# Chapter 2: Statistical Learning.

- goal: develope an acurate model to predict.

- Observe quantitative response $Y$
  
  $p$ different predictors $X = (X_1, X_2, \dots, X_p)$.

Assumption: $\quad y = f(x) + \varepsilon \qquad \varepsilon$: independent of X.

$\downarrow$ fixed but unknown

$\longrightarrow$ random error term. $\qquad E(\varepsilon) = 0$.

- $f$: represents systematic info that X provides about Y.

- 2.1.1 why estimate $f$.
  
  ① Prediction  ② Inference.

① Prediction.

$$\hat{y} = \hat{f}(x). \qquad \hat{f}: \text{black box.}$$

$\longrightarrow$ Reducible : improve $\hat{f}$ so that it is close to $f$. $\to$ use correct statistical technique.

Accuracy of $\hat{Y}$

$\longrightarrow$ Irreducible. $y = f(x) + \boxed{\varepsilon}$.
error.  $\qquad$ cause.

even if $\hat{y} = f(x)$, then $\varepsilon$ error.

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

$\downarrow$ overall expected variation.

$$= \left[f(x) - \hat{f}(x)\right]^2 + Var(\varepsilon)$$

reducible  $\qquad$ irreducible.

② Inference. $f$ cannot be treated as black box.

1. which predictors are associated with the response?
   ↳ identify few important predictors.

2. what is relationship between the response & each predictor?

3. what is complexity of relationship between $y$ and $x_i$? is it linear or more complicated?

### 2.1.2 How Do we Estimate $f$?

$$\hat{f} \text{ s.t } y \approx \hat{f}(x)$$

Parametric            Non - Parametric

① Parametric methods:

1. Assumption about functional form /shape.
$$f(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p \longrightarrow \text{linear model.}$$

2. Procedure to fit the model. /train the model.
   → OLS to estimate $\beta_i$

   "flexible" models cause problem of overfitting

② Non - parametric methods:
   - no functional form assumption
   - very large no of observations required.

doubt: - "level of smoothness". → may lead to overfitting.

### 2.1.3 The trade off betⁿ Prediction Accuracy & model interpretability

   - if interested in inference then: restrictive model.
   - (relationship with each variable is easy to understand).

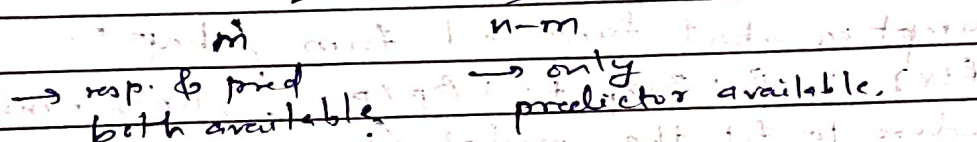## 2.1.4 Supervised Versus Unsupervised Learning.

| Supervised Learning | Unsupervised Learning. |
|---|---|
| • Response variable present. | • Response Variable absent. |
| • prediction/ inference. | • Understand relationship betn predictors / find patterns. |
| • Linear reg, logistic regres, GAM, boosting , etc . | • cluster analysis : whether problems fall into relatively distinct groups . |
| | └ all ; classified in correct group . |

* Semi-supervised learning : $n$ observ.

$$m \qquad n-m$$

→ resp. & pred both available

→ only predictor available.

### 2.1.5 Regression vs Classification.

### 2.2 Assessing model accuracy ;

1. MSE : mean Squared Error.

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{f}(x_i))^2$$

→ one calculated from test data.

→ We , choose model with low test MSE.

— The model with least training MSE may not have low test MSE .

— Fundamental property: model flexibility ↑ , training MSE ↓ but test MSE may not .

— over fitting : low train MSE , but high test MSE .

└→ We are finding too much of patterns — computing patterns due to errors — which are not present in test data.

— Cross validation : Method of estimating test MSE using training data.
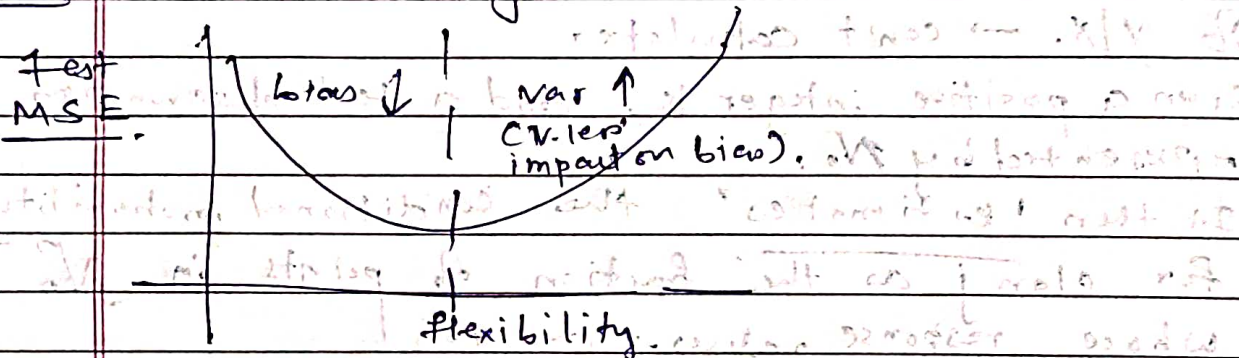
**2.2.2 : Bias Variance Trade off.**

$$E(MSE) =$$

Expected test MSE = Bias + Var + Var($\epsilon$).

$$\Rightarrow \geq Var(\epsilon).$$

— more flexible statistical methods have higher variance, but less bias.

How is flexibility measured?



Test MSE.

↓bias ↓ │ Var ↑
CV-les'
impact on bias).

flexibility.

**2.2.3 : The classification Setting.**

error rate $= \frac{1}{n} \sum_{i=1}^{n} I(y_i \neq \hat{y}_2)$ ⟶ computes ~~fraction of~~ incorrect classificat.

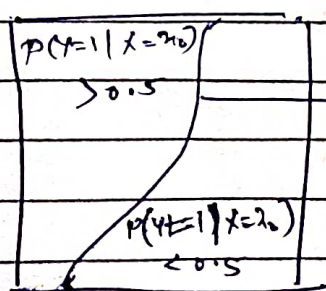test error rate $=$ Ave $(I(y_0 \neq \hat{y}_i))$.

**The Bayes Classifier.**

assigns each observation to the most likely class, given its predictor values.

$x_0 \rightarrow$ class $j$ ~~for~~ if largest $P(Y=j \mid X=x_0)$.

_Ex. :_ $y_i = \{$ Class 1, class 2 $\}$

$x_0 \rightarrow$ class 1 if $P(Y=1 \mid X=x_0) > 0.5$.



$P(Y=1 \mid X=x_0)$
$> 0.5$

⟶ Bayes decision boundary.

$P(Y=1 \mid X=x_0) =$

$P(Y=2 \mid X=x_0) = 0.5$.

$P(Y=1 \mid X=x_0)$
$< 0.5$

- The Bayes clansifier produces lowest possible test error rate.
" Bayes error rate ".

$$= 1 - E\left[\max_j \left\{ P(Y=j|X) \right\}\right]$$

~~We do not know the~~

## K Nearest Neighbors.

- Baye's clansifier: we do not know conditional distribution of $Y|X$. → can't calculate.
- Given a positive integer $k$, and a test observation represented by $N_o$.
- It then 'estimates' the conditional probability for class $j$ as the fraction of points in $N_o$ whose response values equal to $j$

$$P(Y=j | X=n_o) = \frac{1}{k} \sum_{i \in N_o} I(y_i=j) = p$$

↪ KNN clansifies the test observation $x_o$ to the class with the largest $p$.

- KNN can often produce classifiers that are surprisingly close to the optimal Bayes classifier.