# Feature Selection in Machine Learning for BioMedical Data

Nov 25 2024
Giảng viên: TS. Lưu Phúc Lợi
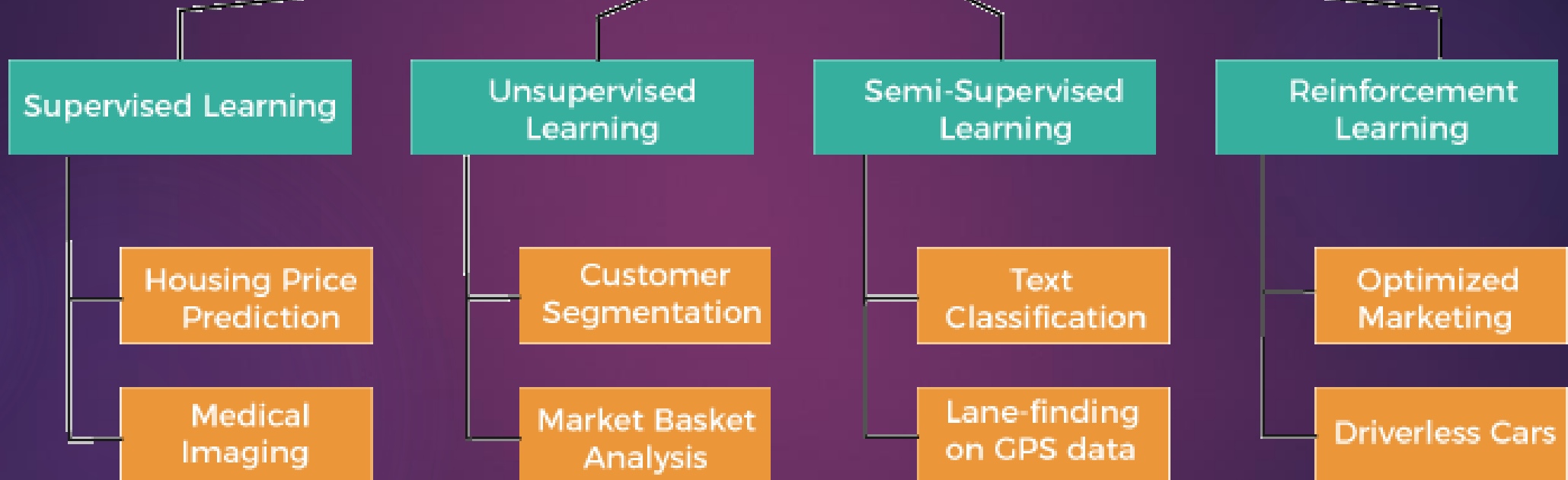Luu.p.loi@googlemail.com

# Content of Lecture 1

1. Introduction to Feature Selection
2. Filter method
3. Wrapper methods
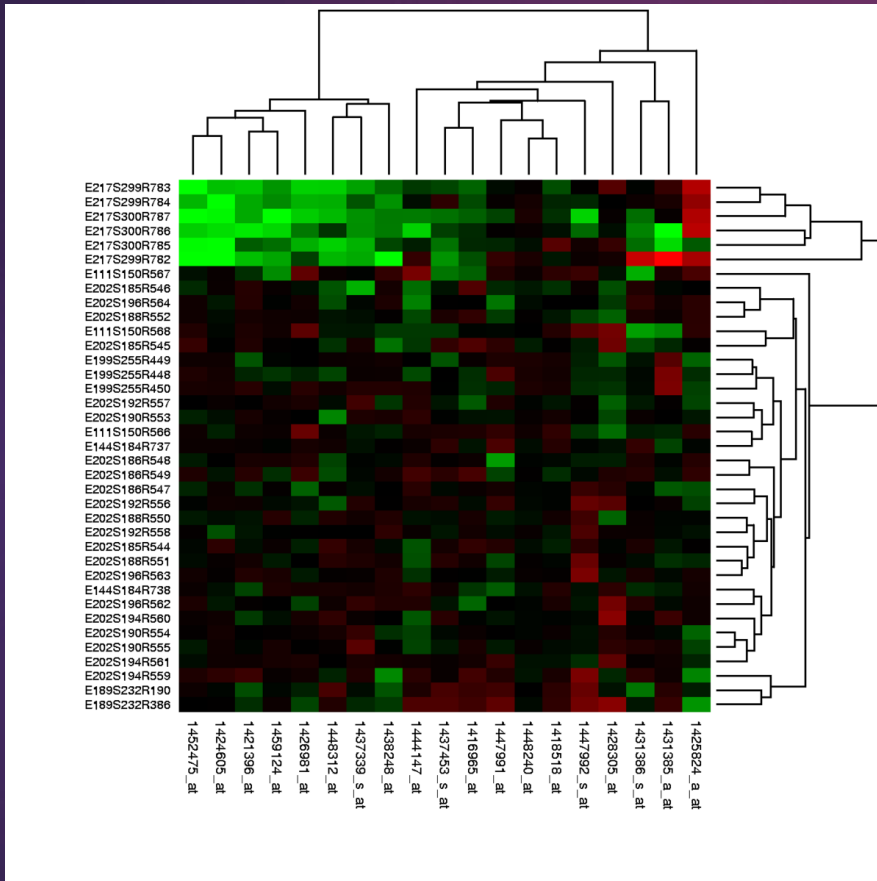4. Embedded method

# Introduction to Feature Selection
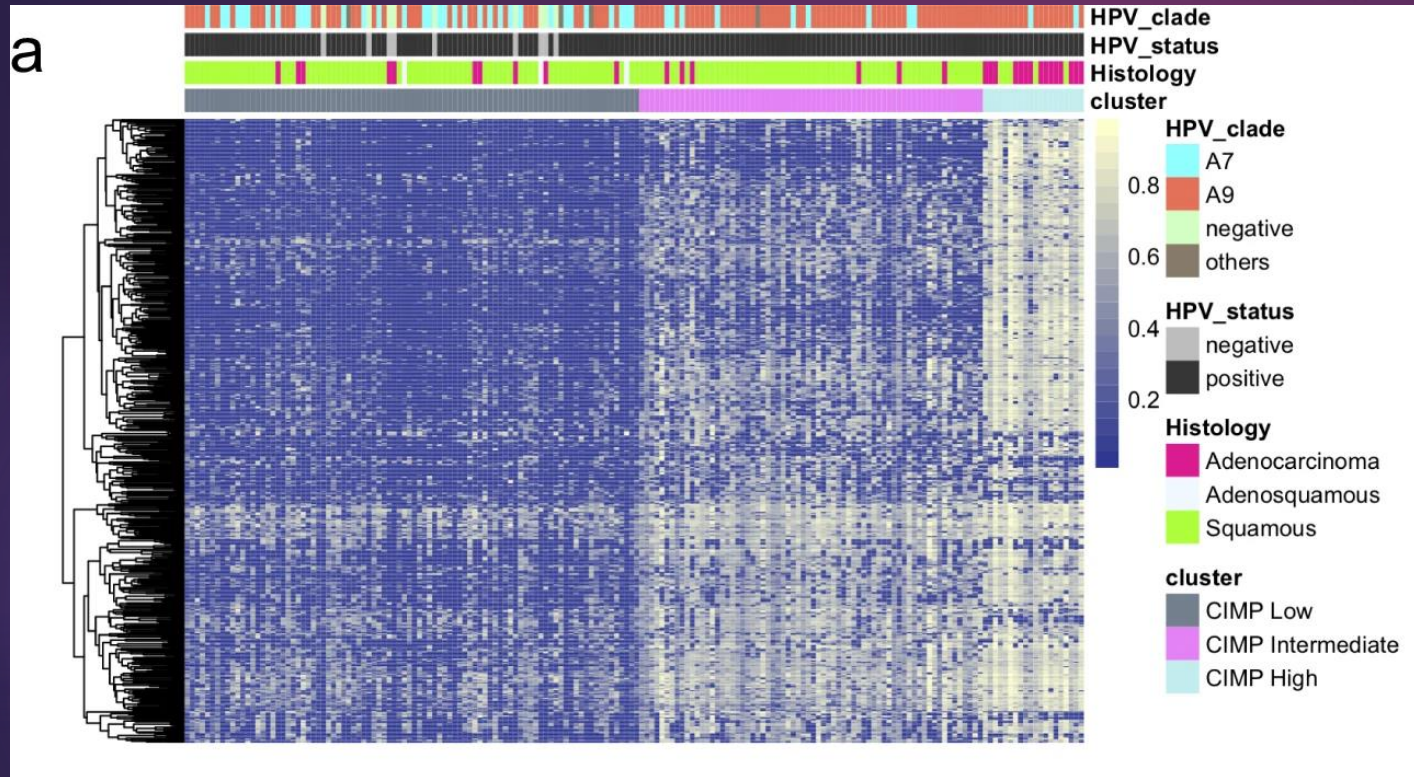
# Supervise Learning: regression or classification

$$
\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & & \ddots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}
$$

$$\mathbb{Y} \qquad\qquad\qquad \mathbb{X} \qquad\qquad\qquad \theta$$

# BioMedical data: gene expression with p >> n



n = number of samples 6, 10, 100, 1k
P = number of genes 20k

https://en.wikipedia.org/wiki/Gene_expression_profiling

# BioMedical data: DNA methylation (p >> n)



n = number of samples 6, 10, 100, 1k
P = number of CpG 28M

# What is feature/variable selection?

▶ Find the features (variables/columns) in X which are important for predicting, and remove the features that are not

▶ Give:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \ldots & x_{1p} \\ 1 & x_{21} & x_{22} & \ldots & x_{2p} \\ \vdots & & \ddots & & \vdots \\ 1 & x_{n1} & x_{n2} & \ldots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

$$\mathbb{Y} \qquad\qquad \mathbb{X} \qquad\qquad \theta$$

| Bias | Age | Height | Hours of sleep |
|------|-----|--------|----------------|
| 1 | 21 | 170.82 | 6 |
| 1 | 19 | 208.78 | 10 |
| 1 | 22 | 158.57 | 10 |
| 1 | 23 | 194.08 | 8 |
| 1 | 19 | 151.22 | 7 |
| ... | ... | ... | ... |
| 1 | 24 | 190.41 | 8 |
| 1 | 24 | 172.04 | 6 |
| 1 | 23 | 159.80 | 10 |
| 1 | 19 | 178.16 | 9 |
| 1 | 18 | 194.08 | 11 |

$$\longrightarrow \begin{bmatrix} 1 & x_{11} & x_{12} & \ldots & x_{1p} \\ 1 & x_{21} & x_{22} & \ldots & x_{2p} \\ \vdots & & \ddots & & \vdots \\ 1 & x_{n1} & x_{n2} & \ldots & x_{np} \end{bmatrix}$$
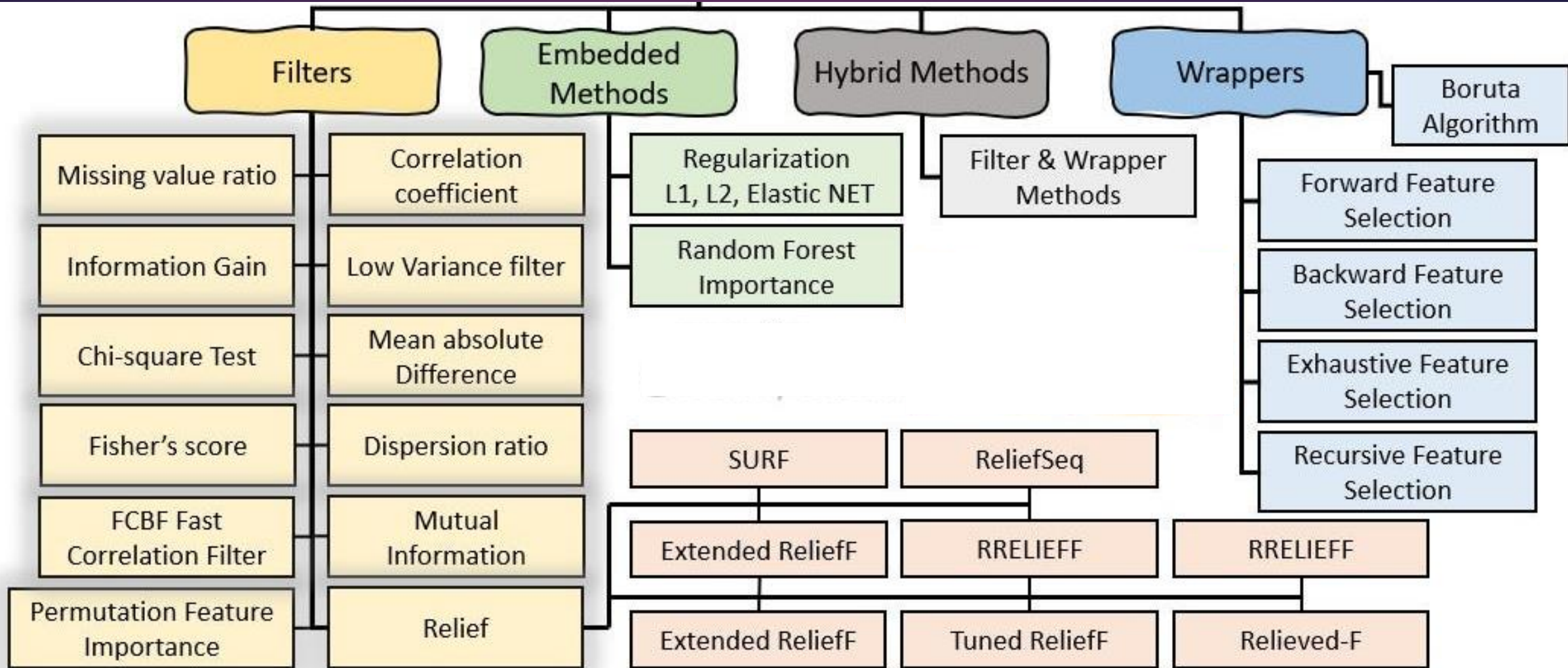
▶ Find the columns in X which are important for predicting y

# Why feature selection?

► Interpretability: Models are more interpretable with fewer features. If you get the same performance with 10 features instead of 500 features, why not use the model with smaller number of features?

► Computation: Models fit/predict faster with fewer columns.

► Data collection: What type of new data should I collect? It may be cheaper to collect fewer columns.

► Fundamental tradeoff: Can I reduce overfitting by removing useless features?

► Feature selection can often result in better performing (less overfit), easier to understand, and faster model.

# How do we carry out feature selection?
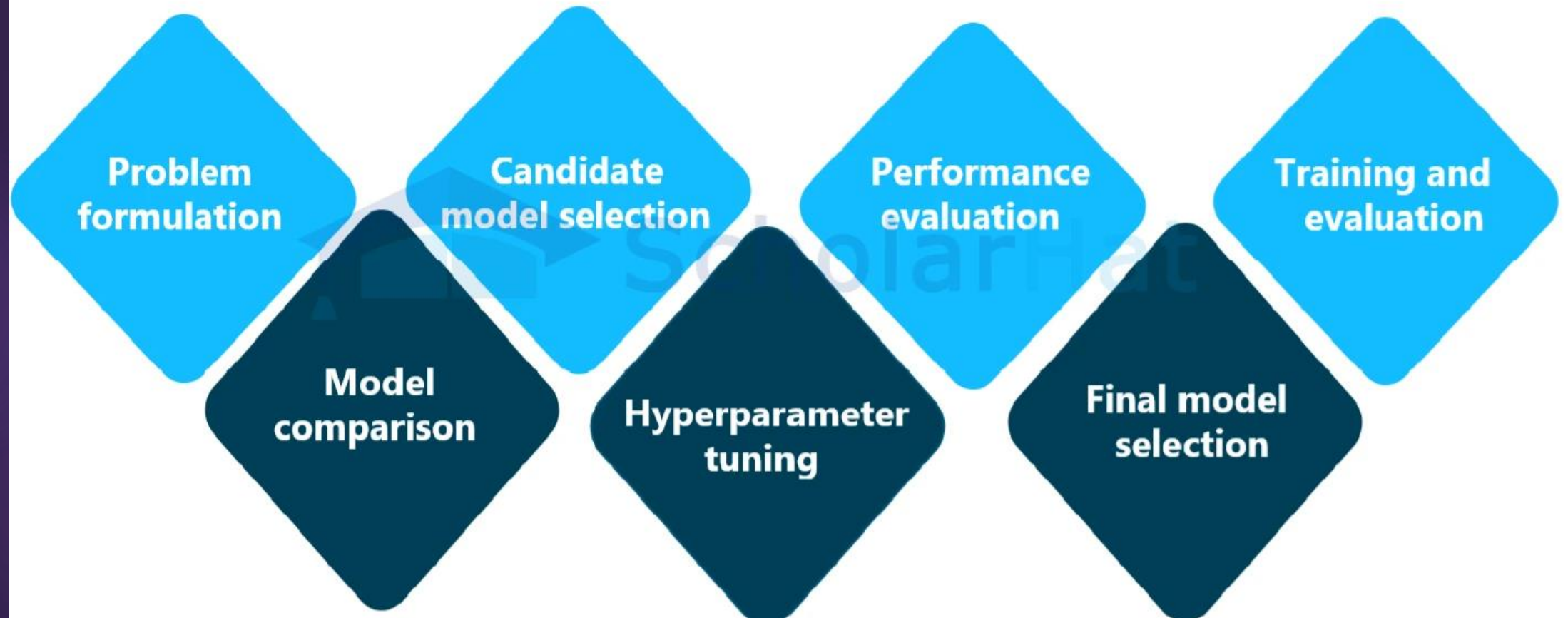## Supervised Feature Selection

# Model Selection vs Feature Selection

▶ Feature Selection is a part of Model Selection
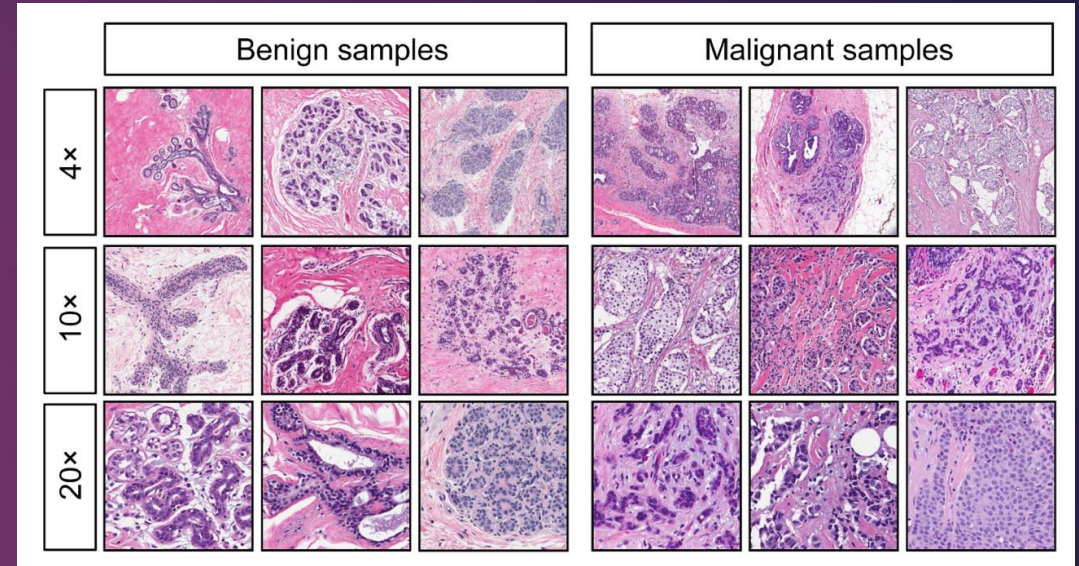
# Model selection: steps

# Model selection: steps
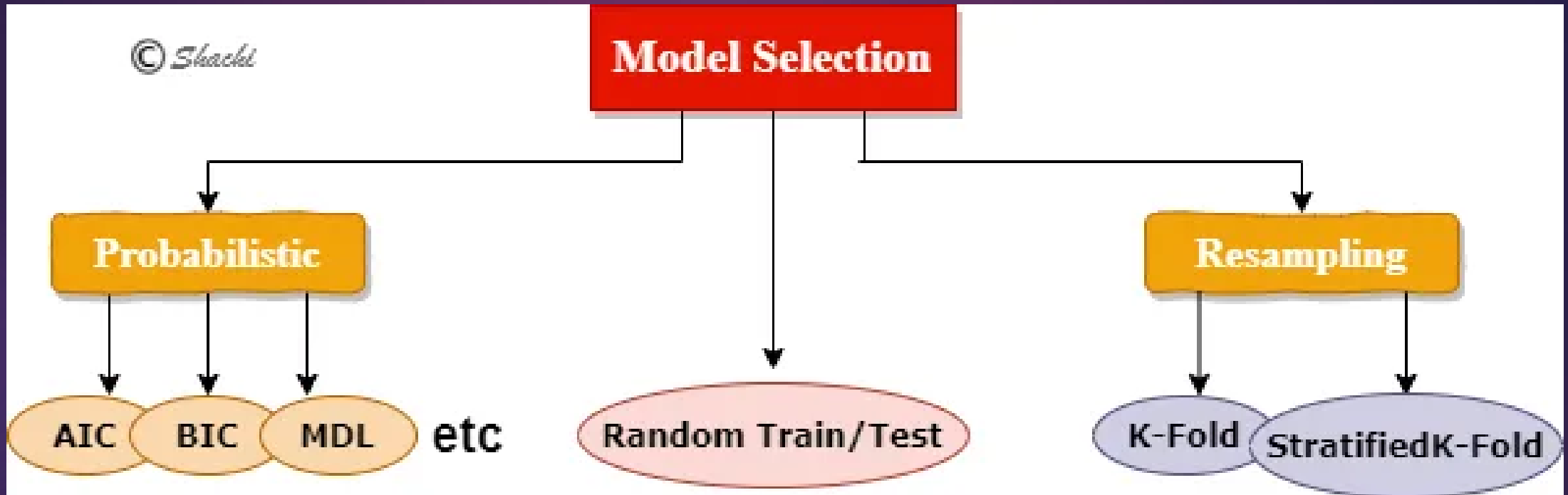
**Stage 1: Selecting the regression model forms**

**Stage 2: Selecting the regression model and the independent variables**

**Stage 3: Fitting the model**

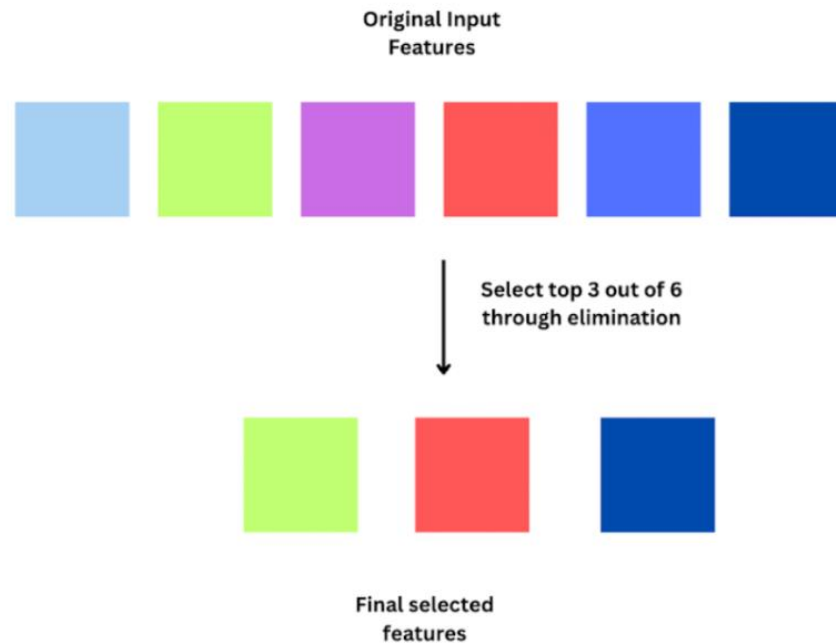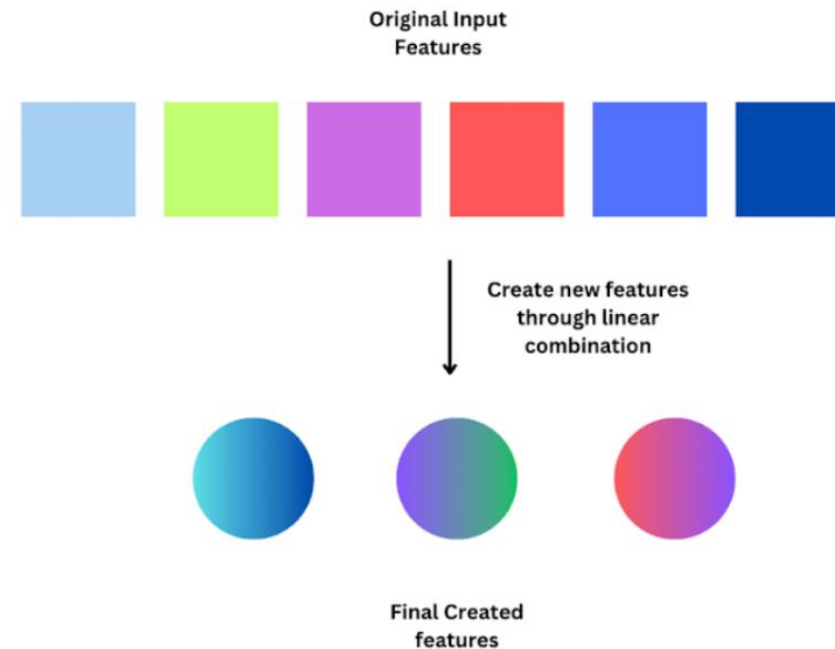**Stage 4: Examining or validation of the applied model**

# Model selection: methods

# Feature Selection vs Feature Extraction/Engineering

# Xin chân thành cảm ơn!

LUU PHUC LOI, PHD

ZALO: 0901802182

LUU.P.LOI@GOOGLEMAIL.COM