

 $MSE(\hat{B}) = E[\hat{B}-B] = \sigma^2(\hat{B}) + B(\hat{B})$ 



# Lecture 7: Linear Classification Methods

Final projects? Groups Topics Proposal week 5 Lecture 20 is poster session, Jacobs Hall Lobby, snacks Final report 15 June.

## What is "linear" classification?

Classification is intrinsically non-linear

It puts non-identical things in the same class, so a difference in input vector sometimes causes zero change in the answer

"Linear classification" means that the part that adapts is linear

The adaptive part is followed by a fixed non-linearity. It may be preceded by a fixed non-linearity (e.g. nonlinear basis functions).



### Representing the target values for classification

For two classes, we use a single valued output that has target values **1** for the "positive" class and **0** (or **-1**) for the other class

For probabilistic class labels the target value can then be **P(t=1)** and the model output can also represent **P(y=1)**.

For **N** classes we often use a vector of N target values containing a single **1** for the correct class and zeros elsewhere.

For probabilistic labels we can then use a vector of class probabilities as the target vector.





## Three approaches to classification

Use discriminant functions directly without probabilities:

Convert input vector into real values. A simple operation (like thresholding) can get the class.

Choose real values to maximize the useable information about the class label that is in the real value.

Vinfer conditional class probabilities:  $p(class = C_k | \mathbf{x})$ 

Compute the conditional probability of each class.

Then make a decision that minimizes some loss function

Compare the probability of the input under separate, classspecific, generative models.

E.g. fit a multivariate **Gaussian** to the input vectors of each class and see which Gaussian makes a test data vector most probable. (Is this the best bet?)



r= 3 11 W/12

# Discriminant functions for N>2 classes

One possibility is to use N two-way discriminant functions.

Each function discriminates one class from the rest. Another possibility is to use N(N-1)/2 two-way discriminant functions

Each function discriminates between two particular classes.

Both these methods have problems





51

Two-way preferences need not be transitive!

## A simple solution (4.1.2)



# Maximum Likelihood and Least Squares (from lecture 3)

### Computing the gradient and setting it to zero yields

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w},\beta) = \beta \sum_{n=1}^{N} \left\{ t_n - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n) \right\} \boldsymbol{\phi}(\mathbf{x}_n)^{\mathrm{T}} = \mathbf{0}.$$

Solving for w,

 $\mathbf{w}_{\mathrm{ML}} = \left( \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \right)^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{t}$ 

where

$$\mathbf{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \overline{\cdot}$$

The Moore-Penrose

pseudo-inverse,  $\Phi^{\dagger}$ .

#### LSQ for classification

Each class  $C_k$  is described by its own linear model so that

$$y_k(\mathbf{x}) = \mathbf{w}_k^{\mathrm{T}} \mathbf{x} + w_{k0} \quad \checkmark \tag{4.13}$$

where  $k = 1, \ldots, K$ . We can conveniently group these together using vector notation so that (4.14)

$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^{\mathrm{T}} \widetilde{\mathbf{x}}$$

Consider a training set  $\{x_n, t_n\}, n = 1 \dots N$ Define X and T  $T = \begin{bmatrix} t'_1 & t_N \end{bmatrix}$ X = ×2 N

LSQ solution:  

$$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^{\mathrm{T}}\widetilde{\mathbf{X}})^{-1}\widetilde{\mathbf{X}}^{\mathrm{T}}\mathbf{T} = \widetilde{\mathbf{X}}^{\dagger}\mathbf{T}$$
(4.16)  
And prediction  
 $\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^{\mathrm{T}}\widetilde{\mathbf{x}} = \mathbf{T}^{\mathrm{T}}\left(\widetilde{\mathbf{X}}^{\dagger}\right)^{\mathrm{T}}\widetilde{\mathbf{x}}.$ 
(4.17)

ton K P

K

# Using "least squares" for classification It does not work as well as better methods, but it is easy: It reduces classification to least squares regression.



# PCA don't work well



# picture showing the advantage of Fisher's linear discriminant



When projected onto the line joining the class means, the classes are not well separated.



Fisher chooses a direction that makes the projected classes much tighter, even though their projected means are less far apart.

### Math of Fisher's linear discriminants

What linear transformation is best for discrimination?

The projection onto the vector separating the class means seems sensible:

But we also want small variance within each  $s_1^2 = \sum (y_n - m_1)$ class:

$$\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$$

 $n \varepsilon C_1$ 

 $y = \mathbf{w}^T \mathbf{x}$ 

$$s_2^2 = \sum_{n \varepsilon C_2} (y_n - m_2)$$

Fisher's objective function is:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \stackrel{\text{between}}{\longleftarrow}$$
within

### More math of Fisher's linear discriminants

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

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

$$\mathbf{S}_{W} = \sum_{n \in C_{1}} (\mathbf{x}_{n} - \mathbf{m}_{1}) (\mathbf{x}_{n} - \mathbf{m}_{1})^{T} + \sum_{n \in C_{2}} (\mathbf{x}_{n} - \mathbf{m}_{2}) (\mathbf{x}_{n} - \mathbf{m}_{2})^{T}$$

*Optimal solution*:  $\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$ 

# We have probalistic classification!



## Probabilistic Generative Models for Discrimination (Bishop p 196)

Use a generative model of the input vectors for each class, and see which model makes a test input vector most probable. The posterior probability of class 1 is:

$$p(C_{1} | \mathbf{x}) = \frac{p(C_{1})p(\mathbf{x} | C_{1})}{p(C_{1})p(\mathbf{x} | C_{1}) + p(C_{0})p(\mathbf{x} | C_{0})} = \frac{1}{1 + e^{-z}}$$
  
where  $z = \ln \frac{p(C_{1})p(\mathbf{x} | C_{1})}{p(C_{0})p(\mathbf{x} | C_{0})} = \ln \frac{p(C_{1} | \mathbf{x})}{1 - p(C_{1} | \mathbf{x})}$ 

z is called the **logit** and is given by the **log odds** 

### An example for continuous inputs

Assume input vectors for each class are Gaussian, all classes have the same covariance matrix.

$$p(\mathbf{x} | C_k) = a \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mathbf{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \mathbf{\mu}_k) \right\}$$

For two classes,  $C_1$  and  $C_0$ , the posterior is a logistic:

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
  

$$\mathbf{w} = \mathbf{\Sigma}^{-1}(\mathbf{\mu}_1 - \mathbf{\mu}_0)$$
  

$$w_0 = -\frac{1}{2}\mathbf{\mu}_1^T \mathbf{\Sigma}^{-1} \mathbf{\mu}_1 + \frac{1}{2}\mathbf{\mu}_0^T \mathbf{\Sigma}^{-1} \mathbf{\mu}_0 + \ln \frac{p(C_1)}{p(C_0)}$$

$$z = ln \left( \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} \right)$$
  
=  $ln P(x|C_1) - ln (x|C_2) + ln(c)$   
=  $-\frac{1}{2}(x-\mu)^T z^{-1}(x-\mu) + \frac{1}{2}(x-\mu)^T z^{-(x-\mu_2)} + ccnst$   
=  $\mu_1^T \overline{z}^{-1} x - \mu_2^T \overline{z}^{-1} x + \mu_1^T \overline{z}^{-1} \mu_1 + \mu_2 \overline{z}^{-1} \mu_2 + ccnst$   
=  $W^T x + W_0$   
 $W^T = \mu_1^T \overline{z}^{-1} - \mu_2^T \overline{z}^{-1}$   
 $W = \overline{z}^{-1}(\mu_1 - \mu_2)$ 

### The role of the inverse covariance matrix

If the Gaussian is spherical no need to worry about the covariance matrix.

So, start by transforming the data space to make the Gaussian spherical

This is called "whitening" the data.

It pre-multiplies by the matrix **square root of the inverse covariance matrix**.

In transformed space, the weight vector is the difference between transformed means.

$$\mathbf{w} = \mathbf{\Sigma}^{-1} (\mathbf{\mu}_1 - \mathbf{\mu}_0)$$
  
gives the same value  
for  $\mathbf{w}^T \mathbf{x}$  as:  
$$\mathbf{w}_{aff} = \mathbf{\Sigma}^{-\frac{1}{2}} \mathbf{\mu}_1 - \mathbf{\Sigma}^{-\frac{1}{2}} \mathbf{\mu}_0$$
  
and  $\mathbf{x}_{aff} = \mathbf{\Sigma}^{-\frac{1}{2}} \mathbf{x}$   
gives for  $\mathbf{w}_{aff}^T \mathbf{x}_{aff}$ 

aff = affine







The decision surface is planar when the covariance matrices are the same and quadratic when not.

#### Bernoulli distribution Random variable $x \in \{0,1\}$ Coin flipping: heads=1, tails=0 $p(x = 1|\mu) = \mu$ $\checkmark$ Bernoulli Distribution $\operatorname{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x}$ $\mathbb{E}[x] = \mu$ $\operatorname{var}[x] = \mu(1-\mu)$

ML for Bernoulli

Given: 
$$\mathcal{D} = \{x_1, \dots, x_N\}, m \text{ heads (1), } N - m \text{ tails (0)}$$
  
 $p(\mathcal{D}|\mu) = \prod_{n=1}^{N} p(x_n|\mu) = \prod_{n=1}^{N} \mu_{\cdot n}^{x_n} (1-\mu)^{1-x_n}$   
 $\ln p(\mathcal{D}|\mu) = \sum_{n=1}^{N} \ln p(x_n|\mu) = \sum_{n=1}^{N} \{x_n \ln \mu + (1-x_n) \ln(1-\mu)\}$   
 $\mu_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} x_n = \frac{m}{N}$ 

### The logistic function

The output is a smooth function of the inputs and the weights.





Its odd to express it in terms of y.



Logistic regression (Bishop 205)  

$$p(C_{1}|x) = \sigma(w^{T}x).$$
Observations  $[X_{n}, \pm_{n}] = t_{n} \in [0, 1]$ 
Likelihood  

$$y = \sigma(w^{T}x)$$

$$p(y|x, w) = \operatorname{Bern}(g_{1}|u) = \mu^{g}(i-\mu)^{(-\mu)}$$

$$p(T|x, w) = \prod_{n}^{N} g_{n}^{\pm_{n}} (i-\mu_{n})^{l-\pm_{n}}$$
Log-likelihood  

$$E_{w} = -\ln(p(T|x, w)) = -\sum_{n}^{N} (t_{n} \ln g_{n} + (i-t_{n})\ln(i-g_{n}))$$
Minimize -log like  
Derivative  

$$\nabla_{w}E_{w} = -\sum_{n}^{N} [t_{n} - \frac{1}{g_{n}} + \frac{1-t_{n}}{1-g_{n}} \cdot (-1)] g_{n}(1-g_{n}) X_{n}$$

$$= \sum_{n}^{N} (t_{n} - \frac{M}{g_{n}}) X_{n}$$

$$W_{L_{n}} = W_{L_{n}} - \eta S E_{w}$$

### Logistic regression (page 205)

When there are only two classes we can model the conditional probability of the positive class as

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
 where  $\sigma(z) = \frac{1}{1 + \exp(-z)}$ 

If we use the right error function, something nice happens: The gradient of the logistic and the gradient of the error function cancel each other:

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}), \qquad \nabla E(\mathbf{w}) = \sum_{n=1}^{N} (y_n - t_n) \mathbf{x}_n$$