# Pattern Analysis and Machine Intelligence
Statistical Learning

# Statistical Learning Outline

o What Is Statistical Learning?

- Why estimate f?

- How do we estimate f?

- The trade-off between prediction accuracy & model interpretability

$$X \rightarrow \boxed{f} \rightarrow Y/G$$

o Some important taxonomies (I expect you'll know this by heart!)

- Prediction vs. Inference

- Parametric vs. Non Parametric models

- Regression vs. Classification problems
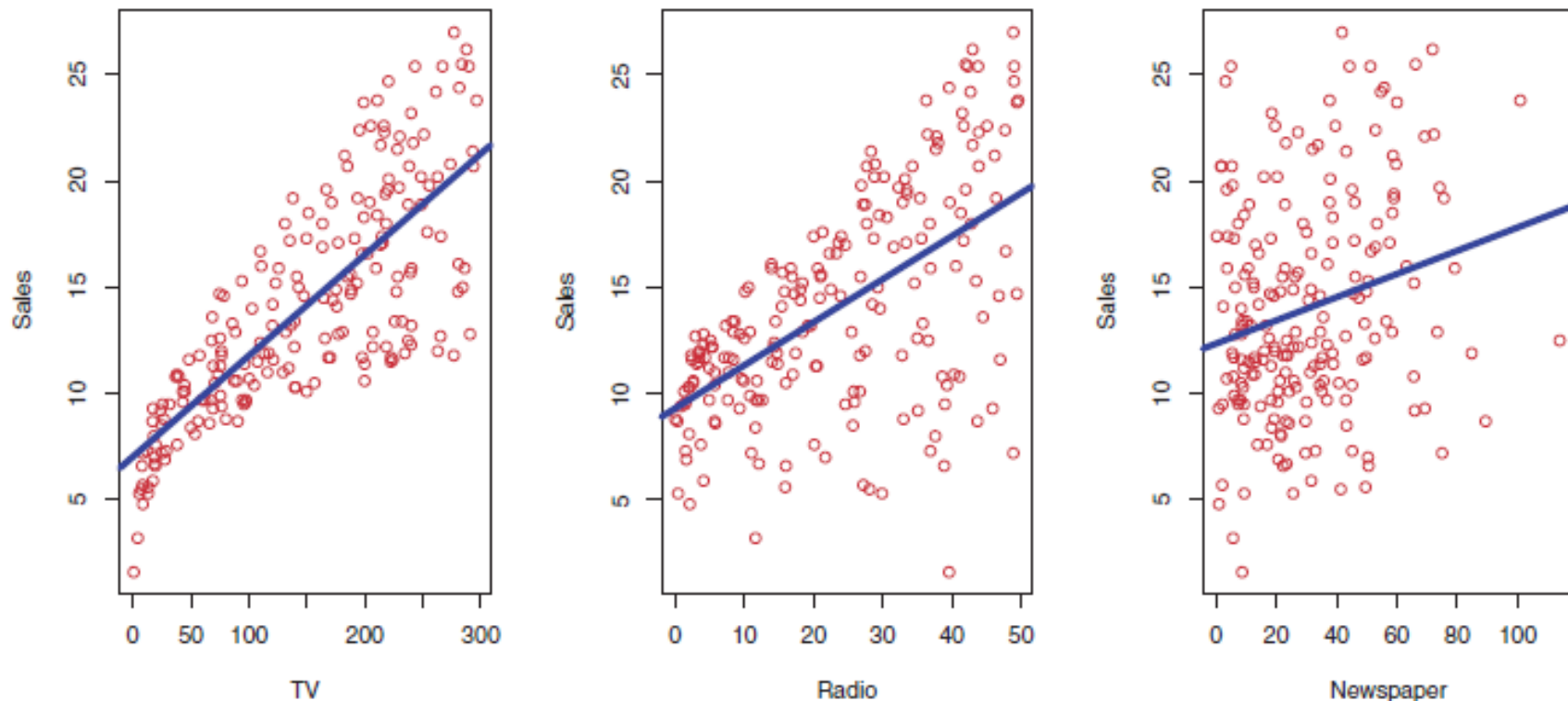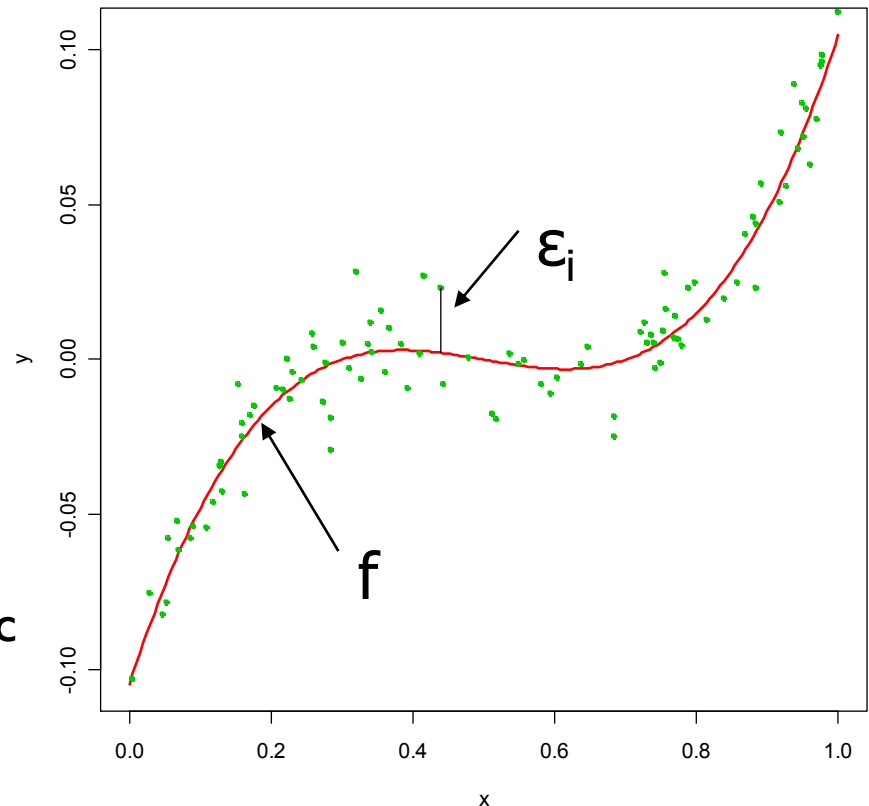
- Supervised vs. Unsupervised learning

- …

FIGURE 2.1. *The* Advertising *data set. The plot displays* sales, *in thousands of units, as a function of* TV, radio, *and* newspaper *budgets, in thousands of dollars, for 200 different markets. In each plot we show the simple least squares fit of* sales *to that variable, as described in Chapter 3. In other words, each blue line represents a simple model that can be used to predict* sales *using* TV, radio, *and* newspaper, *respectively.*

○ Suppose we observe $Y_i$ and $X_i = (X_{i1}, ..., X_{ip})$ for $i = 1, ..., n$

- Assume a relationship exists between Y and at least one of the observed X's
- Assume we can model this relationship as

$$Y_i = f(\mathbf{X}_i) + \varepsilon_i$$

- $f$ : unknown function systematic
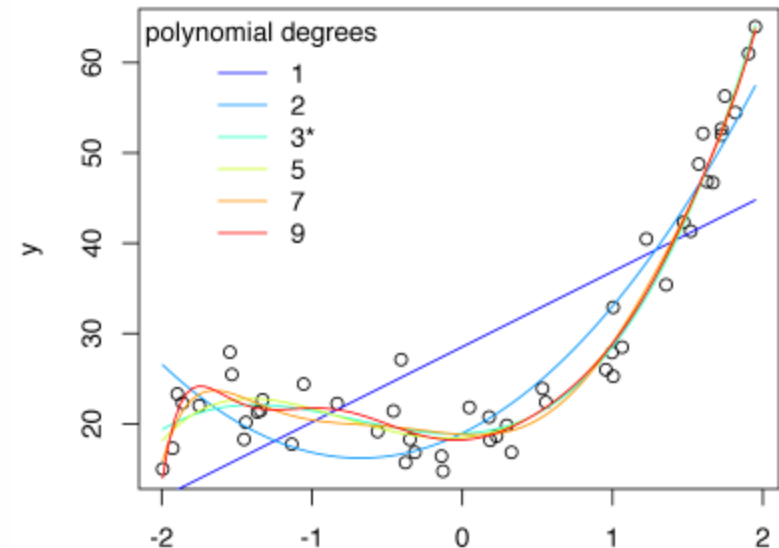- $\varepsilon_i$ : zero mean random error



○ The term statistical learning refers to using the data to "learn" $f$

o The error our estimate will have has two components

$$Y_i = f(\mathbf{X}_i) + \varepsilon_i$$

- *Reducible error* due to the choice of *f (model complexity)*



**X** → [gears] → **Y/G**

- *Irreducible error* due to the presence of $\varepsilon_i$ in the training set

# Because noise matters …



sd=0.001

sd=0.005

sd=0.01

sd=0.03

POLITECNICO DI MILANO

o The error our estimate will have has two components
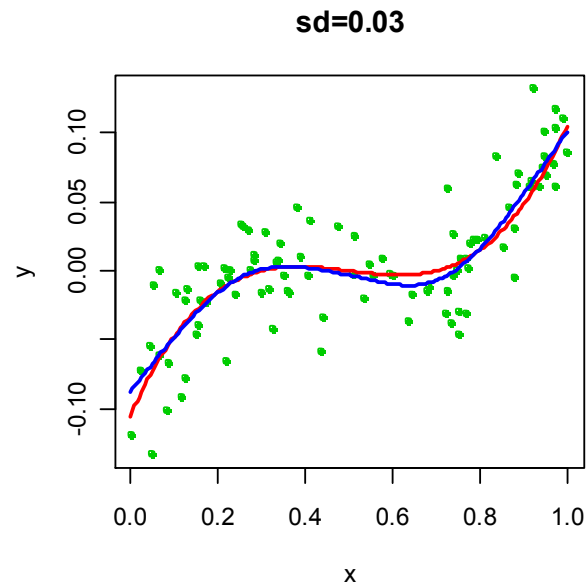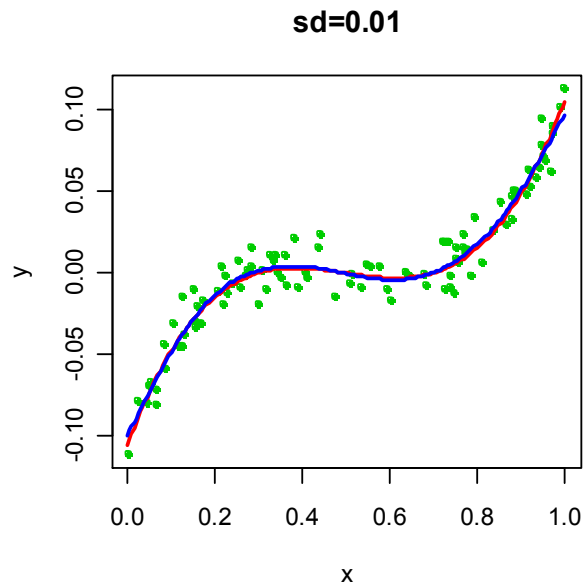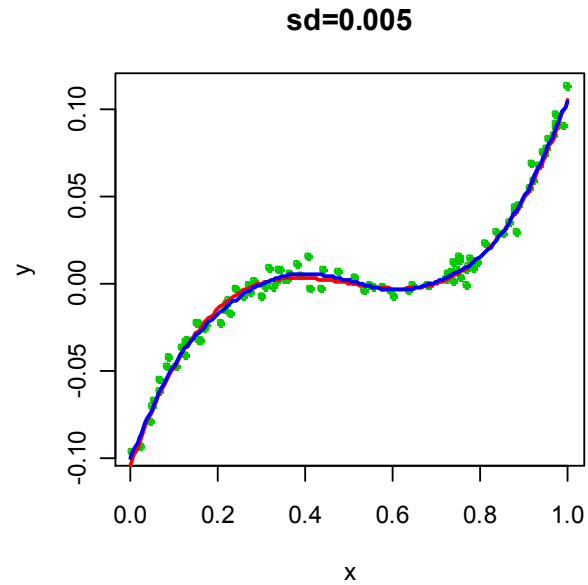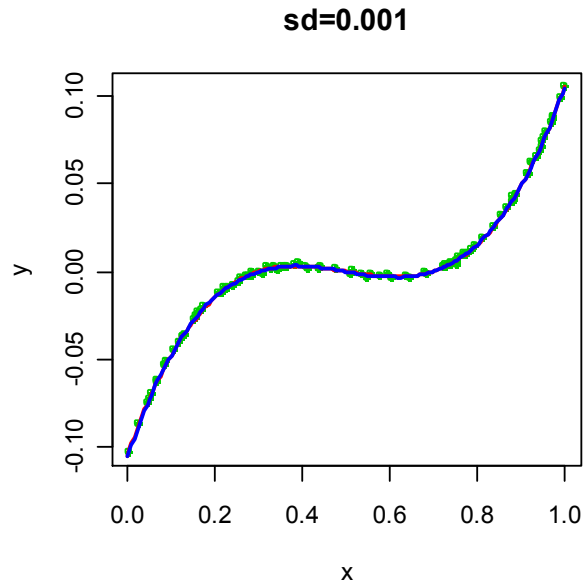
$$Y_i = f(\mathbf{X}_i) + \varepsilon_i$$

- *Reducible error* due to the choice of *f (model complexity)*
- *Irreducible error* due to the presence of $\varepsilon_i$ in the training set
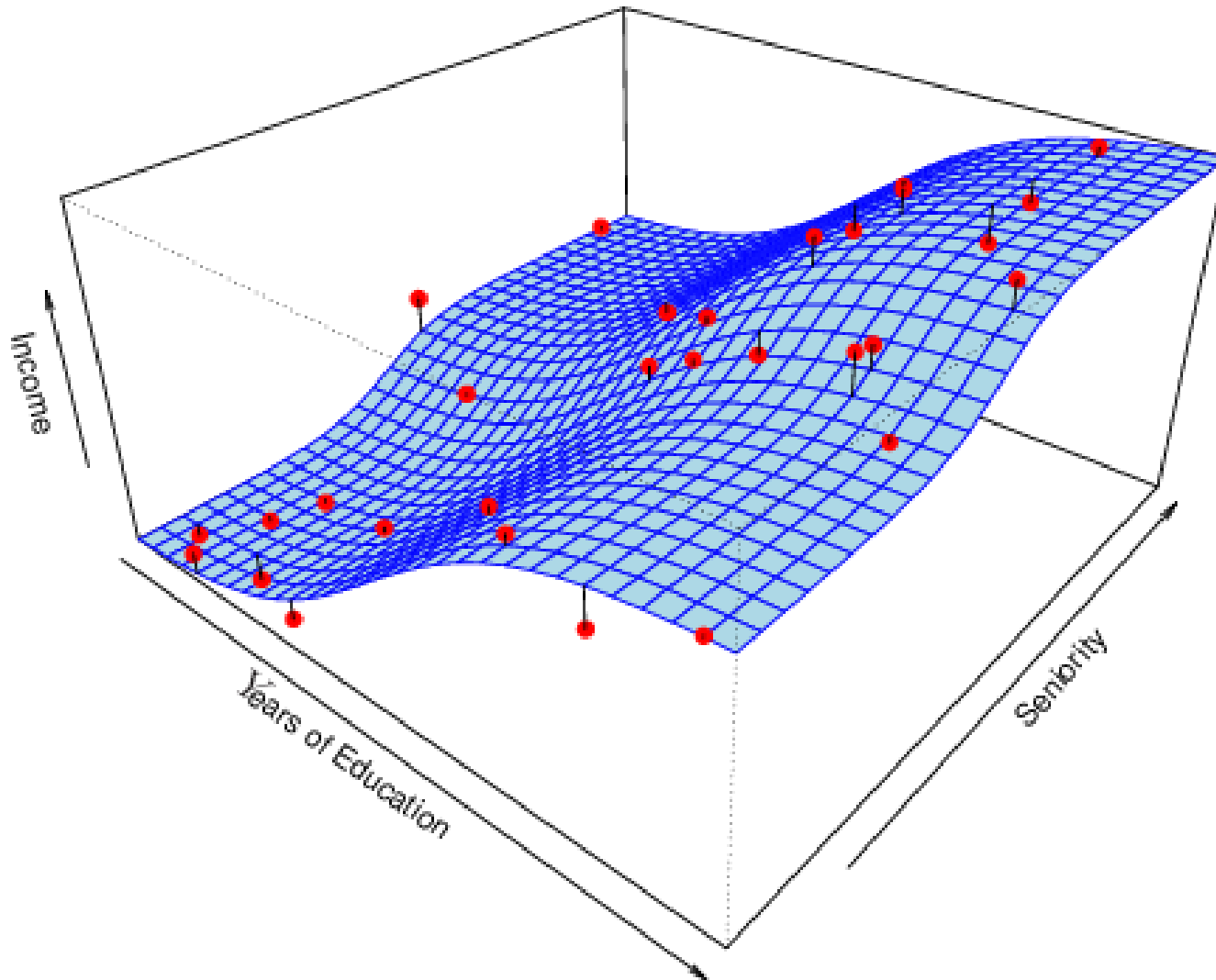
o Let assume $\hat{f}$ and $\mathbf{X}$ fixed for the time being

$$\hat{Y} = \hat{f}(X)$$

$$E(Y - \hat{Y})^2 = E[f(X) + \epsilon - \hat{f}(X)]^2$$

$$= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

o Can you derive this?

$$E(Y - \hat{Y})^2 = E[f(X) + \epsilon - \hat{f}(X)]^2$$

$$= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$



$Y_1 = f(X) + \varepsilon_1$

$Y = f(X)$

$\hat{Y} = \hat{f}(X)$

$Y_2 = f(X) + \varepsilon_2$

$X$

$$E[(Y - \hat{Y})^2] =$$

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

$$= E[f(X)^2 + \varepsilon^2 + \hat{f}(X)^2 - 2 \cdot \varepsilon \cdot f(X)$$
$$\quad - 2 \cdot \varepsilon \cdot \hat{f}(X) - 2 \cdot f(X) \cdot \hat{f}(X)] =$$

$$= f(X)^2 + E[\varepsilon^2] + \hat{f}(X)^2 - 2 \cdot E[\varepsilon] \cdot f(X)$$
$$\quad - 2 \cdot E[\varepsilon] \cdot \hat{f}(X) - 2 \cdot f(X) \cdot \hat{f}(X) =$$

$$= f(X)^2 + E[\varepsilon^2] + \hat{f}(X)^2 - 2 \cdot f(X) \cdot \hat{f}(X) =$$

$$= (f(X)^2 + \hat{f}(X)^2 - 2 \cdot f(X) \cdot \hat{f}(X)) + E[\varepsilon^2] =$$

$$= (f(X) - \hat{f}(X))^2 + E[\varepsilon^2] - 0 =$$

$$= (f(X) - \hat{f}(X))^2 + Var(\varepsilon)$$

○ Function *f* might also involve multiple variables …

# Why do we estimate *f* ?

o There are 2 reasons for estimating *f*

- Prediction
- Inference

X → [ ⚙ ] → Y/G

o Prediction

- If we can produce a good estimate for *f* (and the variance of ε is not too large) we can make accurate predictions for the response, **Y/G**, based on a new value of **X**.

o Inference

- We may be interested in the type of relationship between **Y/G** and the **X**'s to control/influence **Y/G**.
  - Which particular predictors actually affect the response?
  - Is the relationship positive or negative?
  - Is the relationship a simple linear one or is it more complicated etc.?

# Examples for Prediction & Inference

o Direct Mail Prediction

- ▪ Interested in predicting how much money an individual will donate based on observations from 90,000 people on which we have recorded over 400 different characteristics.

- ▪ Don't care too much about each individual characteristic.

- ▪ Just want to know: For a given individual should I send out a mailing?
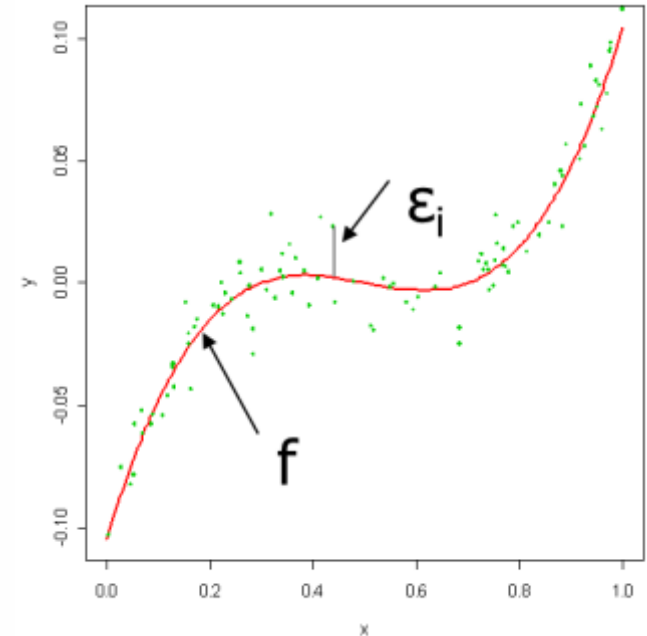
o Medium House Price

- ▪ Which factors have the biggest effect on the response

- ▪ How big the effect is.

- ▪ Want to know: how much impact does a river view have on the house value

# How Do We Estimate *f*?

o We have observed a set of <u>*training data*</u>

$$\{(\mathbf{X}_1, Y_1), (\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_n, Y_n)\}$$

o Use statistical method/model to estimate *f* so that for any (X, Y)

$$Y \approx \hat{f}(X)$$



o Statistical methods/models are usually divided in
  - Parametric Methods/Models
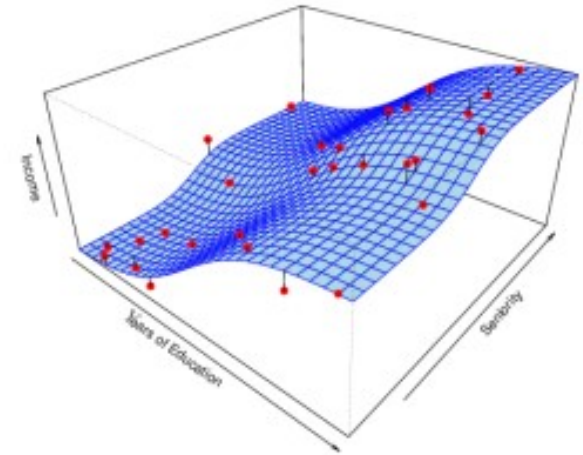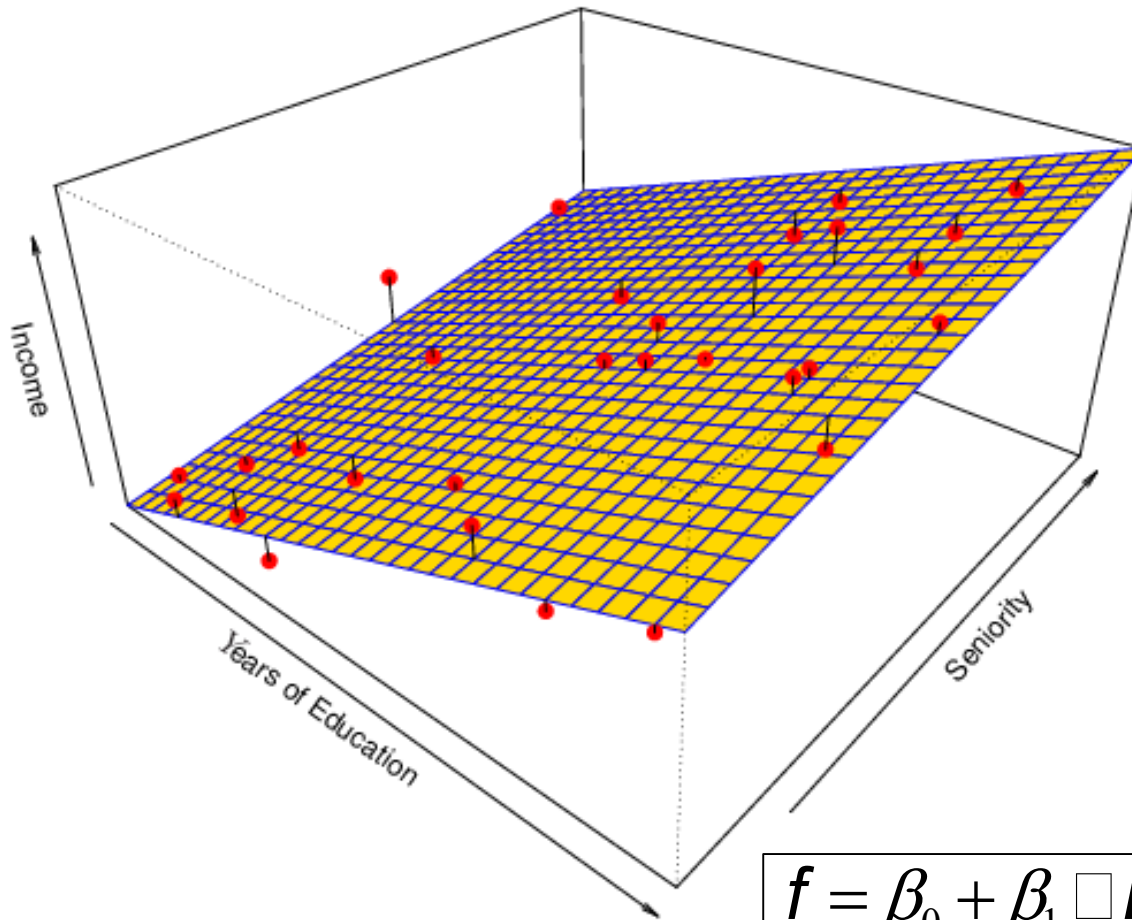  - Non-parametric Methods/Models

o Parametric methods leverage on an assumption about the model underlining *f*

  ▪ They reduce the problem of estimating *f* down to the one of estimating a set of parameters

  ▪ They involve a two-step model based approach

o STEP 1: Make some assumption about the functional form of *f*, i.e. come up with a model (e.g., a linear model)

$$f(\mathbf{X}_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}$$

o STEP 2: Use the training data to fit the model, i.e., estimate *f* through the unknown parameters

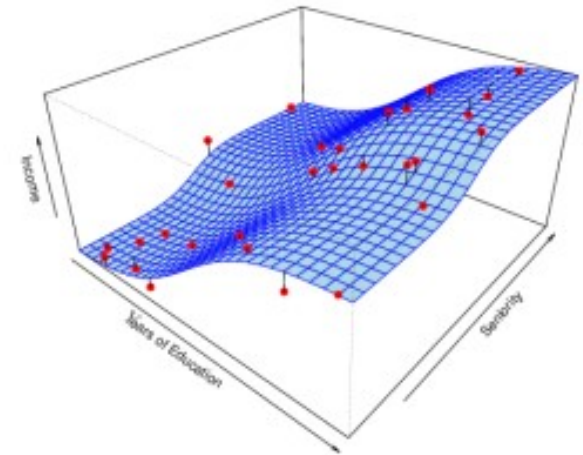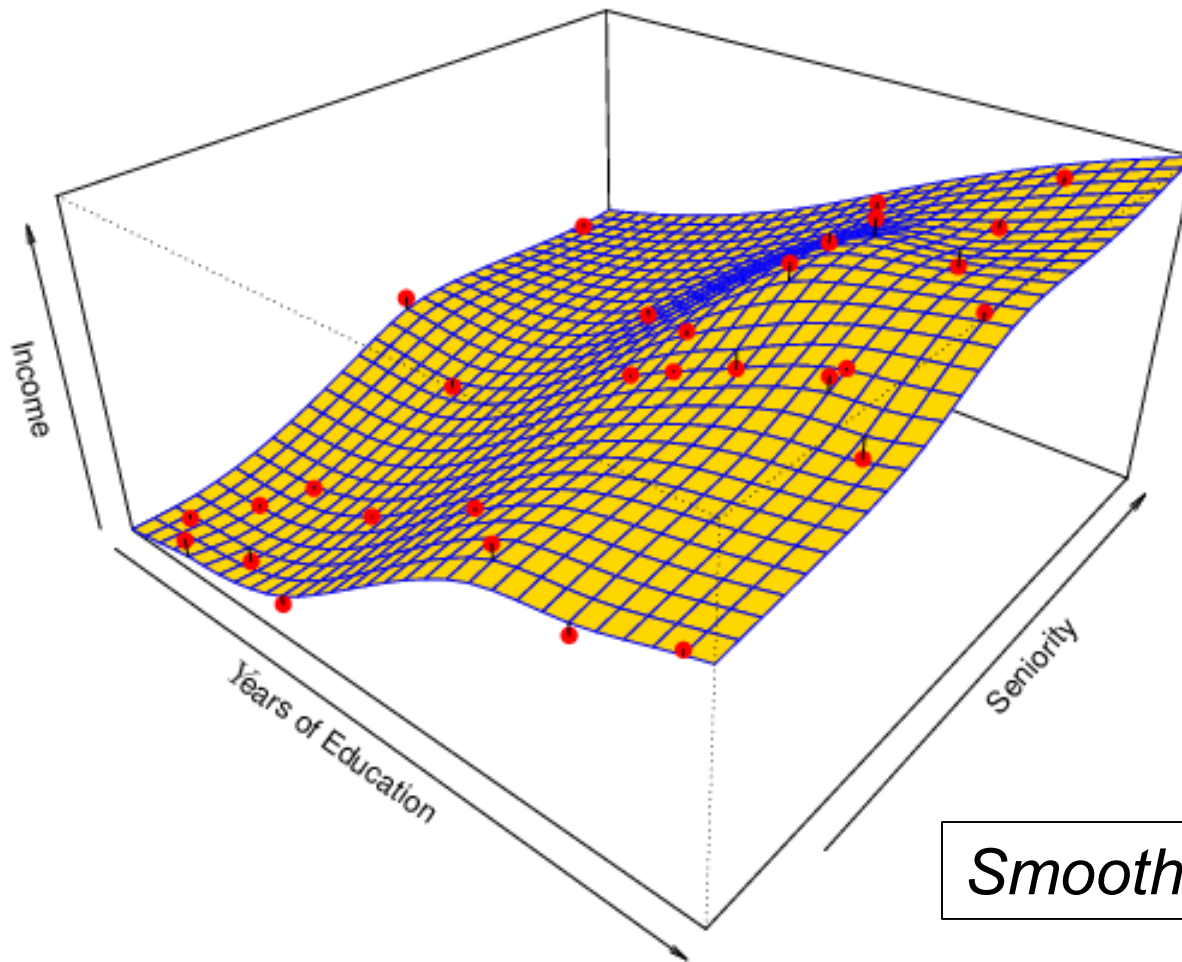$$\beta_0 \quad \beta_1 \quad \beta_2 \quad \dots \quad \beta_p$$

o Parametric methods leverage on an assumption about the model underlining $f$

- They reduce the problem of estimating $f$ down to the one of estimating a set of parameters
- They involve a two-step model based approach

o STEP 1: In this course we will examine far more complicated, and flexible, models for $f$ w.r.t linear ones. In a sense the more flexible the model the more realistic it is.

o STEP 2: The most common approach for estimating the parameters in a linear model is Ordinary Least Squares (OLS), but there are often superior approaches.

$$f = \beta_0 + \beta_1 \cdot Education + \beta_2 \cdot Seniority$$

o Even if the standard deviation is low we will still get a bad answer if we use the wrong model.

# Non-parametric Methods

o Sometimes they are referred as "sample-based" or "instance-based" methods, they do not make explicit assumptions about the functional form of f, they exploit the training data "directly"

o Advantages:
- They accurately fit a wider range of possible shapes of f
- They do not require a "trainining" phase

o Disadvantages:
- A very large number of observations is required to obtain an accurate estimate of $f$
- Higher computational cost at "testing" time
- They accurately fit a wider range of possible shapes of f.

POLITECNICO DI MILANO

# Example: A Thin-Plate Spline Estimate

Smooth thin-plate spline fit

o Non-linear regression methods are more flexible thus they can potentially provide more accurate estimates

POLITECNICO DI MILANO

o   Why not just use a more flexible method if it is more realistic?
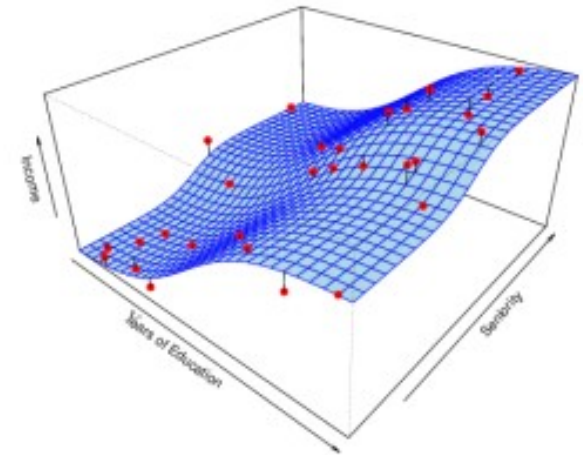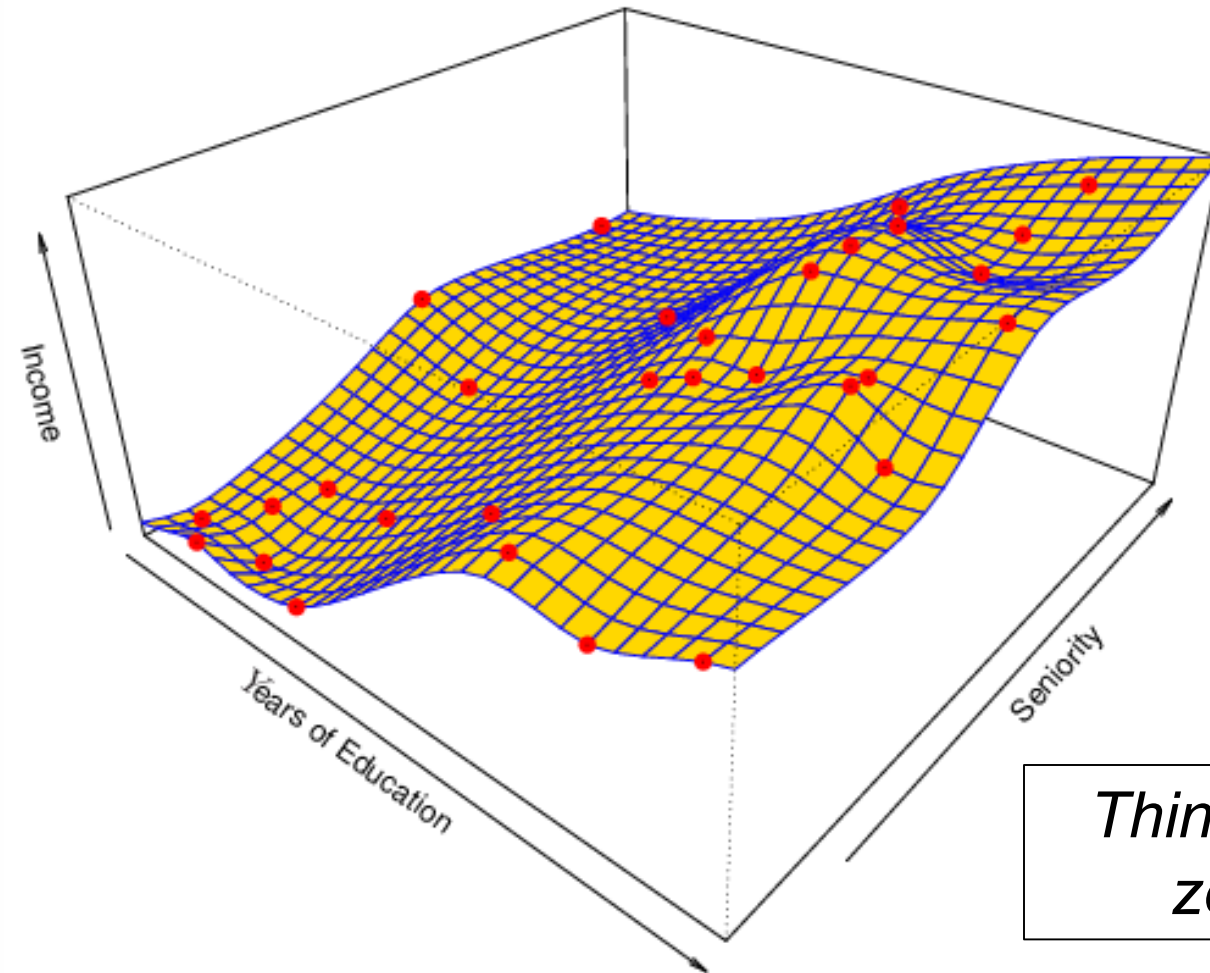
_Reason 1:_

A simple method such as linear regression produces a model which is much easier to interpret (the Inference part is better).

- E.g., in a linear model, $\beta_j$ is the average increase in Y for a one unit increase in $X_j$ holding all other variables constant.

_Reason 2:_

Even if you are only interested in prediction, it is often possible to get more accurate predictions with a simple, instead of a complicated, model.

- This seems counter intuitive but has to do with the fact that it is harder to fit properly a more flexible model.

# A Poor Estimate

*Thin-plate spline fit with zero training error*

o Non-linear regression methods can also be too flexible and produce poor estimates for *f*
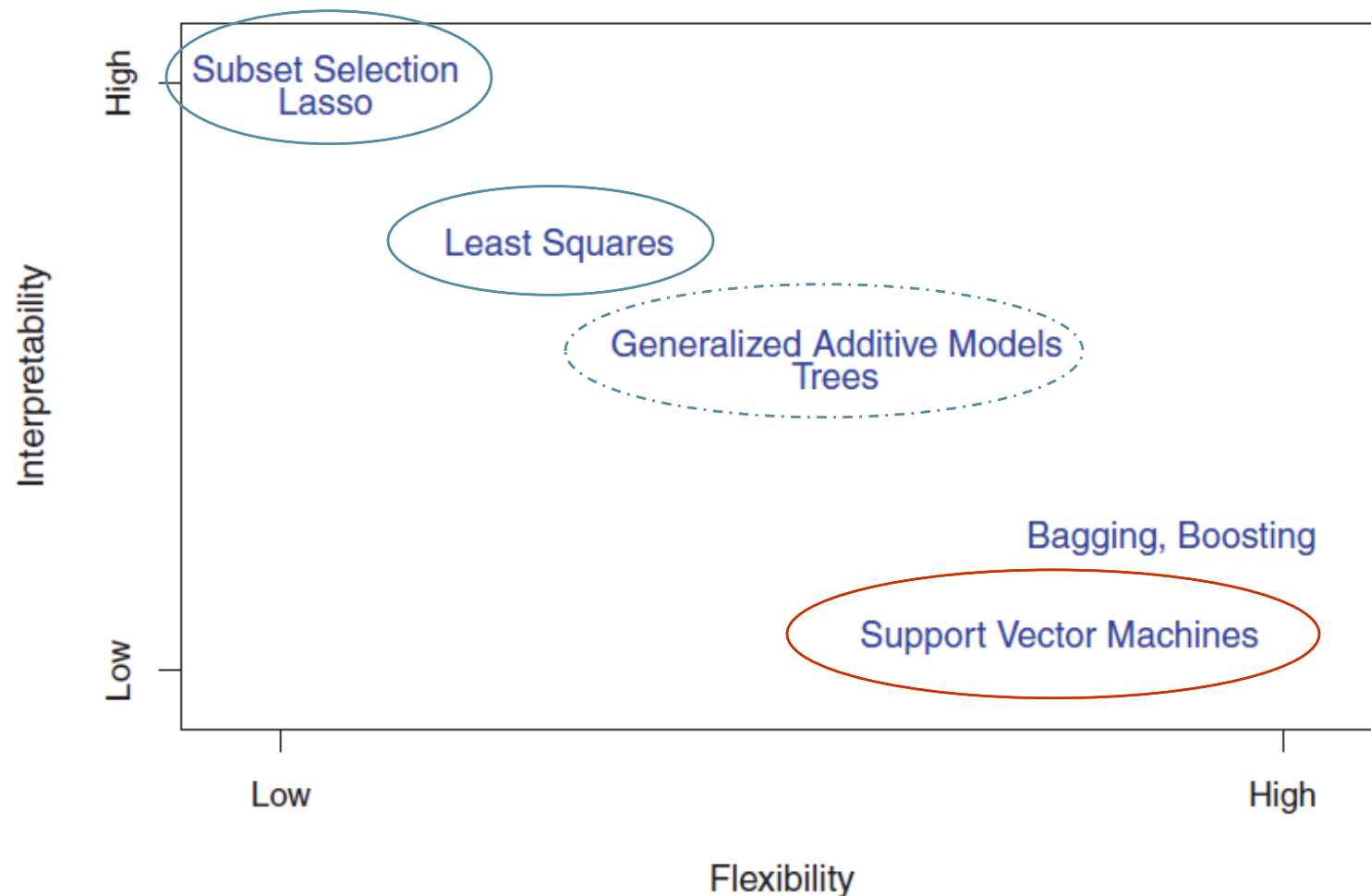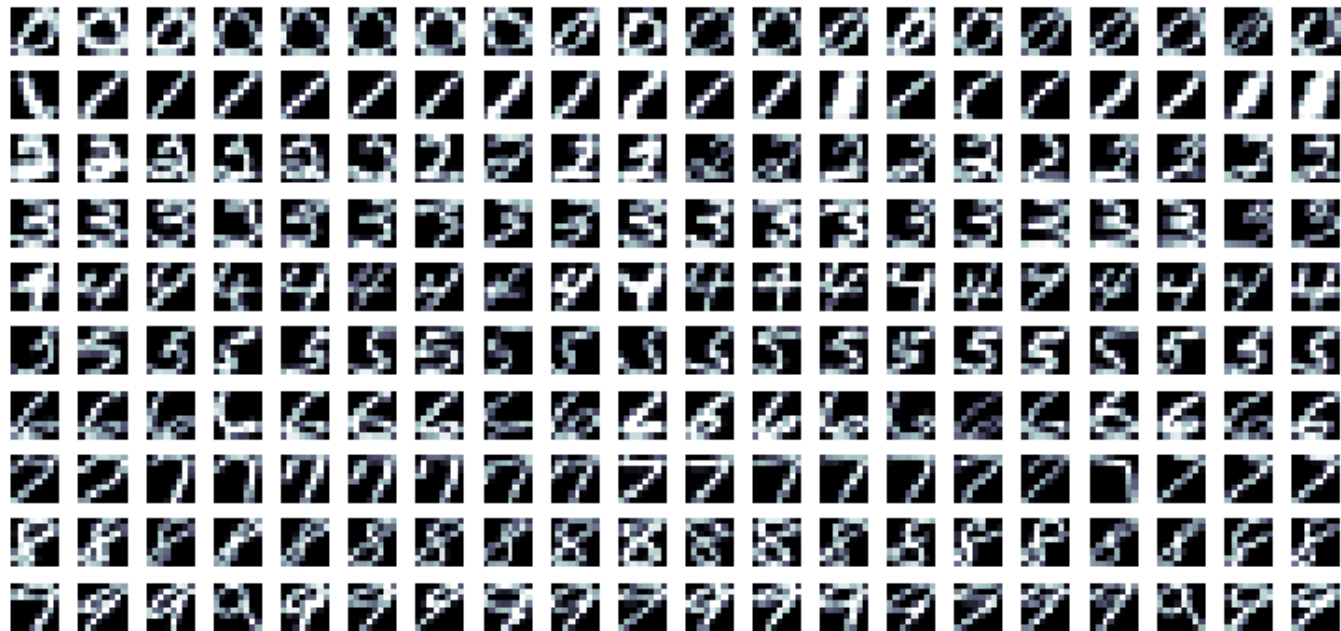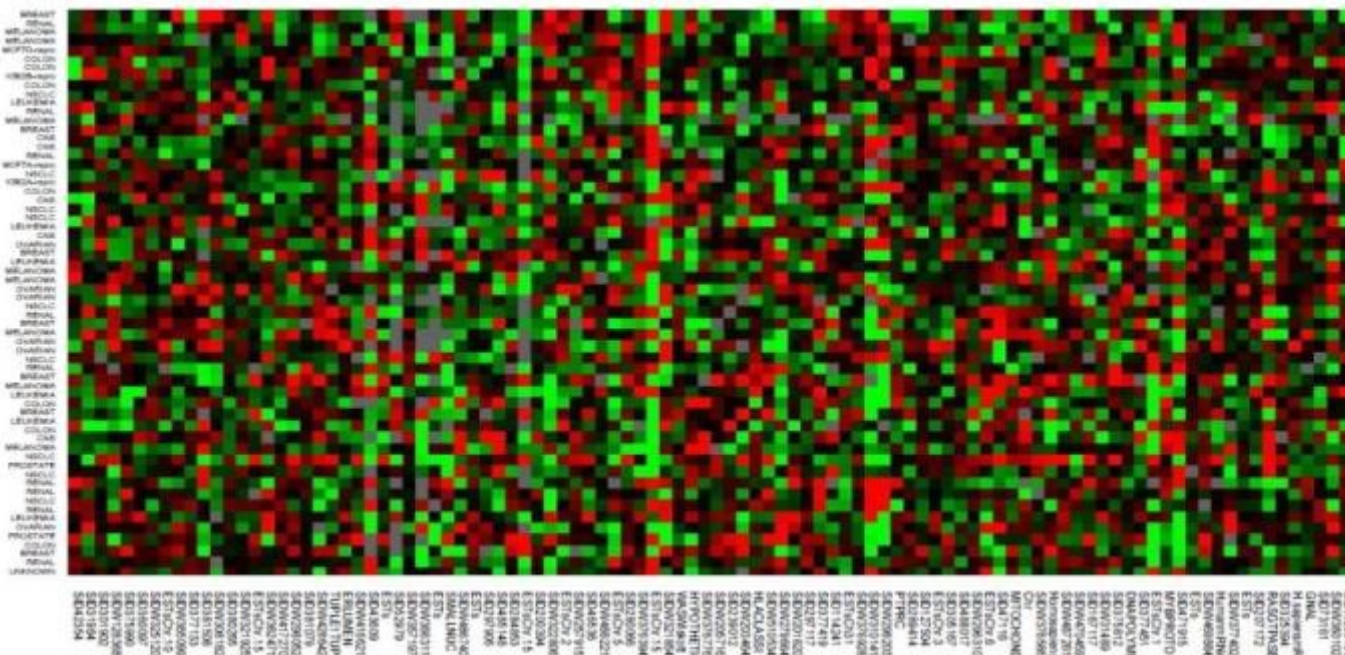
# Flexibility vs Model Interpretability

**FIGURE 2.7.** *A representation of the tradeoff between flexibility and interpretability, using different statistical learning methods. In general, as the flexibility of a method increases, its interpretability decreases.*

# Supervised vs. Unsupervised Learning

o Machine Learning makes usually a clear distinction between

- Supervised Models
- Unsupervised Models

o Supervised Learning:

- Supervised Learning is where both the predictors, $X_i$, and the response, $Y_i$, are observed.

POLITECNICO DI MILANO

# Supervised vs. Unsupervised Learning

o Machine Learning makes usually a clear distinction between

- Supervised Models
- Unsupervised Models

o Unsupervised Learning:

- Only the $X_i$'s are observed and use them to build a high level representation (possibly for modeling some Y)

○ Supervised learning problems can be further divided into

■ **Regression problems cover situations where Y is continuous/numerical**

- Predicting the value of the Dow in 6 months
- Predicting the value of a given house based on various inputs.

■ **Classification problems cover situations where Y is categorical**

- Will the Dow be up (U) or down (D) in 6 months?
- Is this email a SPAM or not?

**TABLE 1.1.** *Average percentage of words or characters in an email message equal to the indicated word or character. We have chosen the words and characters showing the largest difference between* spam *and* email.

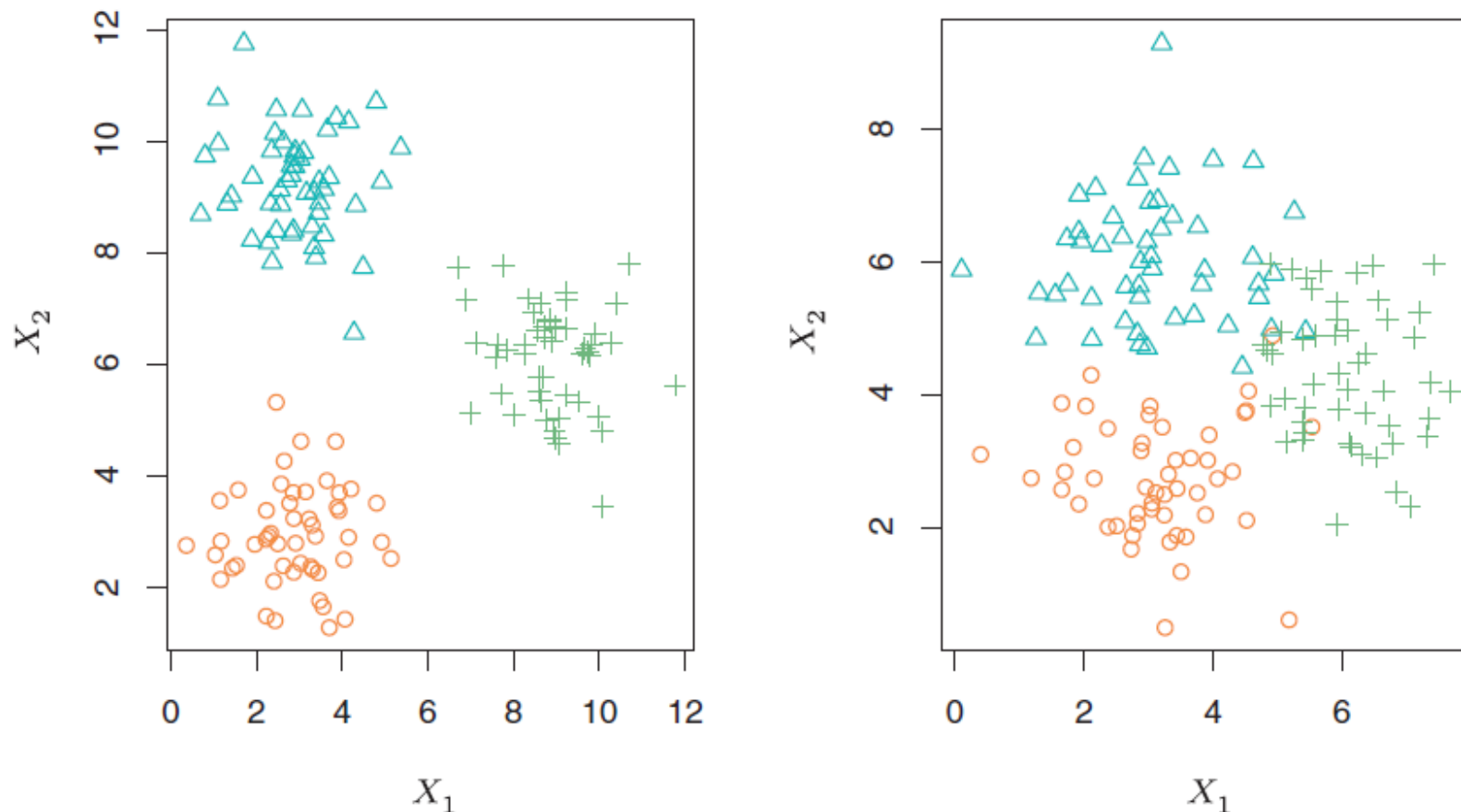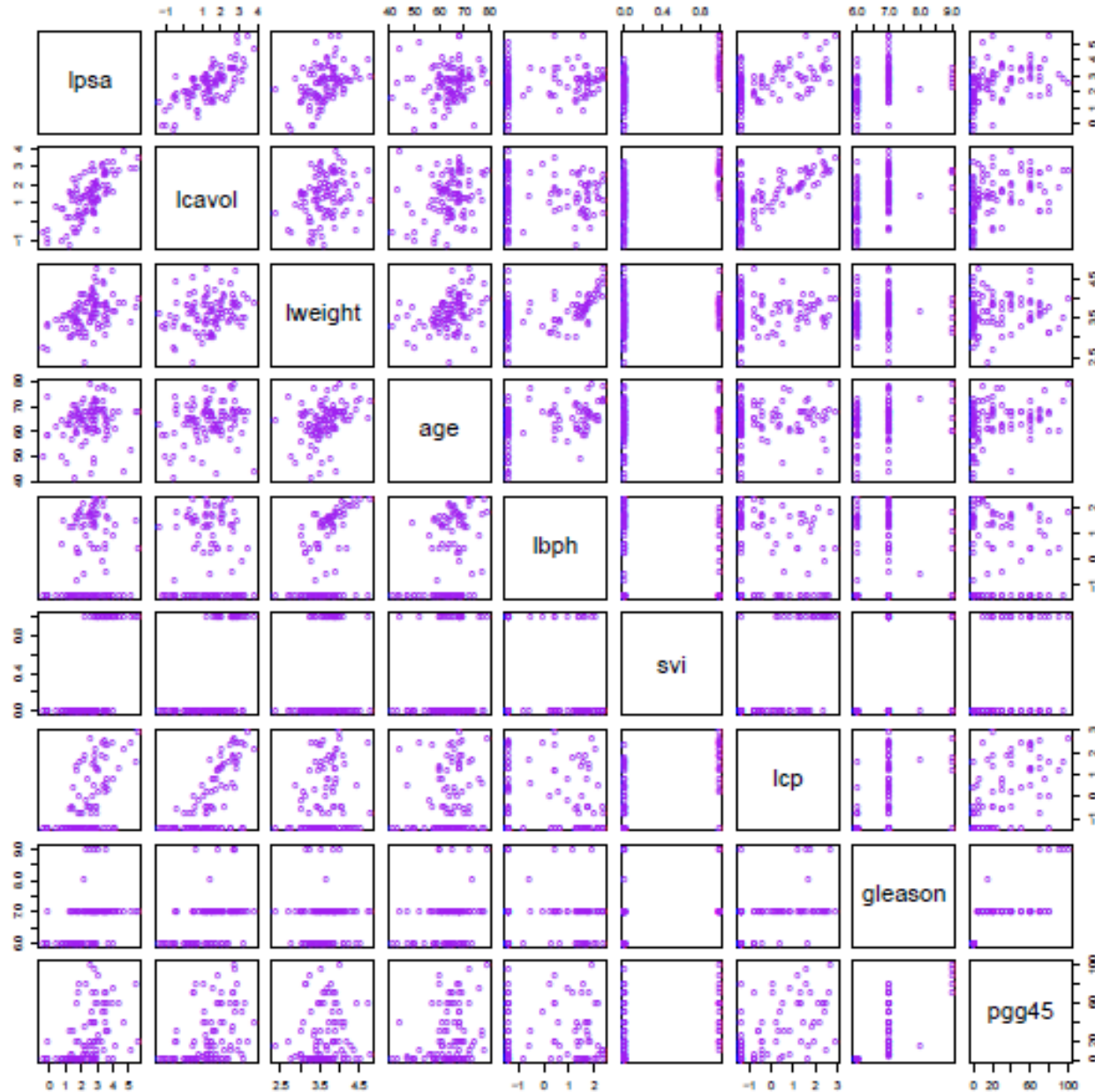|       | george | you  | your | hp   | free | hpl  | !    | our  | re   | edu  | remove |
|-------|--------|------|------|------|------|------|------|------|------|------|--------|
| spam  | 0.00   | 2.26 | 1.38 | 0.02 | 0.52 | 0.01 | 0.51 | 0.51 | 0.13 | 0.01 | 0.28   |
| email | 1.27   | 1.27 | 0.44 | 0.90 | 0.07 | 0.43 | 0.11 | 0.18 | 0.42 | 0.29 | 0.01   |

# A Simple Clustering Example

**FIGURE 2.8.** *A clustering data set involving three groups. Each group is shown using a different colored symbol. Left: The three groups are well-separated. In this setting, a clustering approach should successfully identify the three groups. Right: There is some overlap among the groups. Now the clustering task is more challenging.*

# What about higher dimensions?

POLITECNICO DI MILANO

# Wrap up!

- What Is Statistical Learning?
  - Why estimate f?
  - How do we estimate f?
  - The trade-off between prediction accuracy & model interpretability

$$X \longrightarrow \boxed{f} \longrightarrow Y/G$$

- Some important taxonomies (I expect you'll know this by heart!)
  - Prediction vs. Inference
  - Parametric vs. Non Parametric models
  - Regression vs. Classification problems
  - Supervised vs. Unsupervised learning
  - …