

## HT2015: SC4

### Statistical Data Mining and Machine Learning

Dino Sejdinovic  
Department of Statistics  
Oxford

<http://www.stats.ox.ac.uk/~sejdinov/sdmml.html>

## Supervised Learning

### Unsupervised learning:

- To “extract structure” and postulate hypotheses about data generating process from “unlabelled” observations  $x_1, \dots, x_n$ .
- Visualize, summarize and compress data.

### Supervised learning:

- In addition to the observations of  $X$ , we have access to their response variables / labels  $Y \in \mathcal{Y}$ : we observe  $\{(x_i, y_i)\}_{i=1}^n$ .
- Types of supervised learning:
  - Classification: discrete responses, e.g.  $\mathcal{Y} = \{+1, -1\}$  or  $\{1, \dots, K\}$ .
  - Regression: a numerical value is observed and  $\mathcal{Y} = \mathbb{R}$ .

The goal is to accurately predict the response  $Y$  on new observations of  $X$ , i.e., to **learn a function**  $f: \mathbb{R}^p \rightarrow \mathcal{Y}$ , such that  $f(X)$  will be close to the true response  $Y$ .

## Regression Example: Boston Housing

The original data are 506 observations on 13 variables  $X$ ; medv is the response variable  $Y$ .

```

crim    per capita crime rate by town
zn      proportion of residential land zoned for lots
        over 25,000 sq.ft
indus   proportion of non-retail business acres per town
chas    Charles River dummy variable (= 1 if tract bounds river;
        0 otherwise)
nox     nitric oxides concentration (parts per 10 million)
rm      average number of rooms per dwelling
age     proportion of owner-occupied units built prior to 1940
dis     weighted distances to five Boston employment centers
rad     index of accessibility to radial highways
tax     full-value property-tax rate per USD 10,000
ptratio pupil-teacher ratio by town
b       1000(B - 0.63)^2 where B is the proportion of blacks by town
lstat   percentage of lower status of the population
medv    median value of owner-occupied homes in USD 1000's

```

## Regression Example: Boston Housing

```

> str(X)
'data.frame':  506 obs. of  13 variables:
 $ crim   : num  0.00632 0.02731 0.02729 0.03237 0.06905 ...
 $ zn     : num  18 0 0 0 0 0 12.5 12.5 12.5 12.5 ...
 $ indus  : num  2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 ...
 $ chas   : int   0 0 0 0 0 0 0 0 0 ...
 $ nox    : num  0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 0.524 0.5
 $ rm     : num  6.58 6.42 7.18 7.00 7.15 ...
 $ age    : num  65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9 ...
 $ dis    : num  4.09 4.97 4.97 6.06 6.06 ...
 $ rad    : int   1 2 2 3 3 3 5 5 5 ...
 $ tax    : num  296 242 242 222 222 222 311 311 311 311 ...
 $ ptratio: num  15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
 $ black  : num  397 397 393 395 397 ...
 $ lstat  : num  4.98 9.14 4.03 2.94 5.33 ...

> str(Y)
num[1:506] 24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...

```

Goal: predict median house price  $Y$  given 13 predictor variables  $X$  of a new district.

## Classification Example: Lymphoma

We have gene expression measurements  $X$  of  $n = 62$  patients for  $p = 4026$  genes. For each patient,  $Y \in \{0, 1\}$  denotes one of two subtypes of cancer. Goal: predict cancer subtype given gene expressions of a new patient.

```
> str(X)
'data.frame': 62 obs. of 4026 variables:
 $ Gene 1 : num -0.344 -1.188 0.520 -0.748 -0.868 ...
 $ Gene 2 : num -0.953 -1.286 0.657 -1.328 -1.330 ...
 $ Gene 3 : num -0.776 -0.588 0.409 -0.991 -1.517 ...
 $ Gene 4 : num -0.474 -1.588 0.219 0.978 -1.604 ...
 $ Gene 5 : num -1.896 -1.960 -1.695 -0.348 -0.595 ...
 $ Gene 6 : num -2.075 -2.117 0.121 -0.800 0.651 ...
 $ Gene 7 : num -1.875 -1.818 0.317 0.387 0.041 ...
 $ Gene 8 : num -1.539 -2.433 -0.337 -0.522 -0.668 ...
 $ Gene 9 : num -0.604 -0.710 -1.269 -0.832 0.458 ...
 $ Gene 10 : num -0.218 -0.487 -1.203 -0.919 -0.848 ...
 $ Gene 11 : num -0.340 1.164 1.023 1.133 -0.541 ...
 $ Gene 12 : num -0.531 0.488 -0.335 0.496 -0.358 ...

> str(Y)
num [1:62] 0 0 0 1 0 0 1 0 0 0 ...
```

## Loss function

- Suppose we made a prediction  $\hat{Y} = f(X) \in \mathcal{Y}$  based on observation of  $X$ .
- How good is the prediction? We can use a **loss function**  $L: \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}^+$  to formalize the quality of the prediction.
- Typical loss functions:

- **Misclassification loss** (or **0-1 loss**) for classification

$$L(Y, f(X)) = \begin{cases} 0 & f(X) = Y \\ 1 & f(X) \neq Y \end{cases}.$$

- **Squared loss** for regression

$$L(Y, f(X)) = (f(X) - Y)^2.$$

- Many other choices are possible, e.g., **weighted misclassification loss**.
- In classification, if estimated probabilities  $\hat{p}(k)$  for each class  $k \in \mathcal{Y}$  are returned, **log-likelihood loss** (or **log loss**)  $L(Y, \hat{p}) = -\log \hat{p}(Y)$  is often used.

## Risk

- paired observations  $\{(x_i, y_i)\}_{i=1}^n$  viewed as i.i.d. realizations of a random variable  $(X, Y)$  on  $\mathcal{X} \times \mathcal{Y}$  with joint distribution  $P_{XY}$

### Risk

For a given loss function  $L$ , the **risk**  $R$  of a learned function  $f$  is given by the expected loss

$$R(f) = \mathbb{E}_{P_{XY}} [L(Y, f(X))],$$

where the expectation is with respect to the true (unknown) joint distribution of  $(X, Y)$ .

- The risk is unknown, but we can compute the **empirical risk**:

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)).$$

## The Bayes Classifier

- What is the optimal classifier if the joint distribution  $(X, Y)$  were known?
- The density  $g$  of  $X$  can be written as a mixture of  $K$  components (corresponding to each of the classes):

$$g(x) = \sum_{k=1}^K \pi_k g_k(x),$$

where, for  $k = 1, \dots, K$ ,

- $\mathbb{P}(Y = k) = \pi_k$  are the class probabilities,
- $g_k(x)$  is the conditional density of  $X$ , given  $Y = k$ .
- The **Bayes classifier**  $f_{\text{Bayes}}: x \mapsto \{1, \dots, K\}$  is the one with minimum risk:

$$\begin{aligned} R(f) &= \mathbb{E} [L(Y, f(X))] = \mathbb{E}_X [\mathbb{E}_{Y|X} [L(Y, f(X)) | X]] \\ &= \int_{\mathcal{X}} \mathbb{E} [L(Y, f(X)) | X = x] g(x) dx \end{aligned}$$

- The minimum risk attained by the Bayes classifier is called **Bayes risk**.
- Minimizing  $\mathbb{E}[L(Y, f(X)) | X = x]$  separately for each  $x$  suffices.

## The Bayes Classifier

- Consider the 0-1 loss.
- The risk simplifies to:

$$\begin{aligned}\mathbb{E}[L(Y, f(X)) | X = x] &= \sum_{k=1}^K L(k, f(x)) \mathbb{P}(Y = k | X = x) \\ &= 1 - \mathbb{P}(Y = f(x) | X = x)\end{aligned}$$

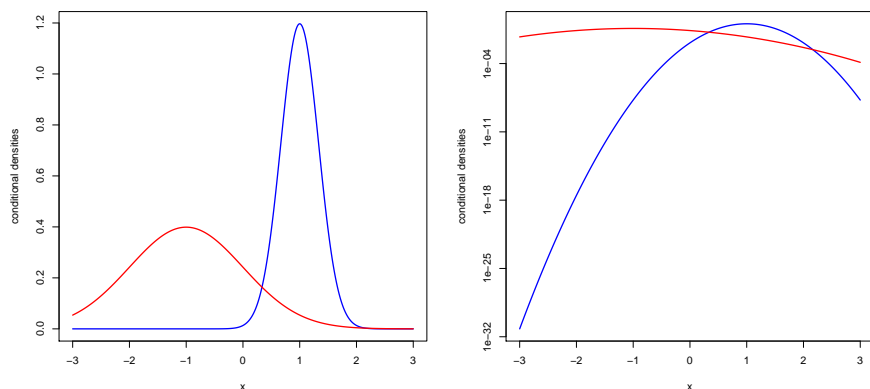
- The risk is minimized by choosing the class with the greatest posterior probability:

$$\begin{aligned}f_{\text{Bayes}}(x) &= \arg \max_{k=1, \dots, K} \mathbb{P}(Y = k | X = x) \\ &= \arg \max_{k=1, \dots, K} \frac{\pi_k g_k(x)}{\sum_{j=1}^K \pi_j g_j(x)} = \arg \max_{k=1, \dots, K} \pi_k g_k(x).\end{aligned}$$

- The functions  $x \mapsto \pi_k g_k(x)$  are called **discriminant functions**. The discriminant function with maximum value determines the predicted class of  $x$ .

## The Bayes Classifier: Example

How do you classify a new observation  $x$  if now the standard deviation is still 1 for class 1 but  $1/3$  for class 2?

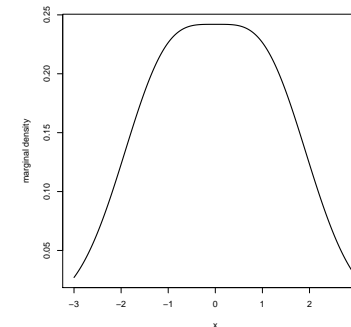
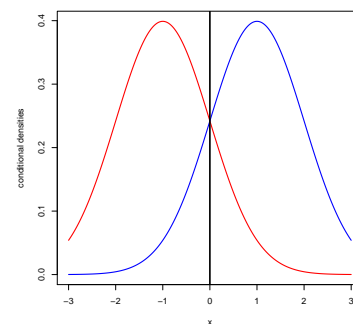


Looking at density in a log-scale, optimal classification is to select class 2 if and only if  $x \in [0.34, 2.16]$ .

## The Bayes Classifier: Example

A simple two Gaussians example: Suppose  $X \sim \mathcal{N}(\mu_Y, 1)$ , where  $\mu_1 = -1$  and  $\mu_2 = 1$  and assume equal priors  $\pi_1 = \pi_2 = 1/2$ .

$$g_1(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x+1)^2}{2}\right) \quad \text{and} \quad g_2(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-1)^2}{2}\right).$$



Optimal classification is  $f_{\text{Bayes}}(x) = \arg \max_{k=1, \dots, K} \pi_k g_k(x) = \begin{cases} 1 & \text{if } x < 0, \\ 2 & \text{if } x \geq 0. \end{cases}$

## Plug-in Classification

- The Bayes Classifier chooses the class with the greatest posterior probability

$$f_{\text{Bayes}}(x) = \arg \max_{k=1, \dots, K} \pi_k g_k(x).$$

- We know neither the conditional densities  $g_k$  nor the class probabilities  $\pi_k$ !
- The **plug-in classifier** chooses the class

$$f(x) = \arg \max_{k=1, \dots, K} \hat{\pi}_k \hat{g}_k(x),$$

- where we plugged in
  - estimates  $\hat{\pi}_k$  of  $\pi_k$  and  $k = 1, \dots, K$  and
  - estimates  $\hat{g}_k(x)$  of conditional densities,
- Linear Discriminant Analysis** is an example of plug-in classification.

## Linear Discriminant Analysis

- **LDA** is the most well-known and simplest example of plug-in classification.
- Assume multivariate normal conditional density  $g_k(x)$  for each class  $k$ :

$$X|Y = k \sim \mathcal{N}(\mu_k, \Sigma),$$

$$g_k(x) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(x - \mu_k)^\top \Sigma^{-1}(x - \mu_k)\right),$$

- each class can have a **different mean**  $\mu_k$ ,
- all classes share the **same covariance**  $\Sigma$ .
- For an observation  $x$ , the  $k$ -th log-discriminant function is

$$\log \pi_k g_k(x) = c + \log \pi_k - \frac{1}{2}(x - \mu_k)^\top \Sigma^{-1}(x - \mu_k)$$

The quantity  $(x - \mu_k)^\top \Sigma^{-1}(x - \mu_k)$  is the squared **Mahalanobis distance** between  $x$  and  $\mu_k$ .

- If  $\Sigma = I_p$  and  $\pi_k = \frac{1}{K}$ , LDA simply chooses the class  $k$  with the nearest (in the Euclidean sense) class mean.

## Parameter Estimation

- How to estimate the parameters of the LDA model?
- We can achieve this by maximum likelihood (EM algorithm is not needed here since the class variables  $y_i$  are observed!).
- Let  $n_k = \#\{j : y_j = k\}$  be the number of observations in class  $k$ .

$$\begin{aligned} \ell(\pi, (\mu_k)_{k=1}^K, \Sigma) &= \log p((x_i, y_i)_{i=1}^n | \pi, (\mu_k)_{k=1}^K, \Sigma) = \sum_{i=1}^n \log \pi_{y_i} g_{y_i}(x_i) \\ &= c + \sum_{k=1}^K \sum_{j:y_j=k} \log \pi_k - \frac{1}{2} \left( \log |\Sigma| + \sum_{j:y_j=k} (x_j - \mu_k)^\top \Sigma^{-1}(x_j - \mu_k) \right) \end{aligned}$$

ML estimates:

$$\hat{\pi}_k = \frac{n_k}{n} \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{j:y_j=k} x_j$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^K \sum_{j:y_j=k} (x_j - \hat{\mu}_k)(x_j - \hat{\mu}_k)^\top$$

- Note: the ML estimate of  $\Sigma$  is biased. For an unbiased estimate we need to divide by  $n - K$ .

## Linear Discriminant Analysis

- Expanding the term  $(x - \mu_k)^\top \Sigma^{-1}(x - \mu_k)$ ,

$$\begin{aligned} \log \pi_k g_k(x) &= c + \log \pi_k - \frac{1}{2} (\mu_k^\top \Sigma^{-1} \mu_k - 2\mu_k^\top \Sigma^{-1} x + x^\top \Sigma^{-1} x) \\ &= c' + \log \pi_k - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \mu_k^\top \Sigma^{-1} x \end{aligned}$$

- Setting  $a_k = \log(\pi_k) - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k$  and  $b_k = \Sigma^{-1} \mu_k$ , we obtain

$$\log \pi_k g_k(x) = c' + a_k + b_k^\top x$$

i.e. a **linear** discriminant function in  $x$ .

- Consider choosing class  $k$  over  $k'$ :

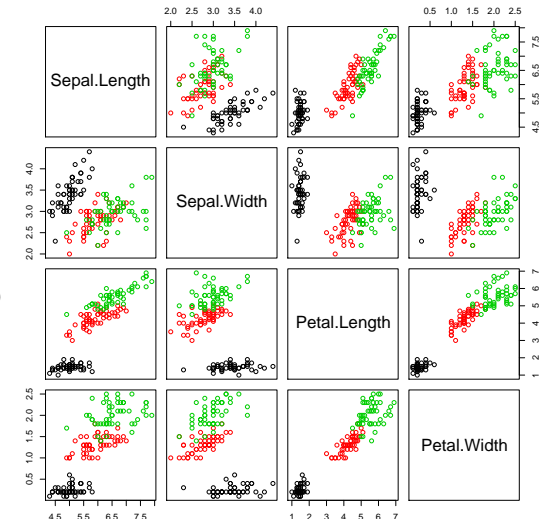
$$a_k + b_k^\top x > a_{k'} + b_{k'}^\top x \quad \Leftrightarrow \quad a_\star + b_\star^\top x > 0$$

where  $a_\star = a_k - a_{k'}$  and  $b_\star = b_k - b_{k'}$ .

- The Bayes classifier thus partitions  $\mathcal{X}$  into regions with the same class predictions via **separating hyperplanes**.
- The Bayes classifier under these assumptions is more commonly known as the **LDA classifier**.

## Iris Dataset

```
library(MASS)
data(iris)
##save class labels
ct <- unclass(iris$Species)
##pairwise plot
pairs(iris[,1:4], col=ct)
```

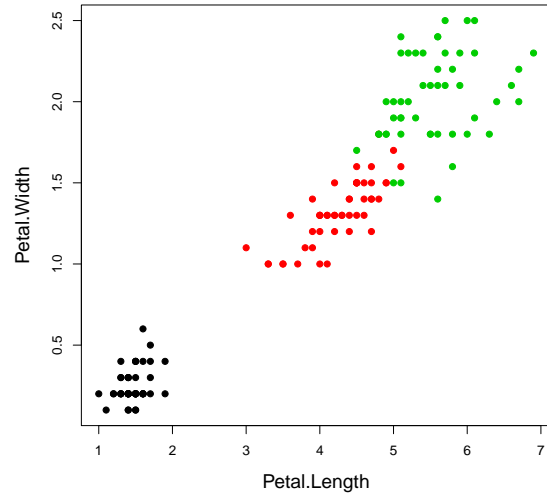


- Note: the ML estimate of  $\Sigma$  is biased. For an unbiased estimate we need to divide by  $n - K$ .

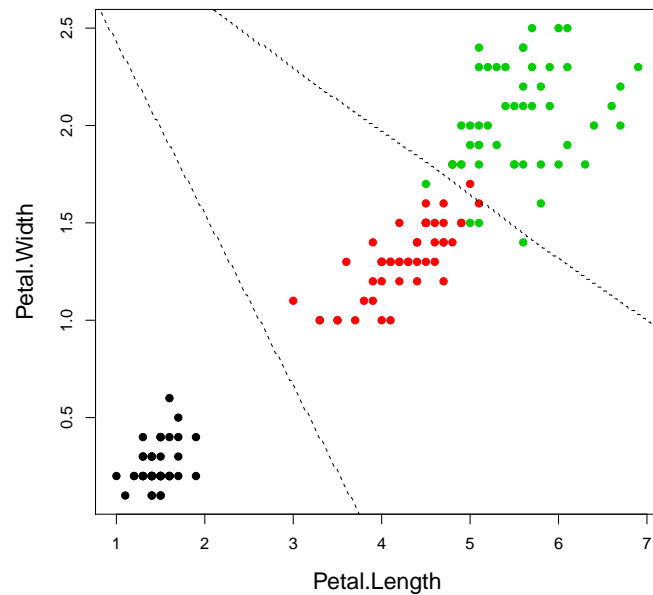
## Iris Dataset

Just focus on two predictor variables.

```
iris.data <- iris[,3:4]
plot(iris.data,col=ct,pch=20,cex=1.5,cex.lab=1.4)
```



## Iris Dataset



## Iris Dataset

Computing and plotting the LDA boundaries.

```
##fit LDA
iris.lda <- lda(x=iris.data,grouping=ct)

##create a grid for our plotting surface
x <- seq(0,8,0.02)
y <- seq(0,3,0.02)
m <- length(x)
n <- length(y)
z <- as.matrix(expand.grid(x,y),0)

##classes are 1,2 and 3, so set contours at 1.5 and 2.5
iris.ldp <- predict(iris.lda,z)$class
contour(x,y,matrix(iris.ldp,m,n),
        levels=c(1.5,2.5), add=TRUE, d=FALSE, lty=2)
```