

## Basics of Statistical Learning

### Notations

- $X \rightarrow$  Predictors, Independent variables
- $Y \rightarrow$  Response, Dependent variables
- Input  $X \rightarrow$  Output

$$X = (X_1, X_2, \dots, X_p) \rightarrow Y$$

General form:

$$Y = f(x) + \varepsilon$$

- $f(x)$  : Fixed unknown function
- $\varepsilon$  : Error term
- We use [Statistical Learning](#) to estimate  $f(x)$

### Reasons to estimate $f(x)$

#### Prediction

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

- $\hat{f}$  estimate of  $f$
- $\hat{Y}$  result of the prediction

There are two types of errors :

1. Reducible Error  $\hat{f} \rightarrow f$
2. Irreducible Error  $Y$  is a Function of  $\varepsilon$  too

$$\begin{aligned} E(Y - \hat{Y})^2 &= E(f(x) + \varepsilon - \hat{f}(x))^2 \\ &= (f(x) - \hat{f}(x)) + \text{Var}(\varepsilon) \end{aligned}$$

- $f(x) - \hat{f}(x)$  is the reducible Error
- $\text{Var}(\varepsilon)$  is the variance of the irreducible Error

Note: In Prediction we treat  $\hat{f}$  as a black-box, What matters more is that we get  $\hat{Y}$  as close to  $Y$  as we can, Only the prediction matters

#### Inference

Here  $\hat{f}$  cannot be treated as a black-box, We want to know its exact form because we may be interested in answering these questions :

- Which *predictors*  $X_i$  are associated with the *response*  $Y$  , The ones with the most impact on the *response* ?
- What is the relationship between the *response*  $Y$  and each *predictor*  $X_i$  ?
- Can the relation between  $Y$  and each *predictor*  $X_i$  be written as a Linear equation or is it more complex than that ?
  - Because a Linear relationship is more interpretable so we always seek that

### How do we Estimate $f$ ?

- We always assume we observe  $n$  data points
- The observation  $x_{ij}$  represents the value of the  $j$ th predictor
- The response  $y_i$  represents the value of the  $i$ th observation
- So  $\{(x_1, y_1), (x_2, y_2), \dots, (x_i, y_i)\}$  will be our [Training Data](#) with  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  , with *variables* or *fields*  $1 \rightarrow p$
- We want to find  $\hat{f}$  such that  $Y \approx \hat{f}(x)$  for any observation  $(X, Y) \rightarrow$  unseen data

There are two types of approaches :

#### Parametric Methods

This method requires two steps:

1. make assumption about the Functional Form or Shape of  $f$   
most simple approach is that  $f$  is Linear [Linear Regression](#)

$$f(x) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

$f(x)$  is so simple we only need to estimate  $\beta_0, \dots, \beta_p$

2. Training or Fitting the model to estimate  $\beta_p$  such as:

$$Y \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

#### Disadvantage:

- The model we choose usually won't match the unknown form of  $f$  if the model is too far from  $f$  our estimations will be poor
- We can get around that by choosing more flexible models that fit many forms of  $f$
- More flexible models require more parameters which can lead to [Overfitting](#)
- The more flexible model the more complex it becomes

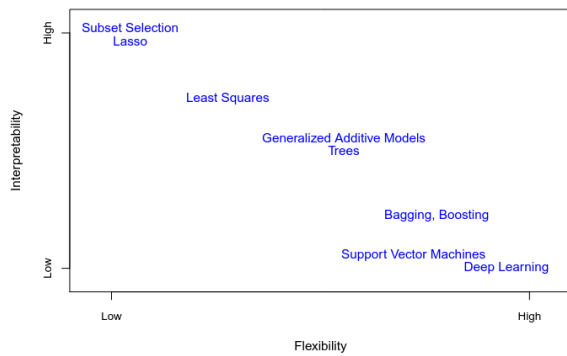
#### Non-Parametric Methods:

- Seek an estimate of  $f$  that gets as close as possible to the Data points
- No assumptions needed
- $f$  fits a wider range of possible shapes and forms
- Avoid the danger of not fitting

#### Disadvantage:

- Since the non-parametric methods don't reduce the number of parameters
- A very big number of Observations data points is needed (Way more than usual) to estimate  $f$

### Trade off : Prediction Accuracy vs Model interpretability



- Interpretability helps when inference cause its easy to understand the relationship between  $Y$  and  $X$
- When Prediction is our goal more flexible models fits more
- **Note** : Sometimes a less flexible models gives better predictions than more flexible ones, it all depends on the type of relationship between  $Y$  and  $X$

## Supervised Versus Unsupervised Learning :

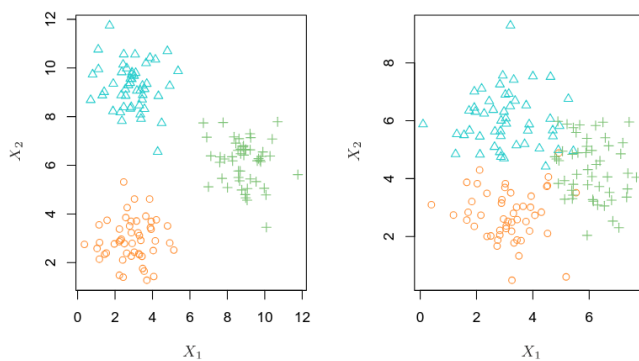
### Supervised Learning:

- For each [Observation](#)  $x_i$  there is a [Response](#), Measurement  $y_i$
- We want to fit a [model](#) that [Observation](#)  $\leftrightarrow$  [Response](#)
- With the aim to predict future unseen responses [Prediction](#)  
Or better understanding of the relationship between Observations  $x_i$  and the response  $y_i$  [inference](#)

### Unsupervised Learning :

Deals with more challenging situations

- We have the [Observation](#)  $x_i$  without the [Response](#)  $y_i$
- Seek to understand the relationships between the variables or the Observations



## Regression Vs Classification Problems :

It all depend on the type of Variables  $1 \rightarrow p$  either

- Quantitative  
Regression Problems
- Qualitative (Categorical)  
Classification Problems
- So Selecting the [Statistical Learning](#) Method depends on  $\rightarrow$  The [Response](#) variable type (quan/quali)
- The predictors  $X$  variable type are way less important on the decision