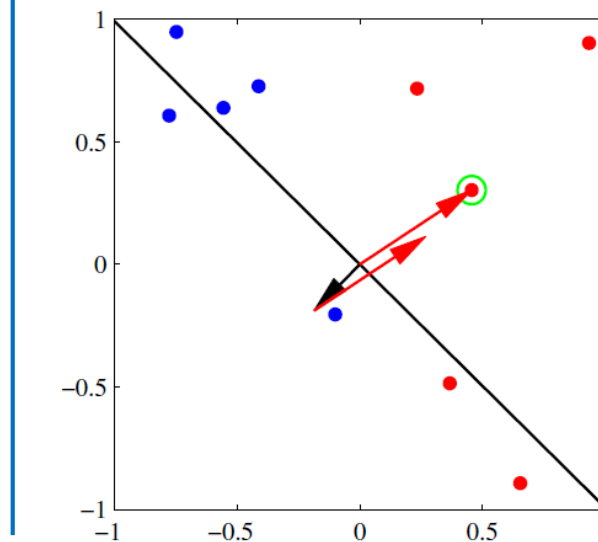
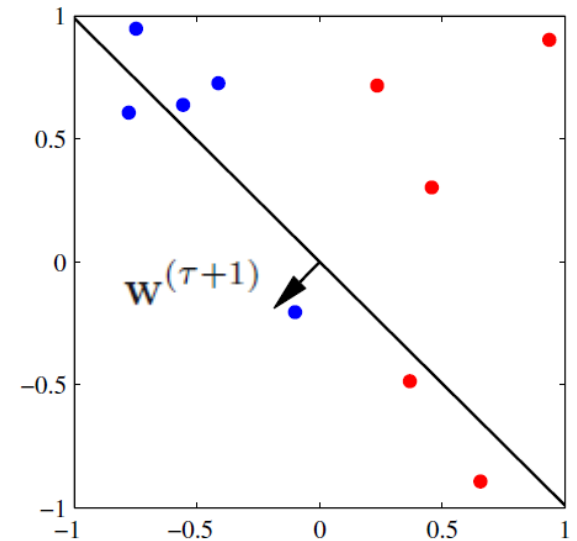
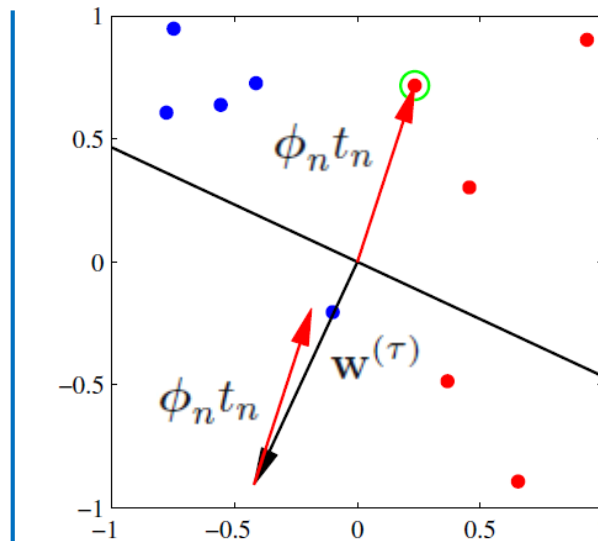
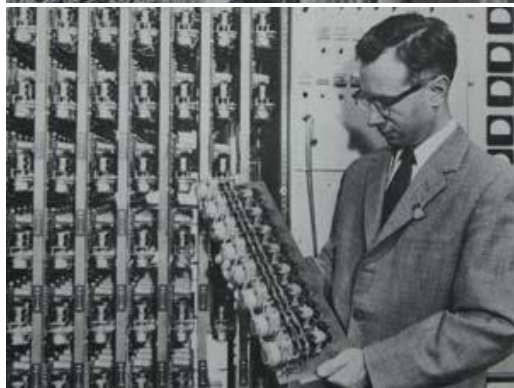
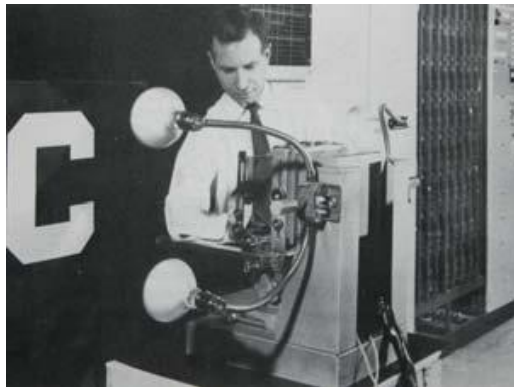
$$-\mathbf{w}^{(\tau+1)T} \phi_n t_n = -\mathbf{w}^{(\tau)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(\tau)T} \phi_n t_n$$

- Perceptron convergence theorem:

- If there exists an exact solution (in other words, if the training data set is linearly separable), then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps.*

The perceptron algorithm

- Analogue hardware implementations (Mark 1):





浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Probabilistic Generative Models



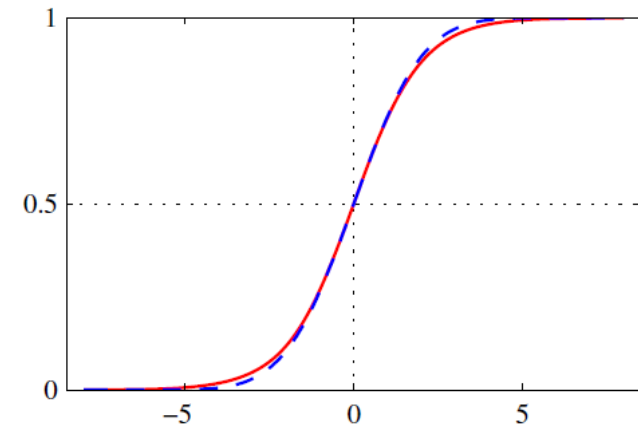
Probabilistic Generative Models

- Compute posterior probabilities by the class-conditional densities and class priors:

$$p(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} = \frac{1}{1 + \exp(-a)} = \sigma(a) \quad \text{logistic sigmoid function}$$

$$a = \ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} \quad \sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$\sigma(-a) = 1 - \sigma(a) \quad a = \ln \left(\frac{\sigma}{1 - \sigma} \right) \quad \text{logit function}$$



- K>2 classes:

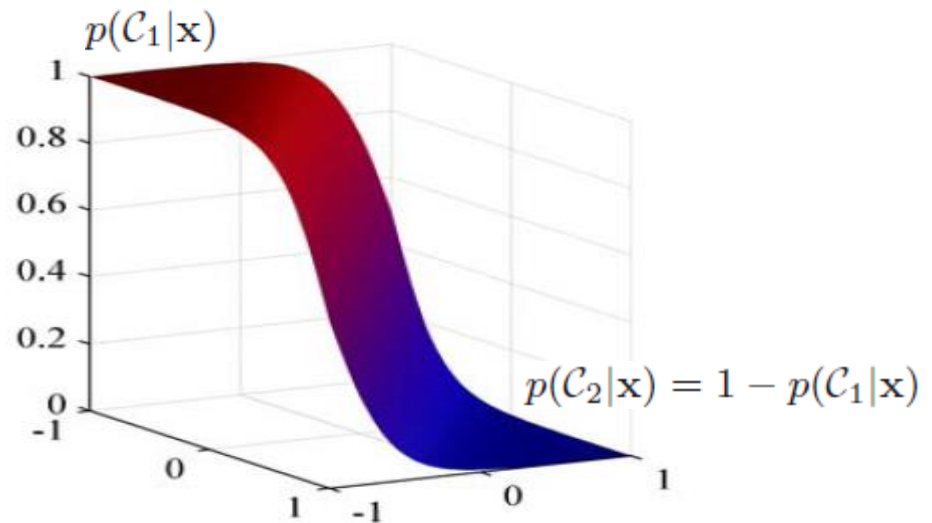
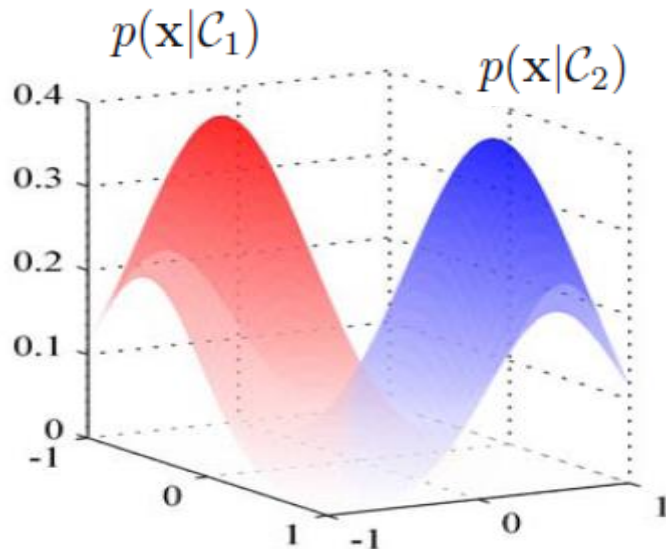
$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_j p(\mathbf{x}|C_j)p(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \ln p(\mathbf{x}|C_k)p(C_k)$$

softmax function (normalized exponential)

Continuous inputs

- Assume:
$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k) \right\}$$

- 2 classes:
$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2)$$
$$w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$



- K classes:*

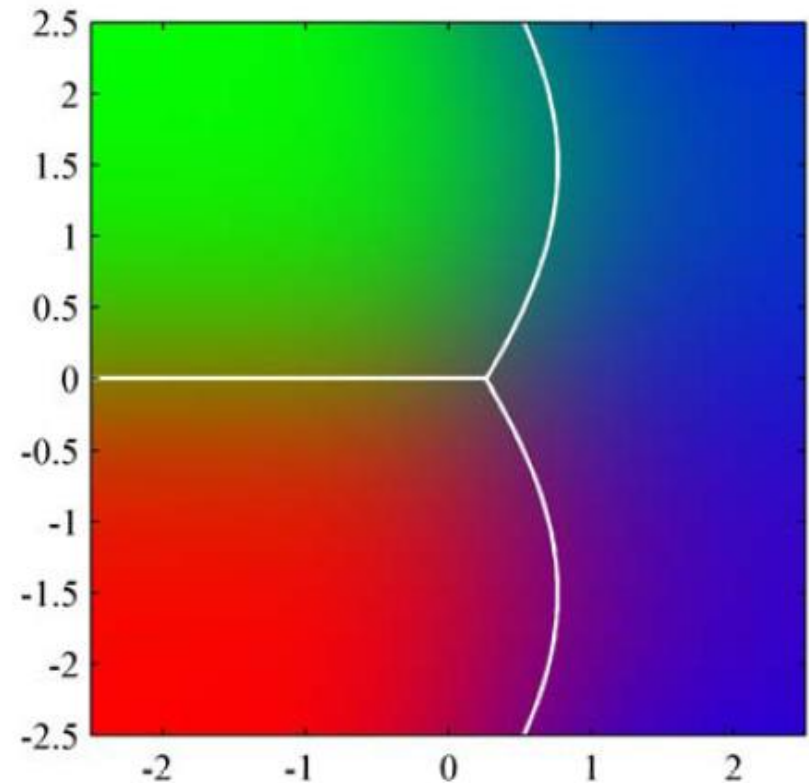
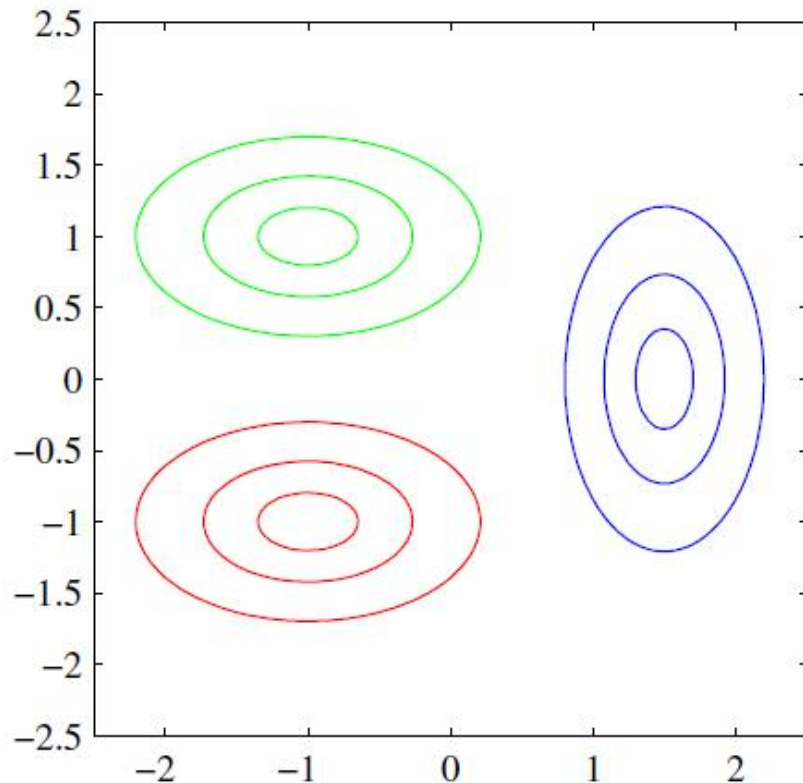
$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad \mathbf{w}_k = \Sigma^{-1} \mu_k \quad w_{k0} = -\frac{1}{2}\mu_k^T \Sigma^{-1} \mu_k + \ln p(\mathcal{C}_k)$$

Continuous inputs

- Assume: $p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k) \right\}$

- K classes with its own covariance matrix (quadratic discriminant):

$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad \mathbf{w}_k = \Sigma_k^{-1} \mu_k \quad w_{k0} = -\frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k + \ln p(\mathcal{C}_k)$$



Maximum likelihood solution for two classes

- *We have assumed (shared covariance matrix):*

$$p(\mathbf{x}_n | \mathcal{C}_k) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \quad \Rightarrow \quad p(\mathbf{x}_n | \mathcal{C}_1) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) \quad p(\mathbf{x}_n | \mathcal{C}_2) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

- *And denote the prior:*

$$p(\mathcal{C}_1) = \pi, \quad p(\mathcal{C}_2) = 1 - \pi$$

- *Hence:*

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}_n | \mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}_n | \mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

- *Now we have an input data set:*

$\{\mathbf{x}_n, t_n\}$ where $n = 1, \dots, N$. $t_n = 1$ denotes class \mathcal{C}_1 and $t_n = 0$ denotes class \mathcal{C}_2

- *Then we estimate the parameters of above model by ML.*
- *The likelihood function:*

$$p(\mathbf{t} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \quad \mathbf{t} = (t_1, \dots, t_N)^T$$

- *The log likelihood:*

$$\ln p(\mathbf{t} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \sum_{n=1}^N \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi) + t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + (1 - t_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})\}$$

Maximum likelihood solution for two classes

$$\ln p(\mathbf{t}|\pi, \mu_1, \mu_2, \Sigma) = \sum_{n=1}^N \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi) + t_n \ln \mathcal{N}(\mathbf{x}_n|\mu_1, \Sigma) + (1 - t_n) \ln \mathcal{N}(\mathbf{x}_n|\mu_2, \Sigma)\}$$

- *Solve π :*
$$\sum_{n=1}^N \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi)\} \quad \Rightarrow \quad \pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$
- *Solve μ_1, μ_2 :*
$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n|\mu_1, \Sigma) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \mu_1)^T \Sigma^{-1} (\mathbf{x}_n - \mu_1) + \text{const}$$
$$\Rightarrow \quad \mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \quad \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$
- *Solve Σ :*
$$-\frac{1}{2} \sum_{n=1}^N t_n \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \mu_1)^T \Sigma^{-1} (\mathbf{x}_n - \mu_1)$$
$$-\frac{1}{2} \sum_{n=1}^N (1 - t_n) \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \mu_2)^T \Sigma^{-1} (\mathbf{x}_n - \mu_2) = -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr} \{ \Sigma^{-1} \mathbf{S} \}$$
$$\mathbf{S} = \frac{N_1}{N} \boxed{\frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mu_1)(\mathbf{x}_n - \mu_1)^T}^{\mathbf{S}_1} + \frac{N_2}{N} \boxed{\frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mu_2)(\mathbf{x}_n - \mu_2)^T}^{\mathbf{S}_2} \quad \Rightarrow \quad \Sigma = \mathbf{S}$$

Maximum likelihood solution for K-classes

- *The likelihood function:* $p(\{\phi_n, \mathbf{t}_n\}|\{\pi_k\}) = \prod_{n=1}^N \prod_{k=1}^K \{p(\phi_n|\mathcal{C}_k)\pi_k\}^{t_{nk}} \quad \sum_k \pi_k = 1$
 - *The log likelihood:* $\ln p(\{\phi_n, \mathbf{t}_n\}|\{\pi_k\}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \{\ln p(\phi_n|\mathcal{C}_k) + \ln \pi_k\}$
 - *Introduce a Lagrange multiplier λ :* $\ln p(\{\phi_n, \mathbf{t}_n\}|\{\pi_k\}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$
-

- *Solve π_k :* $\pi_k = \frac{N_k}{N}$
 - *Solve μ_k :* $\mu_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} \phi_n$
 - *Solve Σ :* $\Sigma = \sum_{k=1}^K \frac{N_k}{N} \mathbf{S}_k \quad \mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} (\phi_n - \mu_k)(\phi_n - \mu_k)^T$
- Assumption:** Each class-conditional density is Gaussian with a shared covariance matrix.



浙江大学

ZheJiang University



人工智能研究所

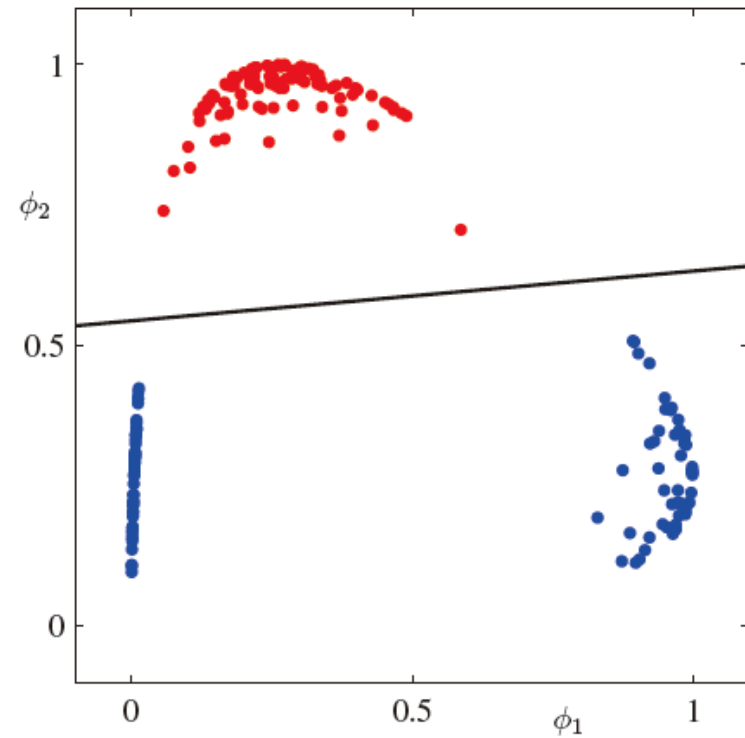
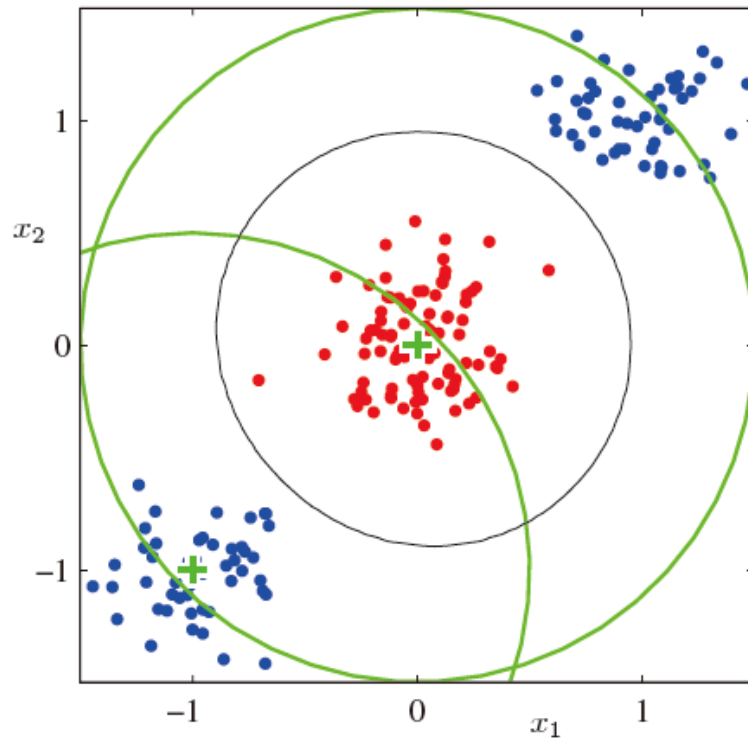
Institute of Artificial Intelligence

Probabilistic Discriminative Models



Fixed basis functions

- Classification models work on feature space instead of original input space by nonlinear basis functions:



Logistic regression

- Logistic regression model:

- Only M parameters need to be estimated.

logistic sigmoid function

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad p(C_2|\phi) = 1 - p(C_1|\phi) \quad \sigma(a) = \frac{1}{1 + \exp(-a)}$$

- For a data set $\{\phi_n, t_n\}$, where $t_n \in \{0, 1\}$ and $\phi_n = \phi(\mathbf{x}_n)$, $n = 1, \dots, N$, the likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad \text{where } \mathbf{t} = (t_1, \dots, t_N)^T \text{ and } y_n = p(C_1|\phi_n).$$

- Cross-entropy error function:

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$y_n = \sigma(a_n)$$

$$a_n = \mathbf{w}^T \phi_n$$

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

$$\left\{ \begin{array}{l} \frac{\partial E}{\partial y_n} = \frac{1 - t_n}{1 - y_n} - \frac{t_n}{y_n} = \frac{y_n(1 - t_n) - t_n(1 - y_n)}{y_n(1 - y_n)} = \frac{y_n - t_n}{y_n(1 - y_n)} \\ \frac{\partial y_n}{\partial a_n} = \frac{\partial \sigma(a_n)}{\partial a_n} = \sigma(a_n)(1 - \sigma(a_n)) = y_n(1 - y_n) \end{array} \right.$$

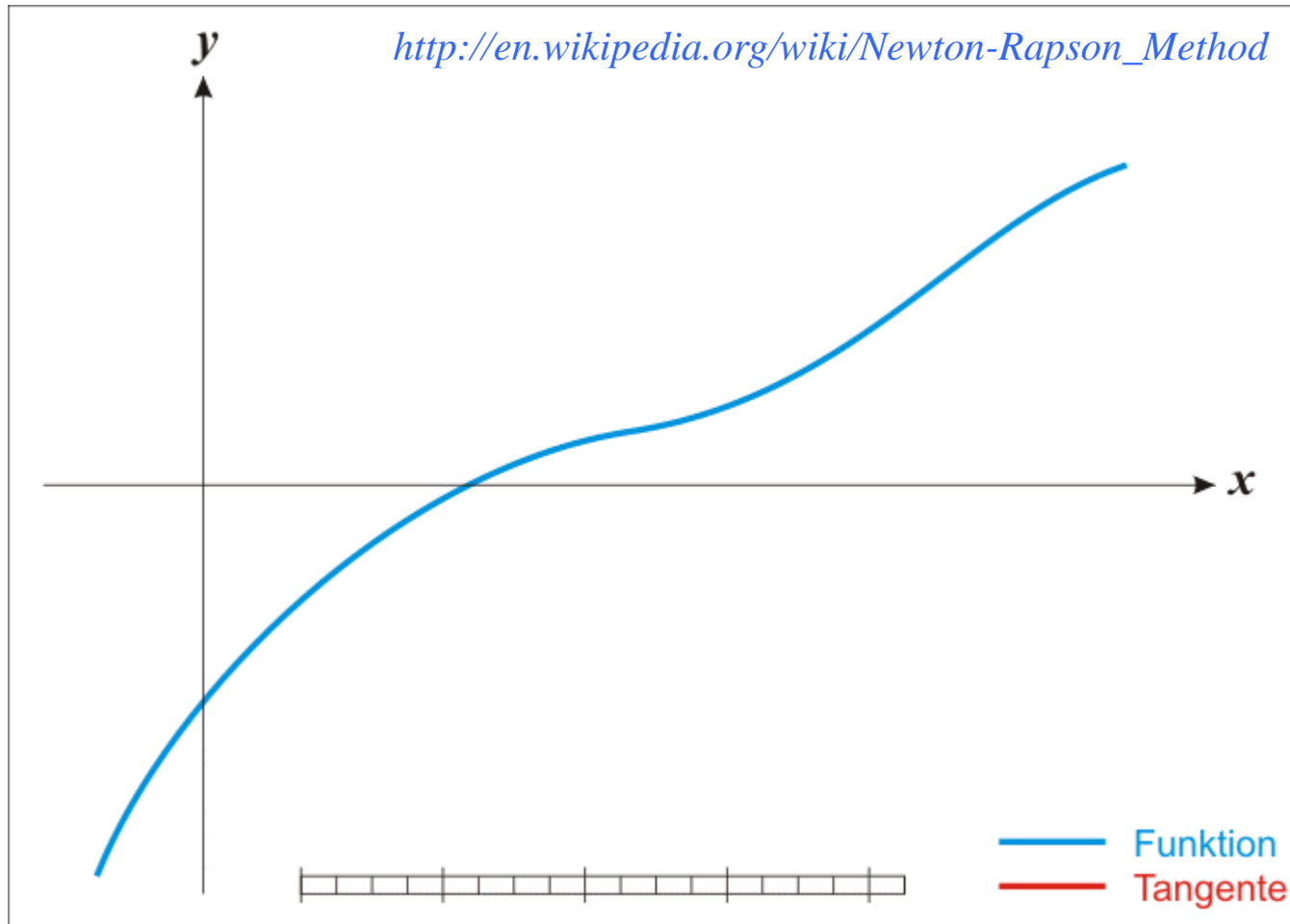
No closed-form solution

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial a_n} \nabla a_n = \sum_{n=1}^N (y_n - t_n) \phi_n = \sum_{n=1}^N (\sigma(\mathbf{w}^T \phi_n) - t_n) \phi_n$$

Iterative reweighted least squares (IRLS) algorithm

- *Newton-Raphson* iterative optimization scheme:

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$



Iterative reweighted least squares (IRLS) algorithm

- *Newton-Raphson* iterative optimization scheme: $\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$

- For linear regression model with the sum-of-squares error function:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t} \quad \mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi$$

$$\Rightarrow \mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - (\Phi^T \Phi)^{-1} \{ \Phi^T \Phi \mathbf{w}^{(\text{old})} - \Phi^T \mathbf{t} \} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- For logistic regression model with cross-entropy error function:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t}) \quad \mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi$$

$$\Rightarrow \mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) = (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}$$

$$\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$$

local linear approximation to the logistic sigmoid function around the current operating point $\mathbf{w}^{(\text{old})}$

$$a_n(\mathbf{w}) \simeq a_n(\mathbf{w}^{(\text{old})}) + \left. \frac{da_n}{dy_n} \right|_{\mathbf{w}^{(\text{old})}} (t_n - y_n) \\ = \phi_n^T \mathbf{w}^{(\text{old})} - \frac{(y_n - t_n)}{y_n(1 - y_n)} = z_n.$$

Multiclass logistic regression

- Models:

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \implies p(\mathcal{C}_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \mathbf{w}_k^T \phi$$

- Likelihood function (1-of-K coding scheme):

$$p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(\mathcal{C}_k|\phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad y_{nk} = y_k(\phi_n)$$

- Cross-entropy error function:

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j) \implies \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n \quad \sum_k t_{nk} = 1$$

Exercise 4.18

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n \implies \mathbf{w}_1, \dots, \mathbf{w}_K \quad \text{Sequential learning algorithm to update } w \text{ one by one}$$

- Solve by IRLS algorithm:

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

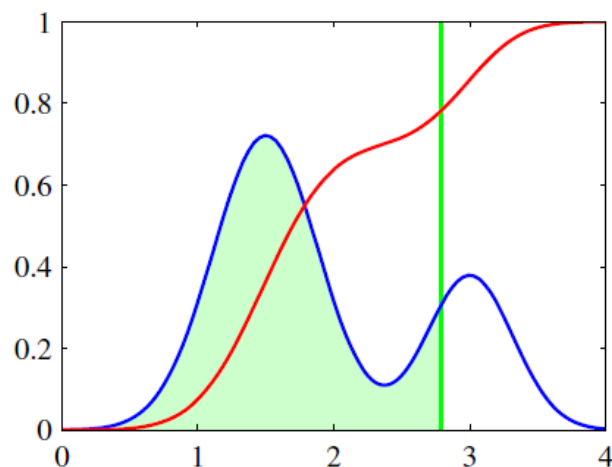
*Block j,k of Hessian matrix (comprise blocks of size M*M)*

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T$$



Another PDM: Probit regression

- Use CDF to construct an activation function:



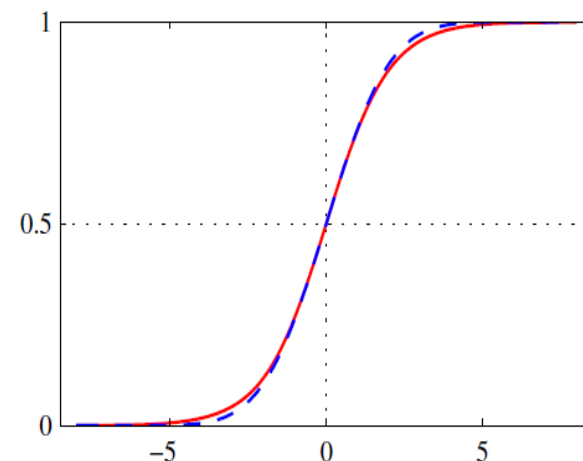
Logistic sigmoid function

$$p(C_1|\phi) = \frac{1}{1 + \exp(\mathbf{w}^T \phi)}$$



$$p(C_1|\phi) = \int_{-\infty}^{\mathbf{w}^T \phi} \mathcal{N}(\theta|0, 1) d\theta$$

Inverse probit function





Canonical link functions

- Canonical link function:

$$p(t|\eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\} \quad y \equiv \mathbb{E}[t|\eta] = -s \frac{d}{d\eta} \ln g(\eta)$$

- Generalized linear model: $y = f(\mathbf{w}^T \phi)$
 - $f(\cdot)$ is the activation function and $f^{-1}(\cdot)$ is the link function
- Log likelihood function:

$$\ln p(\mathbf{t}|\eta, s) = \sum_{n=1}^N \ln p(t_n|\eta, s) = \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} + \text{const}$$

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\eta, s) = \sum_{n=1}^N \left\{ \frac{d}{d\eta_n} \ln g(\eta_n) + \frac{t_n}{s} \right\} \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \nabla a_n = \sum_{n=1}^N \frac{1}{s} \{t_n - y_n\} \psi'(y_n) f'(a_n) \phi_n$$

$$\xrightarrow{f^{-1}(y) = \psi(y)} \quad \nabla E(\mathbf{w}) = \frac{1}{s} \sum_{n=1}^N \{y_n - t_n\} \phi_n$$

For the Gaussian $s = \beta^{-1}$,
for the logistic model $s = 1$.



浙江大学

ZheJiang University



人工智能研究所

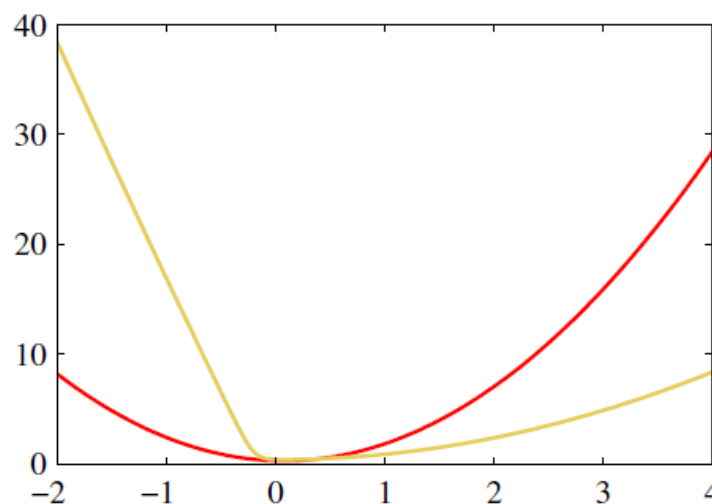
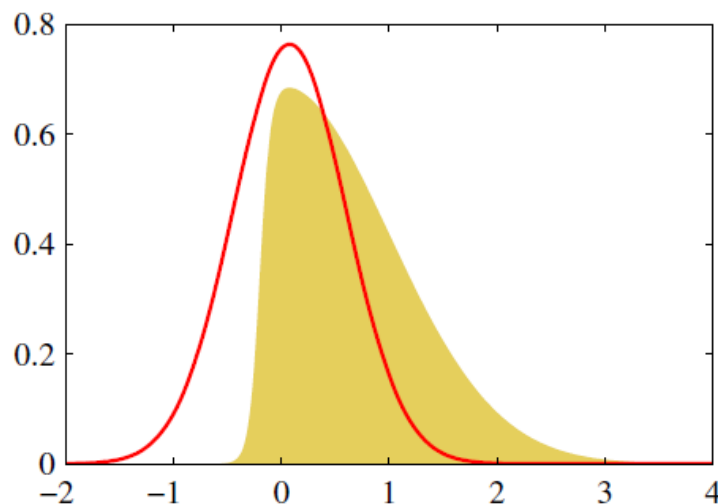
Institute of Artificial Intelligence

The Laplace Approximation



The Laplace Approximation

- Find a Gaussian approximation $q(z)$ which is centred on a mode of the distribution $p(z)$:



$$p(z) = \frac{1}{Z} f(z) \quad \ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A (z - z_0)^2 \quad f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$
$$Z = \int f(z) dz \quad A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0} \quad q(z) = \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$



The Laplace Approximation

$$p(z) = \frac{1}{Z} f(z) \quad f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\} \quad Z = \int f(z) dz$$
$$q(z) = \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\} \quad A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

- Extend to M-dimensional space: $p(\mathbf{z}) = f(\mathbf{z})/Z$

$$Z = \int f(\mathbf{z}) d\mathbf{z} \simeq f(\mathbf{z}_0) \int \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} d\mathbf{z} = f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}}$$

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0)$$

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} \quad \mathbf{A} = - \nabla \nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0}$$

➡ $q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1})$



浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Bayesian Logistic Regression

Bayesian Logistic Regression

- In logistic regression model, we have:

- For a data set $\{\phi_n, t_n\}$, where $t_n \in \{0, 1\}$ and $\phi_n = \phi(\mathbf{x}_n)$, $n = 1, \dots, N$, the likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad \text{where } \mathbf{t} = (t_1, \dots, t_N)^T \text{ and } y_n = p(\mathcal{C}_1|\phi_n).$$

- Now assume prior is Gaussian: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$
- Then posterior distribution over w (*Obviously it's not a Gaussian, but we can find its Gaussian approximation by Laplace approximation framework*):

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad \begin{aligned} y_n &= \sigma(a_n) \\ a_n &= \mathbf{w}^T \phi_n \end{aligned}$$

➡
$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \text{const}$$

- Maximize above posterior distribution to give the MAP solution \mathbf{w}_{MAP} , which defines the mean of the Gaussian.
- The covariance is then give by the inverse of the matrix of second derivatives of the negative log likelihood:

$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n) \phi_n \phi_n^T$$

- So, the Gaussian approximation:

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N^{-1})$$

Bayesian Logistic Regression

- For new input vector x , corresponding feature vector is $\phi(x)$
- Then the predictive distribution for class C_1 :

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_N^{-1})$$

$$p(C_1 | \phi, \mathbf{t}) = \int p(C_1 | \phi, \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad p(C_2 | \phi, \mathbf{t}) = 1 - p(C_1 | \phi, \mathbf{t})$$

- We can rewrite (*Mathematics structural trick*):

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da \quad \delta(\cdot) \text{ is the Dirac delta function}$$

$$\Rightarrow \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) p(a) da \quad \text{where } p(a) = \int \delta(a - \mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$$

$$\mu_a = \mathbb{E}[a] = \int p(a) a da = \int q(\mathbf{w}) \mathbf{w}^T \phi d\mathbf{w} = \mathbf{w}_{\text{MAP}}^T \phi$$

$$\sigma_a^2 = \text{var}[a] = \int p(a) \{a^2 - \mathbb{E}[a]^2\} da = \int q(\mathbf{w}) \{(\mathbf{w}^T \phi)^2 - (\mathbf{m}_N^T \phi)^2\} d\mathbf{w} = \phi^T \mathbf{S}_N^{-1} \phi$$

$$\Rightarrow p(C_1 | \mathbf{t}) = \int \sigma(a) p(a) da = \int \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da$$

$$\text{apply the approximation } \sigma(a) \simeq \Phi(\lambda a) \quad \lambda^2 = \pi/8 \quad \kappa(\sigma^2) = (1 + \pi\sigma^2/8)^{-1/2}$$

$$\Rightarrow p(C_1 | \mathbf{t}) \simeq \int \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right) \simeq \sigma(\kappa(\sigma^2)\mu) = \sigma(\kappa(\sigma_a^2)\mu_a)$$



Next: Kernel Methods and SVM

- HW4:
 - 4.4, 4.5, 4.8, 4.17, 4.18, 4.19
 - Programming work: see website for details.