

Introduction to Machine Learning

Franck JAOTOMBO

Introduction to Data Analysis

- In Data Science, before you do anything, you first start by exploring your dataset.
- The basic steps in data exploration are following:
 1. Univariate Data Analysis
 2. Bivariate Data Analysis
 3. (Modeling)

Step 1 - Univariate Data Analysis

1. If the variables are categorical.
 1. Generate the summary table for each variable.
 2. Plot their Pie Chart
 3. Plot their Bar Chart
2. If the variables are quantitative.
 1. Generate the frequency table for each variable
 2. Plot their histogram
 3. Plot their boxplot

Step 2 – Bivariate Data Analysis

1. If the variables are both categorical.
 1. Generate the contingency table
 2. Check the significance of their relationship with the chi-square test & provide Cramer's V
 3. Plot their side-by-side bar charts
 4. Plot their stacked bar charts
2. If the variables are both quantitative.
 1. Compute the correlation (table)
 2. Check the significance of their relationship with the correlation test & provide the r value
 3. Plot their scatter plot (matrix)
3. If the variables are mixed categorical & quantitative.
 1. Compute the anova table
 2. Check the significance of the difference in values between groups
 3. Plot the grouped boxplots

Modeling : supervised

- Which variables should be selected as the **outcome** (target, response, dependent variable) or variable to be explained from the others (**predictors**, features, independent variables)?
- The answer should be justified with theoretical, managerial or statistical arguments

- The goal is to **find a function** that captures in the best possible way the relationship between the outcome and the predictors
- This process is called “modeling” and statistical learning is one way of addressing it

- One goal of modeling is thus to explain the variability (or variance) in the outcome from the predictors.
- This approach to modeling is associated with “**supervised learning**” in statistical learning.

Modeling : unsupervised

- The variance may also be explained by the existence of subgroups of individual instances or subgroup heterogeneity.
- The process to account for these subgroups is associated with “**unsupervised learning**” in statistical learning.

When the number of variables is too numerous, the relationship between the outcome and the predictors can become too complex.

- To reduce this complexity and to simplify interpretability, dimension reduction is recommended.
- The set of tools to reduce dimensions is also associated with “**unsupervised learning**” in statistical learning.

Hypothesis Testing : A review

Correlation Test

$$H_0: \rho = 0$$

$$H_1: \rho \neq 0$$

$$\text{Statistic of the test} - \text{Student}(n - 2) : t = \frac{r\sqrt{n - 2}}{\sqrt{1 - r^2}}$$

Chi-square Test

$$H_0: \chi^2 = 0$$

$$H_1: \chi^2 > 0$$

$$\text{Statistic of the test} - \text{Chi square}[(r - 1)(c - 1)] : X^2 = \sum_{i,j} \frac{\left[n_{ij} - \frac{r_i c_j}{n}\right]^2}{\frac{r_i c_j}{n}}$$

Anova Test

$$H_0: \mu_1 = \mu_2 = \dots = \mu_c$$

$$H_1: \text{All the } \mu_j \text{ are not equal}$$

$$\text{Statistic of the test} - \text{Fisher}(c - 1, n - c) : F_{stat} = \frac{MSB}{MSW}$$

Linear Correlation Test (Pearson)

- We want to test if there is a significant association between two continuous variables
- The covariance between two variables X and Y indicates if there is an association between the variation of the two variables around their respective means

$$\text{cov}(X, Y) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

- Correlation is a standardized measure of the covariance

$$r = \frac{\text{cov}(X, Y)}{s_X s_Y} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}}$$

s_X and s_Y are the respective standard deviations of X and Y

Chi square

- If the row and column variables are independent: $E(n_{ij}) = \frac{r_i c_j}{n}$

$$X^2 = \sum_{i,j} \frac{\left[n_{ij} - \frac{r_i c_j}{n}\right]^2}{\frac{r_i c_j}{n}}$$

	X_1	X_j	X_c	<i>Total</i>
Y_1	n_{11}		n_{1c}	r_1
Y_i		n_{ij}		r_i
Y_l	n_{r1}		n_{rc}	r_l
<i>Total</i>	c_1	c_j	c_c	n

- r_i and c_j indicate respectively the total (marginal) frequency of row i and column j
- X^2 follows a Chi Square distribution with a degree of freedom = $(r - 1) * (c - 1)$
 - where **r = number of modalities on rows** and **c = number of modalities on columns**
 - We need only to compare X^2 with the threshold values of the Chi Square (χ^2) distribution

- Effect size : $\phi = \sqrt{\frac{X^2}{n}}$ and $V_{Cramer} = \sqrt{\frac{X^2}{n \cdot \min[(r-1), (c-1)]}}$

Analysis of Variance

$$SST = \sum_{j=1}^c \sum_{i=1}^{n_j} (X_i^j - \bar{\bar{X}})^2$$

$$MST = \frac{SST}{n - 1}$$

$$SSB = \sum_{j=1}^c n_j (\bar{X}_j - \bar{\bar{X}})^2$$

$$MSB = \frac{SSB}{c - 1}$$

$$SSW = \sum_{j=1}^c \sum_{i=1}^{n_j} (X_i^j - \bar{X}_j)^2$$

$$MSW = \frac{SSW}{n - c}$$

$$\text{Grand mean : } \bar{\bar{X}} = \sum_{j=1}^c \sum_{i=1}^{n_j} \frac{X_i^j}{n}$$

$$\text{Group mean : } \bar{X}_j = \sum_{i=1}^{n_j} \frac{X_i^j}{n_j}$$

Exercise

- **Explore the <Flourishing> dataset**
 - Read the instructions in the <Flourishing_Case.docx> document

Session 1 – Main Concepts

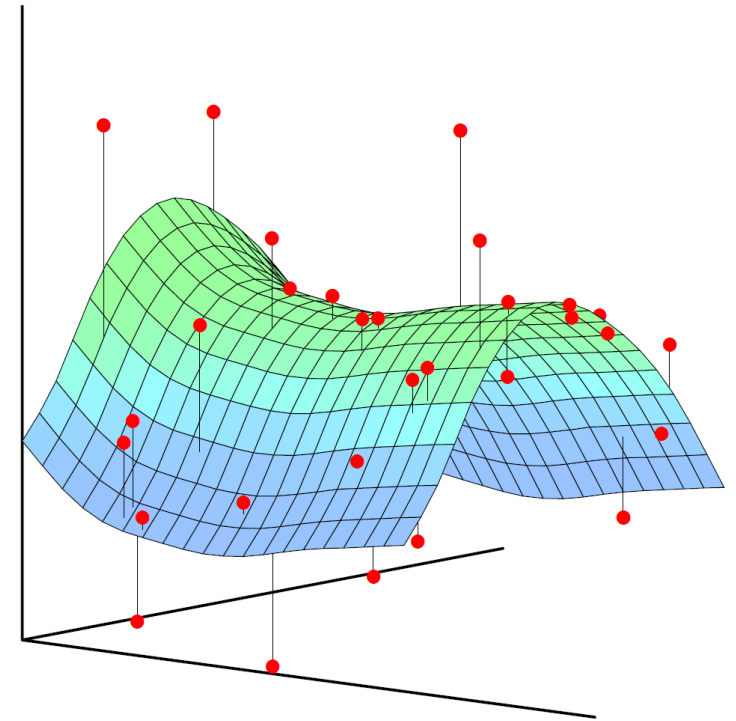
A high level overview of the main concepts used in Machine & Statistical Learning

Reference :

[James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. \(2023\). *An Introduction to Statistical Learning : With Applications in Python* \(1st ed. 2023 edition\). Springer.](#)

Statistical Learning versus Machine Learning

- **Machine learning arose as a subfield of Artificial Intelligence.**
 - Leaning more towards computer science.
 - An algorithmic approach.
- **Statistical learning arose as a subfield of Statistics.**
 - Leaning more towards mathematics and statistics.
 - A modeling approach.
- **There is much overlap** : both fields focus on supervised and unsupervised problems:
 - Machine learning has a greater emphasis on **large scale** applications and **prediction accuracy**.
 - Statistical learning emphasizes models and their interpretability, and **precision** and **uncertainty**.
- **But the distinction has become more and more blurred, and there is a great deal of “cross-fertilization”.**
- **Machine learning has the upper hand in Marketing!**



What Is Statistical Learning?

Starting point

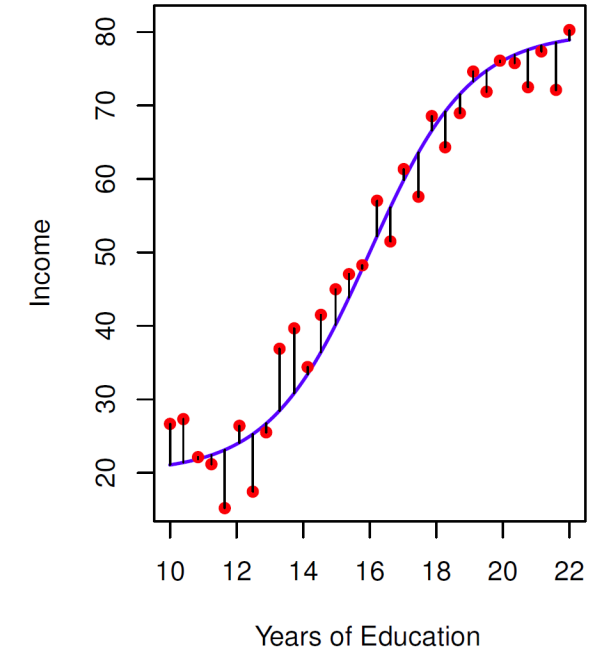
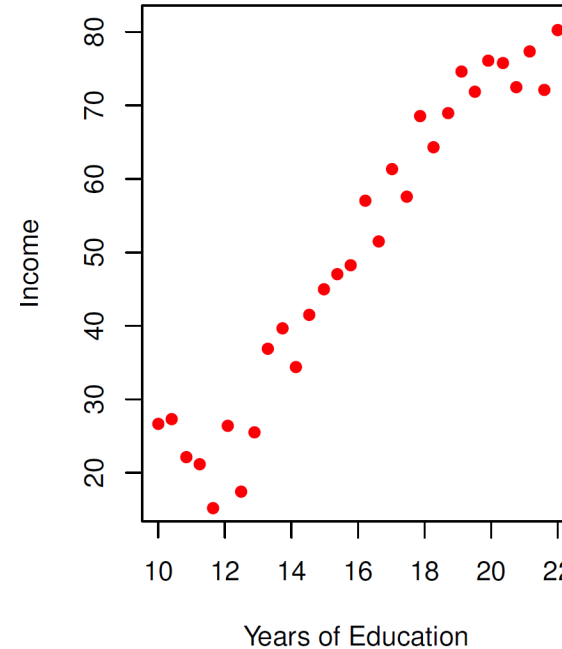
- Several **input** variables $\mathbf{X} = [X_1, X_2, \dots, X_p]$
 - Inputs :: Predictors :: Features :: Independent variables
 - Each predictor X_j has n data points
- One **output** variable Y (with n data points also)
 - Output :: Outcome :: Response :: Dependent variable
- \mathbf{X} and Y are given by the (observed) data
- Some relationship exists between X and Y

$$Y = f(\mathbf{X}) + \epsilon$$

- ϵ : a **random error** term
- f : **systematic information** \mathbf{X} provides about Y

Goal

- Estimating f from the data



In essence, statistical learning refers to a set of approaches for estimating f (James et al, 2023, p.17)

Why estimate f ?

▪ Goal : **Prediction**

- Is this newly admitted patient likely to have a prolonged stay ?
- What is the mostly likely rate of turnover in our organization ?
- Give your own example...

▪ Find \hat{f} – an estimate of f – where

- $\hat{Y} = \hat{f}(X)$ represents the vector of predicted values
- The overall (aggregated) prediction error between Y and \hat{Y} is minimized

▪ Example : minimize $E(Y - \hat{Y})^2$

- Assuming f and X fixed :

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

Reducible error

Irreducible error

▪ Goal : **Inference / Explainability**

- What factors (medical predictors) are most predictive of a prolonged stay ?
- What factors (organizational predictors) are most predictive of the turnover rate ?
- Give your own example...

▪ Relationship between the outcome and the predictors

- Type or nature of the relationship
- Strength of the relationship

Focus on **prediction performance only** raises the issue of the **Black Box Problem**.

Focus on the **explainability alone** raises the issue of **Prediction Reliability**.

How do we estimate f ? Part 1

▪ Parametric methods

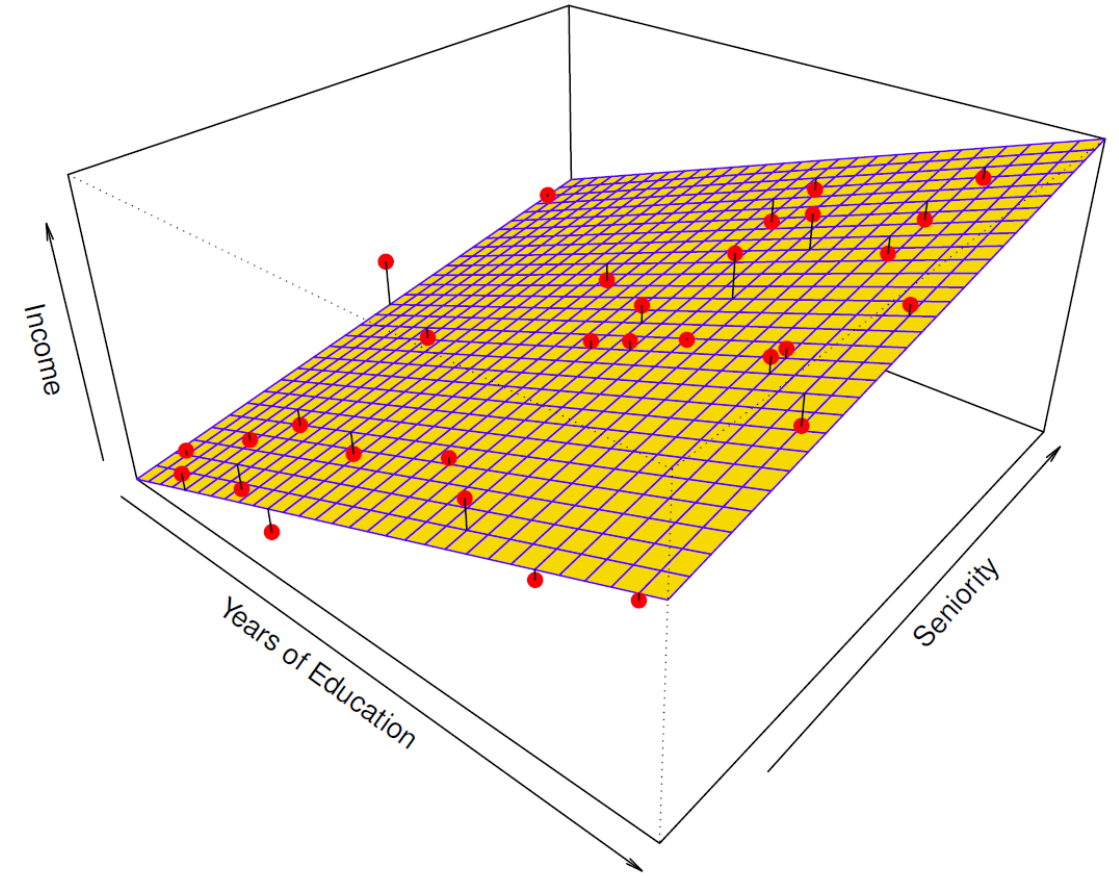
- We make explicit assumptions on the functional relationship between the outcome and the predictors.
- The problem of estimating f is reduced down to estimating a set of parameters.
- Rely on (statistical) modeling
- Example :
 - $Y = f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$
 - The problem of estimating f is reduced to estimating $[\beta_0, \beta_1, \dots, \beta_p]$

▪ Upsides

- Simplifies the estimation to a reduced number of parameters

▪ Downsides

- The model \hat{f} will not match the true f
 - More flexible models may lead to overfitting



How do we estimate f ? Part 2

▪ Non-Parametric methods

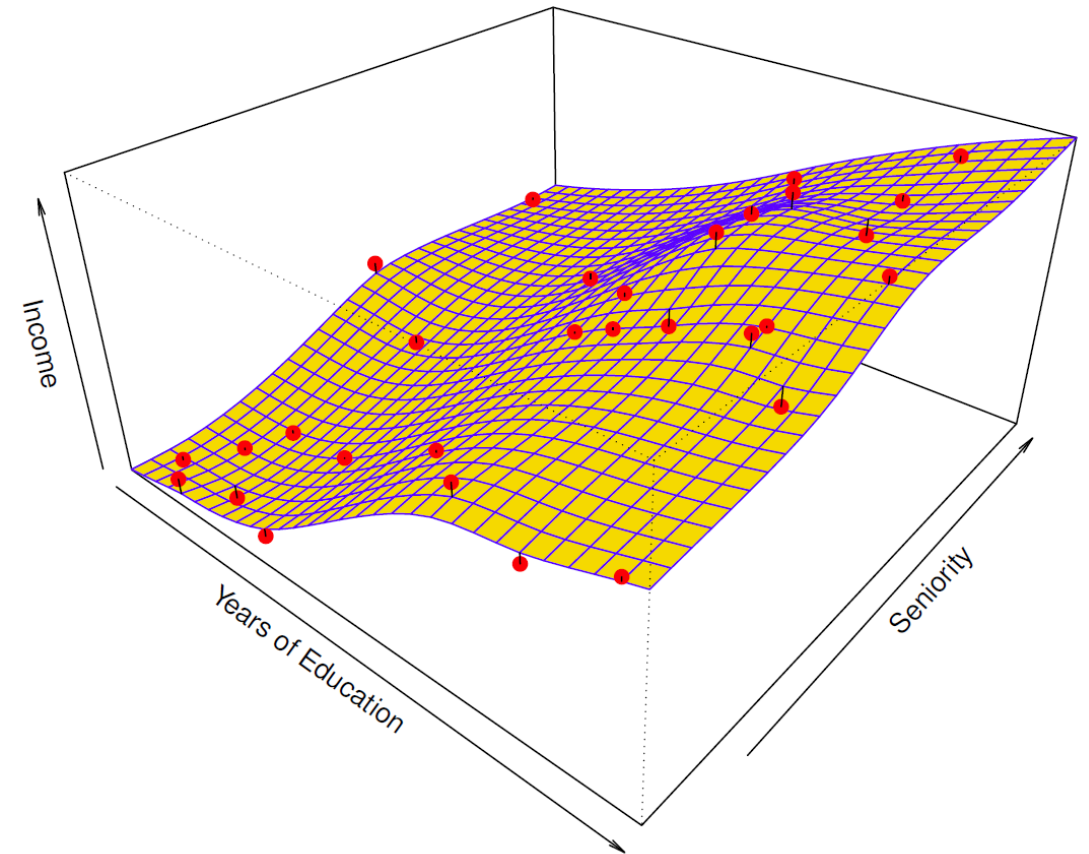
- No explicit assumptions on the functional relationship between the outcome and the predictors.
- Seek an estimate of f as close to the datapoints as possible.
- Rely on an algorithmic approach
- Example :
 - $Y = kNN(X)$: the k – nearest neighbors
 - kNN estimates each datapoint based on the values of the k -nearest neighbors

▪ Upsides

- Have the potential to accurately fit a wider range of possible shapes for f

▪ Downsides

- Large number of observations required to obtain an accurate estimate for f
- May lead easily to *overfitting*



Supervised Learning

Supervised Learning

- Given some observed data (X, Y)
 - Of n data pair points (x_i, y_i)
 - And p variables $\mathbf{X} = [X_1, X_2, \dots, X_p] = [x_{i,j}]$
 - X : input (features) is associated to Y : output (response)
 - We can learn Y from \mathbf{X} from the relationships in the n data points (x_i, y_i)
- Given a new set $\mathbf{x} = (x_\mu)$ of m data points
 - We can predict the corresponding (y_μ)
 - Based on the information learned in the n data points (x_i, y_i)

$$\begin{bmatrix} X_1 & X_2 & \dots & X_j & \dots & X_p & Y \\ x_{1,1} & x_{1,2} & \dots & x_{1,j} & \dots & x_{1,p} & y_1 \\ \vdots & & & \ddots & & & \vdots \\ x_{i,1} & & & \dots & x_{i,j} & \dots & y_i \\ & & & & \dots & & \\ x_{n,1} & x_{n,2} & \dots & x_{n,j} & \dots & x_{n,p} & y_n \end{bmatrix}$$

If the outcome (or output) Y is quantitative, the supervised learning is called a **regression**
If it is categorical, it is called a **classification**

Give you own example of regression and classification

Unsupervised Learning

■ Unsupervised Learning

- Given some observed data (X)
 - Of n data points (x_i)
 - As p variables $\mathbf{X} = [X_1, X_2, \dots, X_p] = [x_{i,j}]$
 - There is no given output in the data
- We look for patterns of similarity
 - Either between the observations (rows)
 - Usually by comparing the « distances » between observations
 - Or between the columns (variables)
 - Usually by looking at « distances » between columns
- Then aggregating those closest in distance
 - Finding sensible interpretations for each group of observations or of variables

$$\begin{bmatrix} X_1 & X_2 & \dots & X_j & \dots & X_p \\ x_{1,1} & x_{1,2} & \dots & x_{1,j} & \dots & x_{1,p} \\ \vdots & & & \ddots & & \\ x_{i,1} & x_{i,2} & \dots & x_{i,j} & \dots & \\ & & & \dots & & \\ x_{n,1} & x_{n,2} & \dots & x_{n,j} & \dots & x_{n,p} \end{bmatrix}$$

Assessing model performance

■ There is no model that is the best under all circumstances

- Performance depends on the model and the data
- Performance should be compared :
 - Between models
 - On the same dataset
 - The dataset used to train each model and to estimate their performance should not be the same

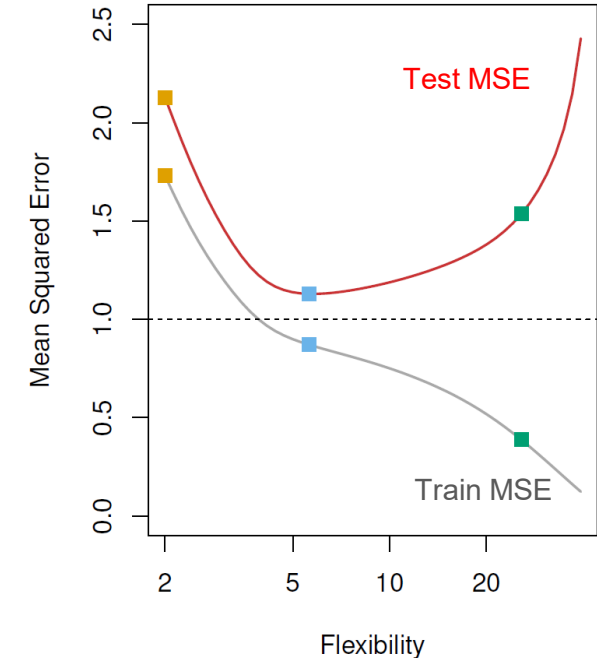
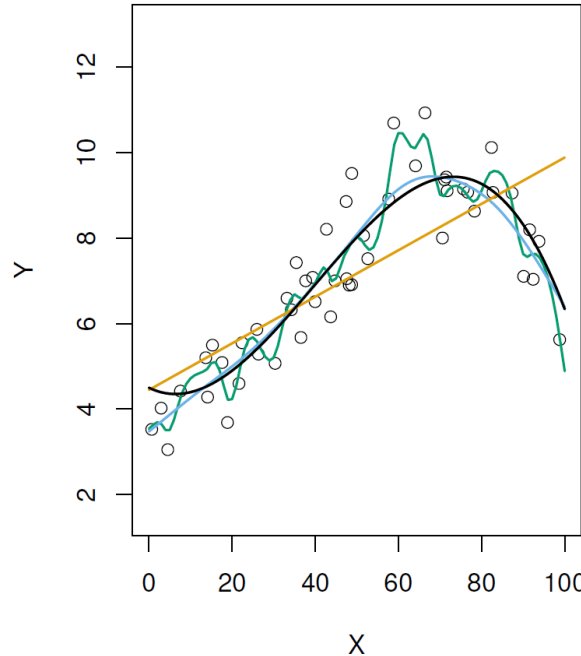
■ Simplest way to estimate model performance : aggregated error of prediction

- Measuring performance of a regression model
 - Aggregated distance between actual and predicted

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

- Measuring classification performance
 - Aggregated counts of misclassification

$$Accuracy = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$



When the performance of a model (here measured in MSE) is **much higher on the test data** than on the training data, we are **overfitting** the data

Bias-Variance Trade-Off

▪ Bias

- Error introduced by approximating a real-life problem

$$\text{bias}(\hat{y}) = E(\hat{y}) - y$$
- More flexible model tends to result in less bias (less prediction errors)

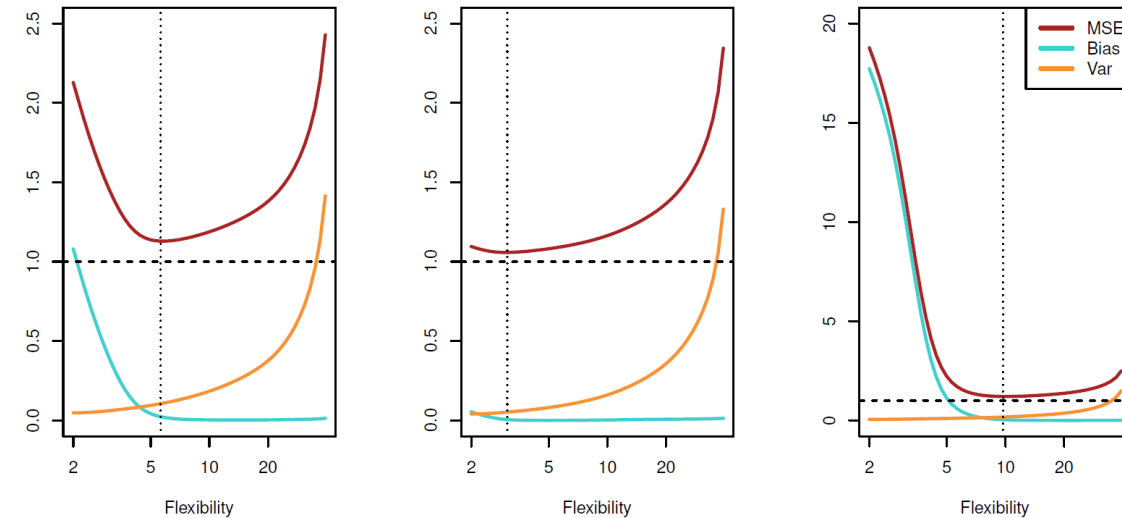
▪ Variance

- Amount by which \hat{f} would change if we estimated it from another dataset
- More flexible model tends to result in higher variance (less reliability)

▪ Expected MSE

$$E(y_0 - f(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{bias}(f(x_0))]^2 + \text{Var}(\varepsilon)$$

- The goal is to find a model that minimizes the bias and the variance simultaneously



The minimum value for the red curve represents the **Bias-Variance Trade Off** for respectively : **medium** flexibility, **low** flexibility and **high** flexibility

Class Exercise

- **Hands on...**

- Apply the principles discovered in this session on the <Flourishing> dataset