

## STATISTICAL LEARNING

- Linear Regression
- Linear Discriminant Analysis (LDA, Fisher, 1936) → QDA
- Logistic Regression (1940s)
- Generalized Linear Models (1970s, Nelder)
- Classification and Regression trees (Breiman and Friedman, 1980s)  
→ Extended to random forests (1995)
- Machine Learning  
→ Support vector machines, neural networks, clustering, principal component analysis

Individual methods, details, theory, ...

Model accuracy, model selection, ....

BOOK James et al (2013)

Hastie et al (2001) The elements of Statistical learning

ASSESSMENT 3 Homeworks  
Final exam

TUTORATO Monday at 14:30 (remotely)  
Starting on 8<sup>th</sup> March

## Ch 2 What is Statistical Learning?

Example Wage data: wage and other variables of 3000 workers in US  
Objective Predict wage from a number of factors (age, education, ...)

WAGE DATASET : Statistical learning problem

We could formalize this problem:

$y$  : response / dependent variable / outcome

$x = (x_1, \dots, x_p)$  : p predictors / features / covariates / independent variables

In general:

$$y = \underbrace{f(x)}_{\substack{\text{fixed} \\ \text{but unknown} \\ \text{function}}} + \underbrace{\epsilon}_{\text{error (random, stochastic)}} \quad E[\epsilon] = 0 \quad \epsilon \perp x$$

(systematic / deterministic)

Statistical Learning : a set of approaches to estimate  $f$ !

## ② Why estimating $f$ ?

Two main reasons

PREDICTION

INFERENCE

### PREDICTION

Predict  $y$  when you only observe  $x$

Since  $\varepsilon$  has average zero,

$$\hat{y} = \hat{f}(x)$$

with  $\hat{f}$  our estimate of  $f$

In this context,  $\hat{f}$  can be a black-box.

The accuracy of  $\hat{y}$  as a prediction of  $y$  can be evaluated by:

$$E[(y - \hat{f}(x))^2 | x=x] = E[(f(x) + \varepsilon - \hat{f}(x))^2] =$$

$\uparrow$   
fixed

$$\begin{aligned}
 &= E[(f(x) + \epsilon - \hat{f}(x))^2] \\
 &= E[(f(x) - \hat{f}(x))^2] + \underbrace{E[\epsilon^2]}_{\substack{\text{“}E[\epsilon^2] - (E[\epsilon])^2 \\ 0}} + 2E((f(x) - \hat{f}(x))\epsilon) \\
 &= (f(x) - \hat{f}(x))^2 + \text{Var}(\epsilon) + 2(f(x) - \hat{f}(x)) \cdot \underbrace{E(\epsilon)}_0 \\
 &= \underbrace{(f(x) - \hat{f}(x))_+^2}_{\text{reducible error}} + \underbrace{\text{Var}(\epsilon)}_{\text{irreducible error}}
 \end{aligned}$$

## INFERENCE

Understanding how  $y$  changes in terms of  $X$

- ① Which predictors are most important?
- ② What is the relationship between  $y$  and  $X_j$ ?
- ③ Is the relationship between  $y$  and  $X$  linear or non-linear?
- ④ How to estimate  $f$ ?

Given a training data  $(\underline{x}_i, y_i)$ ,  $i=1, \dots, n$  with

$\underline{x}_i = (x_{i1}, \dots, x_{ip})^t$  vector of observations for unit  $i$

$y_i$  response

Different ways :

① Parametric methods

We make an assumption about the functional form of  $f$ :

$$f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

Fit & train the model on training data

→ parameter estimation

$$\hat{\beta}_0, \dots, \hat{\beta}_p \rightarrow \hat{f}(x)$$

Disadvantage: True  $f$  could be far from our parametric assumption

Remark:

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_1^4$$

Polynomial  
Regression

Linearity in the parameters!

## ② Non-parametric Methods

- No explicit assumption about  $f(x)$
- Estimating  $f$  by getting as close as possible to the data, but without being too rough or wiggly



why would we choose the more restrictive models?

- Interpretability (inference)
- generalizability outside of training data (prediction)

#### ④ Bias-Variance Trade off

Let  $(x_0, y_0)$  be a new observation from the same population as our training data.

We are interested in a model  $\hat{f}$  such that  $\hat{f}(x_0)$  is close to  $y_0$ .

Too complex models tend to overfit the training data, not doing well outside of the training data.

Mathematically :

$$E_{x,y}[(y_0 - \hat{f}(x_0))^2]$$

EXPECTED ERROR UNDER  
REPEATED TRAINING SETS

$$\begin{aligned}
 & E[(y_0 - \hat{f}(x_0))^2] \\
 &= E[(f(x_0) + \varepsilon - \hat{f}(x_0))^2] \\
 &= E\left[\left(f(x_0) + \varepsilon - \hat{f}(x_0) + E(\hat{f}(x_0)) - E(\hat{f}(x_0))\right)^2\right] \\
 &= E\left((f(x_0) - E(\hat{f}(x_0)))^2\right) + E(\varepsilon^2) + E[(E(\hat{f}(x_0)) - \hat{f}(x_0))^2] \\
 &\quad \boxed{+ 2 \left[ (f(x_0) - E(\hat{f}(x_0))) \varepsilon \right]} \quad \leftarrow (f(x_0) - E(\hat{f}(x_0))) \underbrace{E[\varepsilon]}_{=0} = 0 \\
 &\quad \boxed{+ 2 \left[ (E(\hat{f}(x_0)) - \hat{f}(x_0)) \varepsilon \right]} \quad \leftarrow \varepsilon \perp x \quad E(E(\hat{f}(x_0)) - \hat{f}(x_0)) \underbrace{E[\varepsilon]}_{=0} = 0 \\
 &\quad \boxed{+ 2 \left[ (f(x_0) - E(\hat{f}(x_0))) \cdot (E(\hat{f}(x_0)) - \hat{f}(x_0)) \right]} \\
 &\quad \quad \quad \underbrace{\text{constant}}_{\text{bias}} \quad \underbrace{E[E(\hat{f}) - \hat{f}]}_{= E(\hat{f}) - E(f)} = E(\hat{f}) - E(f) = 0 \\
 &= \underbrace{[f(x_0) - E(\hat{f}(x_0))]^2}_{\text{bias}} + \underbrace{\text{Var}(\varepsilon)}_{\text{irreducible error}} + \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{variance}}
 \end{aligned}$$

## BIA~~S~~-VARIANCE TRADE OFF

Simple model      ↗ high bias  
                        ↘ low variance

Complex model      ↗ low bias  
                        ↘ high variance



Good model : good trade-off between bias and variance  
→ cross-validation of error

## Supervised vs Unsupervised Learning

Supervised :  $(x, y)$  training data

Unsupervised :  $x$  observed

Supervised Methods :

① Regression methods  $\longleftrightarrow$  quantitative response (eg wage)

② Classification methods  $\longleftrightarrow$  qualitative response  
(categorical)

Objective : Given  $(x_i, y_i)$  training data, find a rule  
 $\hat{y} \leftarrow$  <sup>regression function</sup> ①  
classifier ②  
that allows to predict  $y$  from  $x$ .