# Feature Selection in Machine Learning for BioMedical Data

Nov 07 2025
Giảng viên: TS. Lưu Phúc Lợi
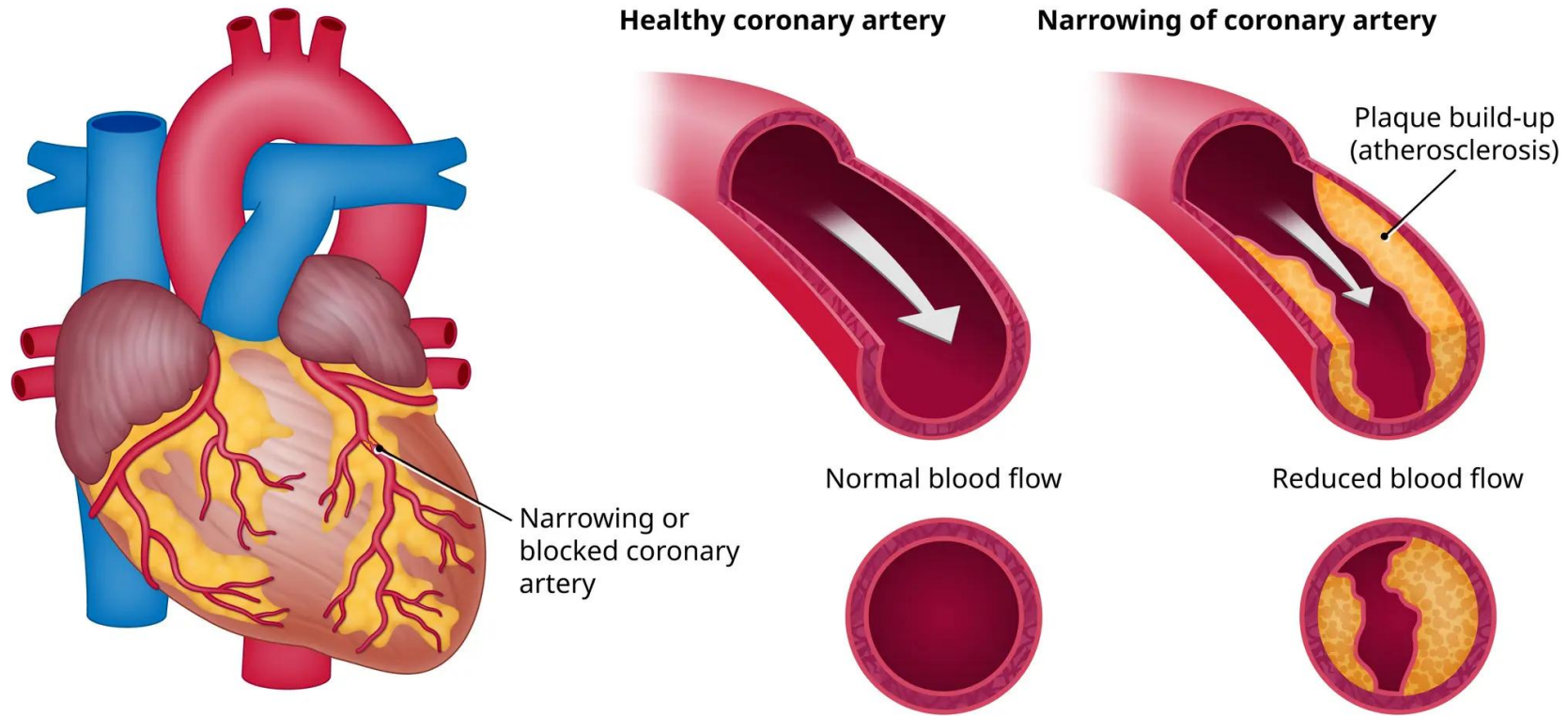Email: Luu.p.loi@googlemail.com
Zalo: 0901802182

# Content

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

# Coronary Artery Disease

**Healthy coronary artery**

**Narrowing of coronary artery**

Plaque build-up (atherosclerosis)

Narrowing or blocked coronary artery

Normal blood flow

Reduced blood flow

ORIGINAL ARTICLE

# A noninvasive method for coronary artery diseases diagnosis using a clinically-interpretable fuzzy rule-based system

Marateb, Hamid Reza; Goudarzi, Sobhan

Author Information ⊙

OPEN

---

## Abstract

### Background:

Coronary heart diseases/coronary artery diseases (CHDs/CAD), the most common form of cardiovascular disease (CVD), are a major cause for death and disability in developing/developed countries. CAD risk factors could be detected by physicians to prevent the CAD occurrence in the near future. Invasive coronary angiography, a current diagnosis method, is costly and associated with morbidity and mortality in CAD patients. The aim of this study was to design a computer-based noninvasive CAD diagnosis system with clinically interpretable rules.

## Materials and Methods:

In this study, the Cleveland CAD dataset from the University of California UCI (Irvine) was used. The interval-scale variables were discretized, with cut points taken from the literature. A fuzzy rule-based system was then formulated based on a neuro-fuzzy classifier (NFC) whose learning procedure was speeded up by the scaled conjugate gradient algorithm. Two feature selection (FS) methods, multiple logistic regression (MLR) and sequential FS, were used to reduce the required attributes. The performance of the NFC (without/with FS) was then assessed in a hold-out validation framework. Further cross-validation was performed on the best classifier.

## Results:

In this dataset, 16 complete attributes along with the binary CHD diagnosis (gold standard) for 272 subjects (68% male) were analyzed. MLR + NFC showed the best performance. Its overall sensitivity, specificity, accuracy, type I error ($\alpha$) and statistical power were 79%, 89%, 84%, 0.1 and 79%, respectively. The selected features were "age and ST/heart rate slope categories," "exercise-induced angina status," fluoroscopy, and thallium-201 stress scintigraphy results.

## Conclusion:

The proposed method showed "substantial agreement" with the gold standard. This algorithm is thus, a promising tool for screening CAD patients.

# Problem statement
## Predicting Heart Disease from the Cleveland Heart Disease Dataset (303 samples x 76 features)

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
| 2 | 63 | 1 | 0 | 145 | 233 | 1 | 2 | 150 | 0 | 2.3 | 2 | 0 | 2 | 0 |
| 3 | 67 | 1 | 3 | 160 | 286 | 0 | 2 | 108 | 1 | 1.5 | 1 | 3 | 1 | 1 |
| 4 | 67 | 1 | 3 | 120 | 229 | 0 | 2 | 129 | 1 | 2.6 | 1 | 2 | 3 | 1 |
| 5 | 37 | 1 | 2 | 130 | 250 | 0 | 0 | 187 | 0 | 3.5 | 2 | 0 | 1 | 0 |
| 6 | 41 | 0 | 1 | 130 | 204 | 0 | 2 | 172 | 0 | 1.4 | 0 | 0 | 1 | 0 |
| 7 | 56 | 1 | 1 | 120 | 236 | 0 | 0 | 178 | 0 | 0.8 | 0 | 0 | 1 | 0 |
| 8 | 62 | 0 | 3 | 140 | 268 | 0 | 2 | 160 | 0 | 3.6 | 2 | 2 | 1 | 1 |
| 9 | 57 | 0 | 3 | 120 | 354 | 0 | 0 | 163 | 1 | 0.6 | 0 | 0 | 1 | 0 |
| 10 | 63 | 1 | 3 | 130 | 254 | 0 | 2 | 147 | 0 | 1.4 | 1 | 1 | 3 | 1 |
| 11 | 53 | 1 | 3 | 140 | 203 | 1 | 2 | 155 | 1 | 3.1 | 2 | 0 | 3 | 1 |
| 12 | 57 | 1 | 3 | 140 | 192 | 0 | 0 | 148 | 0 | 0.4 | 1 | 0 | 2 | 0 |
| 13 | 56 | 0 | 1 | 140 | 294 | 0 | 2 | 153 | 0 | 1.3 | 1 | 0 | 1 | 0 |
| 14 | 56 | 1 | 2 | 130 | 256 | 1 | 2 | 142 | 1 | 0.6 | 1 | 1 | 2 | 1 |
| 15 | 44 | 1 | 1 | 120 | 263 | 0 | 0 | 173 | 0 | 0 | 0 | 0 | 3 | 0 |
| 16 | 52 | 1 | 2 | 172 | 199 | 1 | 0 | 162 | 0 | 0.5 | 0 | 0 | 3 | 0 |
| 17 | 57 | 1 | 2 | 150 | 168 | 0 | 0 | 174 | 0 | 1.6 | 0 | 0 | 1 | 0 |
| 18 | 48 | 1 | 1 | 110 | 229 | 0 | 0 | 168 | 0 | 1 | 2 | 0 | 3 | 1 |
| 19 | 54 | 1 | 3 | 140 | 239 | 0 | 0 | 160 | 0 | 1.2 | 0 | 0 | 1 | 0 |
| 20 | 48 | 0 | 2 | 130 | 275 | 0 | 0 | 139 | 0 | 0.2 | 0 | 0 | 1 | 0 |
| 21 | 49 | 1 | 1 | 130 | 266 | 0 | 0 | 171 | 0 | 0.6 | 0 | 0 | 1 | 0 |
| 22 | 64 | 1 | 0 | 110 | 211 | 0 | 2 | 144 | 1 | 1.8 | 1 | 0 | 1 | 0 |

# Predicting Heart Disease from the Cleveland Heart Disease Dataset
## (303 samples x 13 features)

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
| 2 | 63 | 1 | 0 | 145 | 233 | 1 | 2 | 150 | 0 | 2.3 | 2 | 0 | 2 | 0 |
| 3 | 67 | 1 | 3 | 160 | 286 | 0 | 2 | 108 | 1 | 1.5 | 1 | 3 | 1 | 1 |
| 4 | 67 | 1 | 3 | 120 | 229 | 0 | 2 | 129 | 1 | 2.6 | 1 | 2 | 3 | 1 |
| 5 | 37 | 1 | 2 | 130 | 250 | 0 | 0 | 187 | 0 | 3.5 | 2 | 0 | 1 | 0 |
| 6 | 41 | 0 | 1 | 130 | 204 | 0 | 2 | 172 | 0 | 1.4 | 0 | 0 | 1 | 0 |
| 7 | 56 | 1 | 1 | 120 | 236 | 0 | 0 | 178 | 0 | 0.8 | 0 | 0 | 1 | 0 |
| 8 | 62 | 0 | 3 | 140 | 268 | 0 | 2 | 160 | 0 | 3.6 | 2 | 2 | 1 | 1 |
| 9 | 57 | 0 | 3 | 120 | 354 | 0 | 0 | 163 | 1 | 0.6 | 0 | 0 | 1 | 0 |
| 10 | 63 | 1 | 3 | 130 | 254 | 0 | 2 | 147 | 0 | 1.4 | 1 | 1 | 3 | 1 |
| 11 | 53 | 1 | 3 | 140 | 203 | 1 | 2 | 155 | 1 | 3.1 | 2 | 0 | 3 | 1 |
| 12 | 57 | 1 | 3 | 140 | 192 | 0 | 0 | 148 | 0 | 0.4 | 1 | 0 | 2 | 0 |
| 13 | 56 | 0 | 1 | 140 | 294 | 0 | 2 | 153 | 0 | 1.3 | 1 | 0 | 1 | 0 |
| 14 | 56 | 1 | 2 | 130 | 256 | 1 | 2 | 142 | 1 | 0.6 | 1 | 1 | 2 | 1 |
| 15 | 44 | 1 | 1 | 120 | 263 | 0 | 0 | 173 | 0 | 0 | 0 | 0 | 3 | 0 |
| 16 | 52 | 1 | 2 | 172 | 199 | 1 | 0 | 162 | 0 | 0.5 | 0 | 0 | 3 | 0 |
| 17 | 57 | 1 | 2 | 150 | 168 | 0 | 0 | 174 | 0 | 1.6 | 0 | 0 | 1 | 0 |
| 18 | 48 | 1 | 1 | 110 | 229 | 0 | 0 | 168 | 0 | 1 | 2 | 0 | 3 | 1 |
| 19 | 54 | 1 | 3 | 140 | 239 | 0 | 0 | 160 | 0 | 1.2 | 0 | 0 | 1 | 0 |
| 20 | 48 | 0 | 2 | 130 | 275 | 0 | 0 | 139 | 0 | 0.2 | 0 | 0 | 1 | 0 |
| 21 | 49 | 1 | 1 | 130 | 266 | 0 | 0 | 171 | 0 | 0.6 | 0 | 0 | 1 | 0 |
| 22 | 64 | 1 | 0 | 110 | 211 | 0 | 2 | 144 | 1 | 1.8 | 1 | 0 | 1 | 0 |

# Cleveland Heart Disease Dataset

This database contains 13 attributes and a target variable. It has 8 nominal values and 5 numeric values. The detailed description of all these features are as follows:

1. Age: Patients Age in years (Numeric)

2. Sex: Gender (Male : 1; Female : 0) (Nominal)

3. cp: Type of chest pain experienced by patient. This term categorized into 4 category. 0 typical angina, 1 atypical angina, 2 non- anginal pain, 3 asymptomatic (Nominal)

4. trestbps: patient's level of blood pressure at resting mode in mm/HG (Numerical)

5. chol: Serum cholesterol in mg/dl (Numeric)

6. fbs: Blood sugar levels on fasting > 120 mg/dl represents as 1 in case of true and 0 as false (Nominal)

7. restecg: Result of electrocardiogram while at rest are represented in 3 distinct values 0 : Normal 1: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV) 2: showing probable or definite left ventricular hypertrophyby Estes' criteria (Nominal)

8. thalach: Maximum heart rate achieved (Numeric)

# Cleveland Heart Disease Dataset

This database contains 13 attributes and a target variable. It has 8 nominal values and 5 numeric values. The detailed description of all these features are as follows:

9. exang: Angina induced by exercise 0 depicting NO 1 depicting Yes (Nominal)

10. oldpeak: Exercise induced ST-depression in relative with the state of rest (Numeric)

11. slope: ST segment measured in terms of slope during peak exercise
0: up sloping; 1: flat; 2: down sloping(Nominal)

12. ca: The number of major vessels (0–3)(nominal)

13. thal: A blood disorder called thalassemia
0: NULL 1: normal blood flow 2: fixed defect (no blood flow in some part of the heart) 3: reversible defect (a blood flow is observed but it is not normal(nominal)

14. target: It is the target variable which we have to predict 1 means patient is suffering from heart disease and 0 means patient is normal.
**Variable to be predicted**
Absence (1) or presence (2) of heart disease

# Oral Cavity and Pharyngeal Cancer

# Genome-wide association analyses identify new susceptibility loci for oral cavity and pharyngeal cancer

Corina Lesseur, Brenda Diergaarde, Andrew F Olshan, Victor Wünsch-Filho, Andrew R Ness, Geoffrey Liu, Martin Lacko, José Eluf-Neto, Silvia Franceschi, Pagona Lagiou, Gary J Macfarlane, Lorenzo Richiardi, Stefania Boccia, Jerry Polesel, Kristina Kjaerheim, David Zaridze, Mattias Johansson, Ana M Menezes, Maria Paula Curado, Max Robinson, Wolfgang Ahrens, Cristina Canova, Ariana Znaor, Xavier Castellsagué, ... Paul Brennan ✉ 

+ Show authors

https://www.nature.com/articles/ng.3685

## Abstract

We conducted a genome-wide association study of oral cavity and pharyngeal cancer in 6,034 cases and 6,585 controls from Europe, North America and South America. We detected eight significantly associated loci ($P < 5 \times 10^{-8}$), seven of which are new for these cancer sites. Oral and pharyngeal cancers combined were associated with loci at 6p21.32 (rs3828805, *HLA-DQB1*), 10q26.13 (rs201982221, *LHPP*) and 11p15.4 (rs1453414, *OR52N2–TRIM5*). Oral cancer was associated with two new regions, 2p23.3 (rs6547741, *GPN1*) and 9q34.12 (rs928674, *LAMC3*), and with known cancer-related loci−9p21.3 (rs8181047, *CDKN2B-AS1*) and 5p15.33 (rs10462706, *CLPTM1L*). Oropharyngeal cancer associations were limited to the human leukocyte antigen (HLA) region, and classical HLA allele imputation showed a protective association with the class II haplotype HLA-DRB1*1301−HLA-DQA1*0103−HLA-DQB1*0603 (odds ratio (OR) = 0.59, $P = 2.7 \times 10^{-9}$). Stratified analyses on a subgroup of oropharyngeal cases with information available on human papillomavirus (HPV) status indicated that this association was considerably stronger in HPV-positive (OR = 0.23, $P = 1.6 \times 10^{-6}$) than in HPV-negative (OR = 0.75, $P = 0.16$) cancers.

# SNP - Oral Cavity and Pharyngeal Cancer Dataset

| note | Geographic region: Europe |
|---|---|
| ncase | 2497 |
| pmid | 27749845 |
| nsnp | 7514278 |
| ncontrol | 2928 |

# SNP - Oral Cavity and Pharyngeal Cancer Dataset
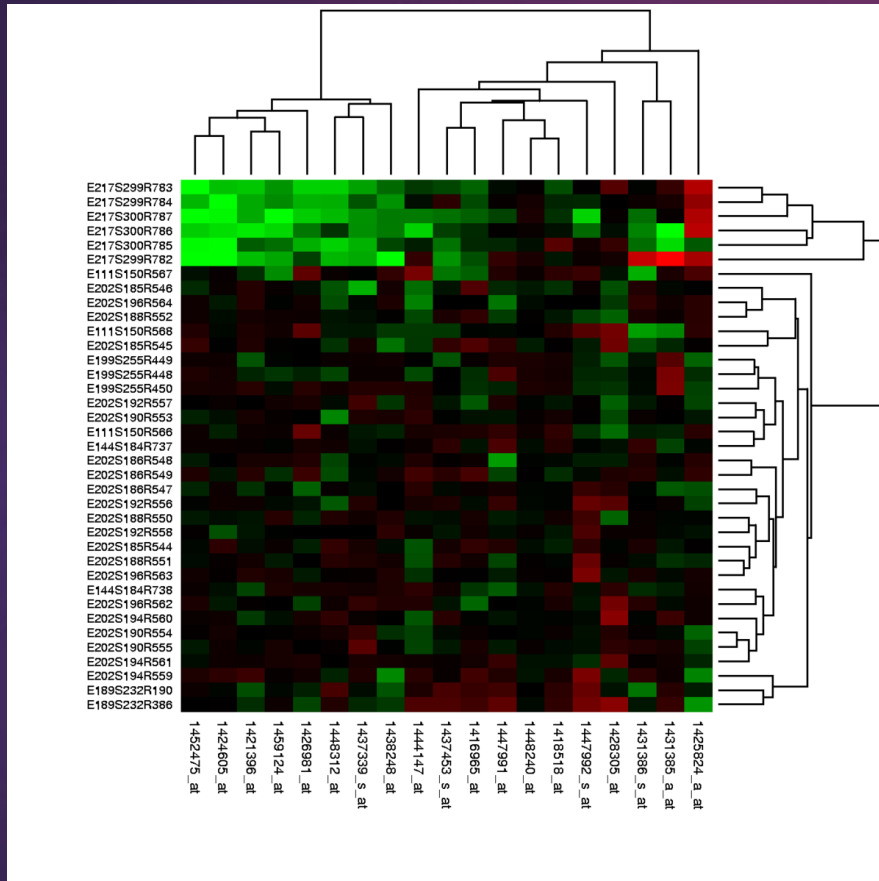


https://opengwas.io/datasets/ieu-b-89#
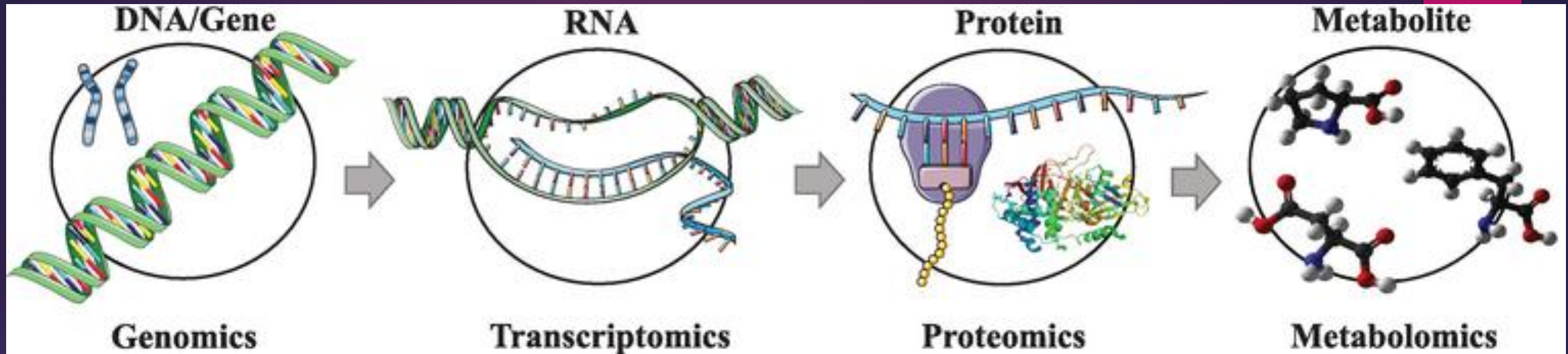
# Introduction to Feature Selection

# 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

# Mối liên kết: Biến thể gen và bệnh di truyền



DNA/Gene — RNA — Protein — Metabolite

Genomics — Transcriptomics — Proteomics — Metabolomics

PAH gene
Ref  …A**T**CGAT…
P1   …A**A**CGAT…

NM_000277.3(PAH):c.971T>A

PAH mRNA
Ref  …A**U**CGAU…
P1   …A**A**CGAU…

NM_000277.3(PAH):c.971T>A

PAH protein
Ref  …**Ile**-Asp…
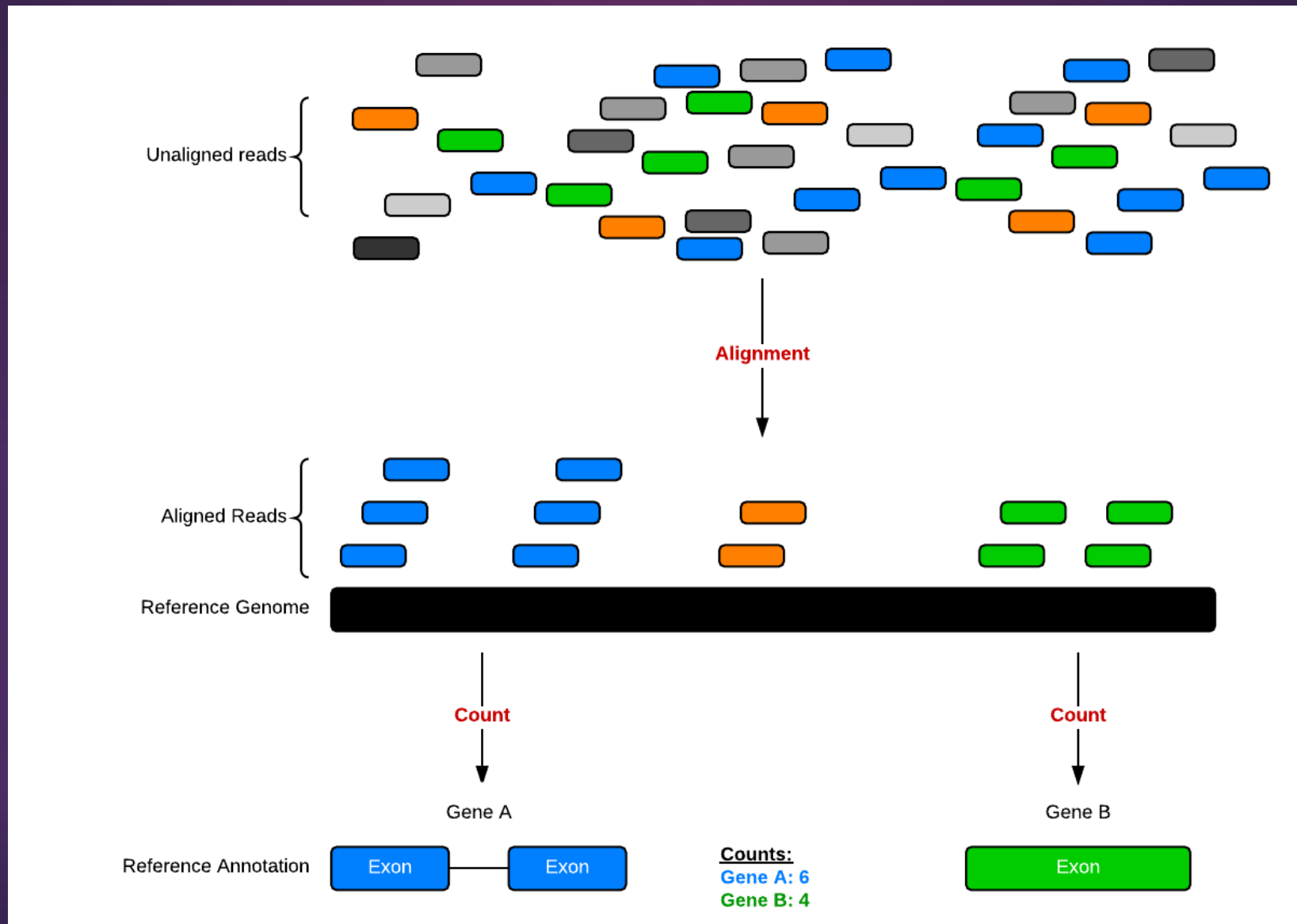P1   …**Asn**-Asp…

NM_000277.3(PAH):p.Ile324Asn

PAH
Ref  Phe → Tyr

**PAH**
P1   Phe ⟶ Tyr

# RNA-seq count table

# RNA-seq count table

```
## # A tibble: 38,694 x 9
##          ensgene SRR1039508 SRR1039509 SRR1039512 SRR1039513 SRR1039516
##            <chr>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
##   1 ENSG00000000003        723        486        904        445       1170
##   2 ENSG00000000005          0          0          0          0          0
##   3 ENSG00000000419        467        523        616        371        582
##   4 ENSG00000000457        347        258        364        237        318
##   5 ENSG00000000460         96         81         73         66        118
##   6 ENSG00000000938          0          0          1          0          2
##   7 ENSG00000000971       3413       3916       6000       4308       6424
##   8 ENSG00000001036       2328       1714       2640       1381       2165
##   9 ENSG00000001084        670        372        692        448        917
## 10 ENSG00000001167        426        295        531        178        740
## # ... with 38,684 more rows, and 3 more variables: SRR1039517 <dbl>,
## #   SRR1039520 <dbl>, SRR1039521 <dbl>
```
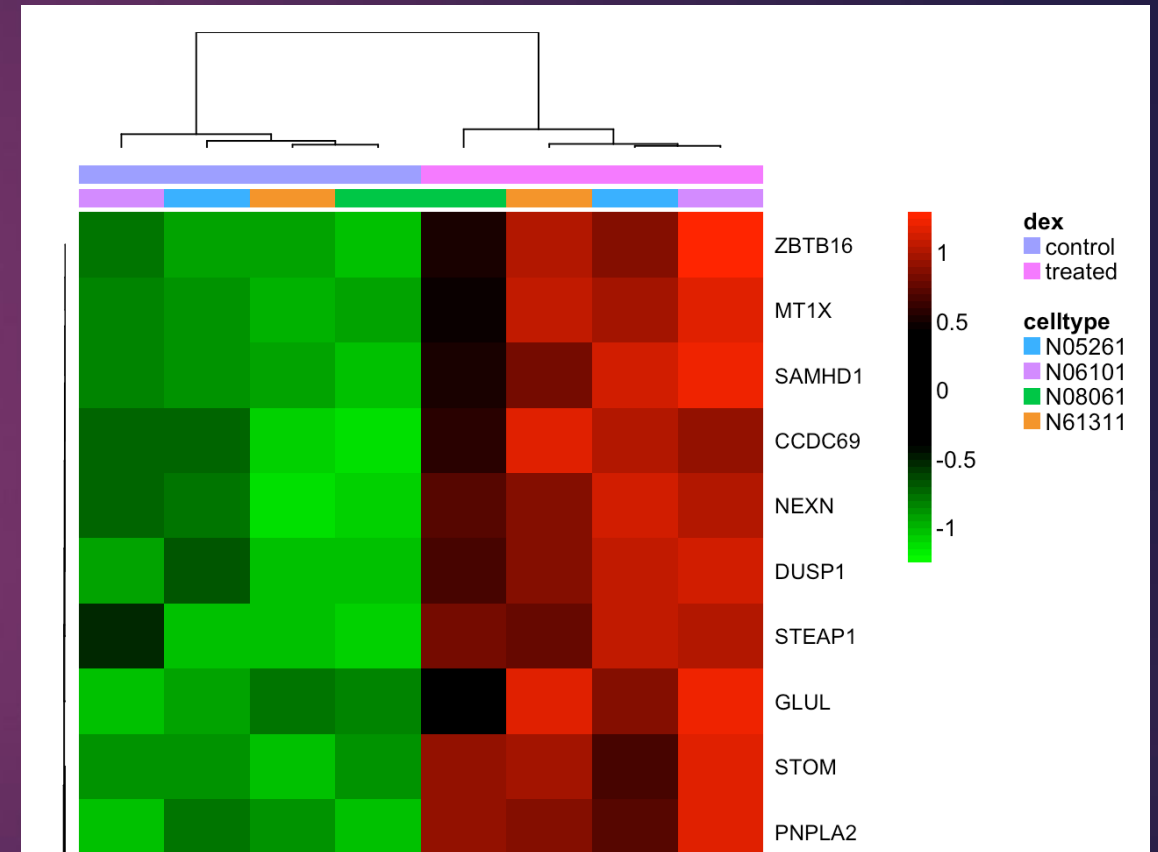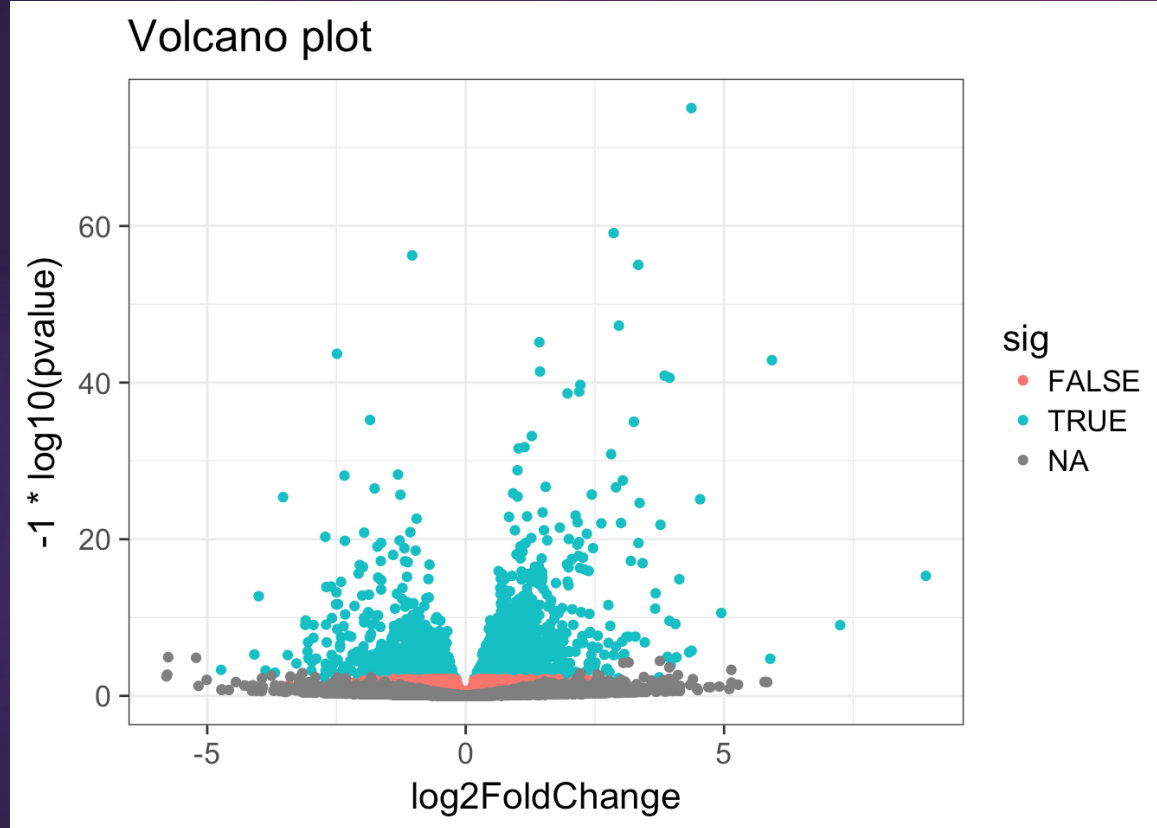
https://bioconnector.github.io/workshops/r-rnaseq-airway.html

# RNA-seq Downstream Analysis

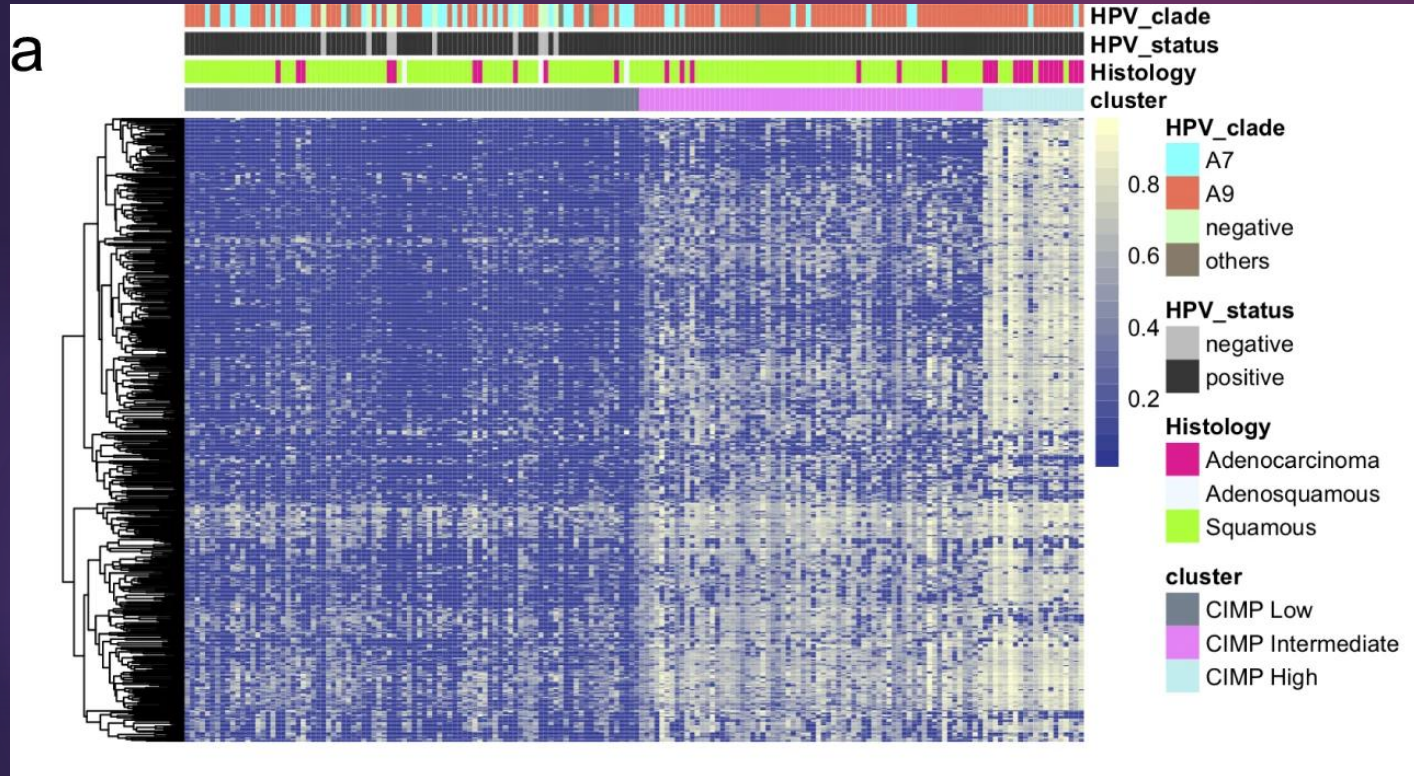# RNA-seq Downstream Analysis



https://bioconnector.github.io/workshops/r-rnaseq-airway.html

# RNA-seq Downstream Analysis

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



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

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



| | A | B chr1:946086 | C chr1:946098 | D chr1:946108 | E chr1:1005383 | F chr1:1005390 | G chr1:1005395 | H chr1:1005397 | I chr1:1005410 | J chr1:1005432 | K chr1:2163528 | L chr1:2163534 | M chr1:2163547 | N chr1:2163550 | O chr1:2163567 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | | | | | | |
| 2 | in3380_10 | 87 | 79 | 92 | 33 | 45 | 39 | 61 | 41 | 54 | 23 | 12 | 1 | 2 | 2 |
| 3 | in3380_11 | 96 | 79 | 95 | 39 | 49 | 48 | 67 | 48 | 62 | 18 | 10 | 0 | 1 | 1 |
| 4 | in3380_12 | 93 | 84 | 94 | 46 | 60 | 55 | 78 | 59 | 67 | 27 | 18 | 3 | 3 | 3 |
| 5 | in3380_13 | 92 | 86 | 95 | 37 | 49 | 42 | 70 | 47 | 65 | 37 | 23 | 5 | 6 | 8 |
| 6 | in3380_14 | 93 | 85 | 95 | 31 | 46 | 40 | 69 | 41 | 59 | 38 | 25 | 5 | 5 | 7 |
| 7 | in3380_15 | 94 | 86 | 96 | 36 | 50 | 41 | 72 | 45 | 56 | 24 | 14 | 5 | 7 | 5 |
| 8 | in3380_16 | 98 | 87 | 94 | 30 | 41 | 36 | 63 | 39 | 49 | 25 | 15 | 2 | 3 | 0 |
| 9 | in3380_17 | 93 | 78 | 94 | 33 | 48 | 40 | 70 | 43 | 58 | 23 | 14 | 2 | 4 | 0 |
| 10 | in3380_18 | 91 | 81 | 91 | 31 | 41 | 38 | 62 | 43 | 61 | 25 | 15 | 1 | 2 | 0 |
| 11 | in3380_19 | 92 | 81 | 93 | 22 | 33 | 34 | 59 | 37 | 62 | 14 | 8 | 1 | 1 | 0 |
| 12 | in3380_1 | 100 | 80 | 100 | 34 | 47 | 39 | 65 | 42 | 58 | 17 | 12 | 0 | 0 | 5 |
| 13 | in3380_22 | 91 | 80 | 91 | 25 | 36 | 31 | 50 | 34 | 50 | 21 | 11 | 2 | 2 | 1 |
| 14 | in3380_23 | 89 | 73 | 93 | 20 | 28 | 27 | 45 | 30 | 61 | 18 | 12 | 1 | 1 | 2 |
| 15 | in3380_24 | 94 | 85 | 95 | 37 | 49 | 42 | 66 | 43 | 59 | 22 | 14 | 2 | 3 | 2 |
| 16 | in3380_25 | 95 | 87 | 97 | 37 | 47 | 43 | 68 | 46 | 58 | 19 | 11 | 2 | 2 | 3 |
| 17 | in3380_26 | 91 | 80 | 93 | 27 | 40 | 38 | 65 | 45 | 63 | 17 | 10 | 1 | 0 | 1 |
| 18 | in3380_27 | 94 | 87 | 95 | 25 | 39 | 38 | 64 | 52 | 72 | 11 | 7 | 1 | 1 | 2 |
| 19 | in3380_28 | 90 | 76 | 94 | 33 | 45 | 39 | 65 | 44 | 53 | 20 | 10 | 0 | 1 | 0 |
| 20 | in3380_29 | 94 | 82 | 94 | 30 | 44 | 35 | 63 | 38 | 52 | 16 | 9 | 1 | 1 | 0 |
| 21 | in3380_2 | 95 | 81 | 95 | 34 | 46 | 41 | 65 | 42 | 56 | 27 | 17 | 3 | 4 | 3 |
| 22 | in3380_3 | 94 | 80 | 95 | 35 | 43 | 39 | 66 | 41 | 61 | 22 | 16 | 1 | 2 | 2 |
| 23 | in3380_4 | 83 | 78 | 91 | 29 | 45 | 36 | 68 | 39 | 63 | 31 | 18 | 0 | 1 | 2 |
| 24 | in3380_5 | 92 | 84 | 93 | 31 | 43 | 37 | 66 | 42 | 60 | 23 | 14 | 2 | 3 | 1 |
| 25 | in3380_6 | 89 | 72 | 93 | 39 | 48 | 43 | 64 | 48 | 64 | 32 | 24 | 7 | 5 | 10 |
| 26 | in3380_7 | 86 | 72 | 80 | 45 | 50 | 46 | 66 | 49 | 65 | 36 | 14 | 1 | 1 | 3 |
| 27 | in3380_8 | 92 | 77 | 88 | 42 | 51 | 45 | 74 | 54 | 64 | 31 | 18 | 3 | 4 | 4 |

Columns chr___: 1254 CpG sites

Columns y:
- 0: normal sample,
- 1: insignificant tumor,
- 2: significant tumor

Rows: 334 samples

# What is feature/variable selection?

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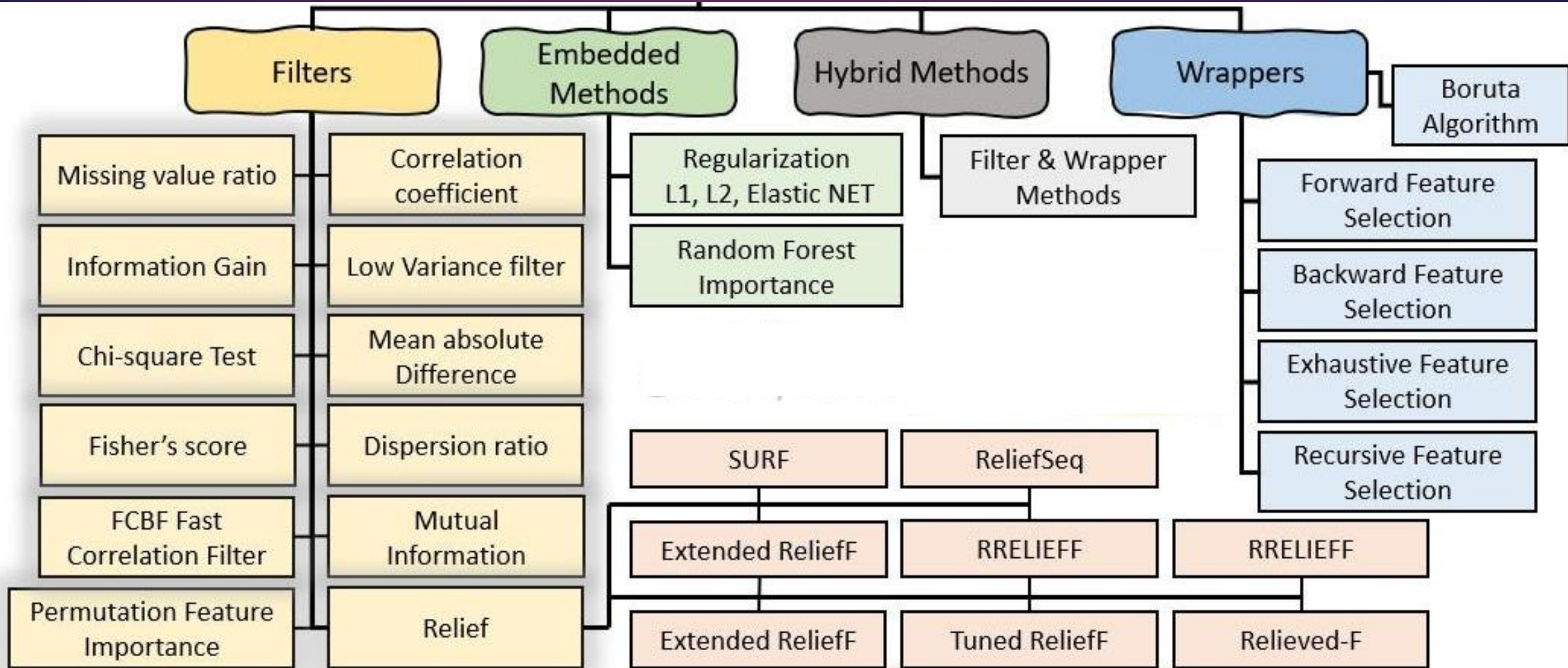
► Give:



► 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.

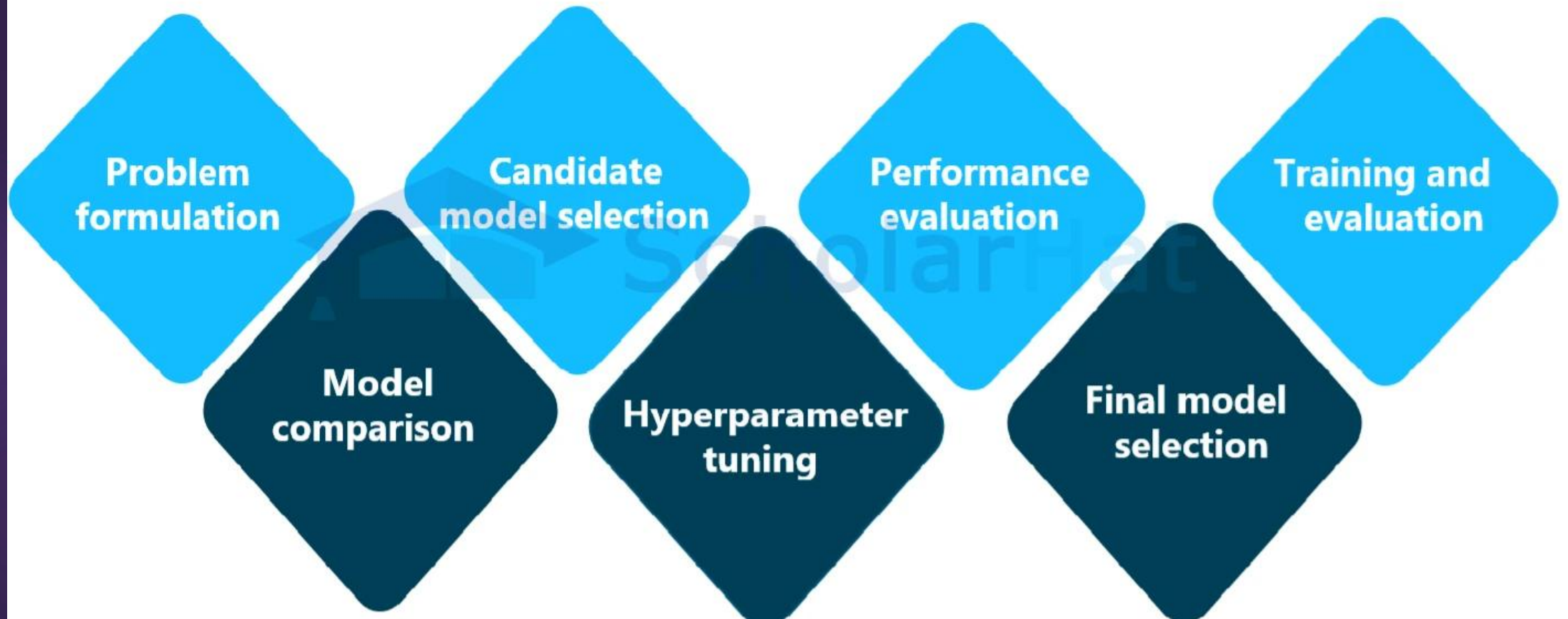# How do we carry out feature selection?
## Supervised Feature Selection

# Model Selection vs Feature Selection

▶ Feature Selection is a part of Model Selection

- **Selecting features (or basis functions)**
  - Linear regression
  - Logistic regression
  - SVMs
- **Selecting parameter value**
  - Prior strength
    - Naïve Bayes, linear and logistic regression
  - Regularization strength
    - Naïve Bayes, linear and logistic regression
  - Decision trees
    - depth, number of leaves
  - Boosting
    - Number of rounds
- More generally, these are called **Model Selection** Problems

# Model selection: steps

# Model selection: steps

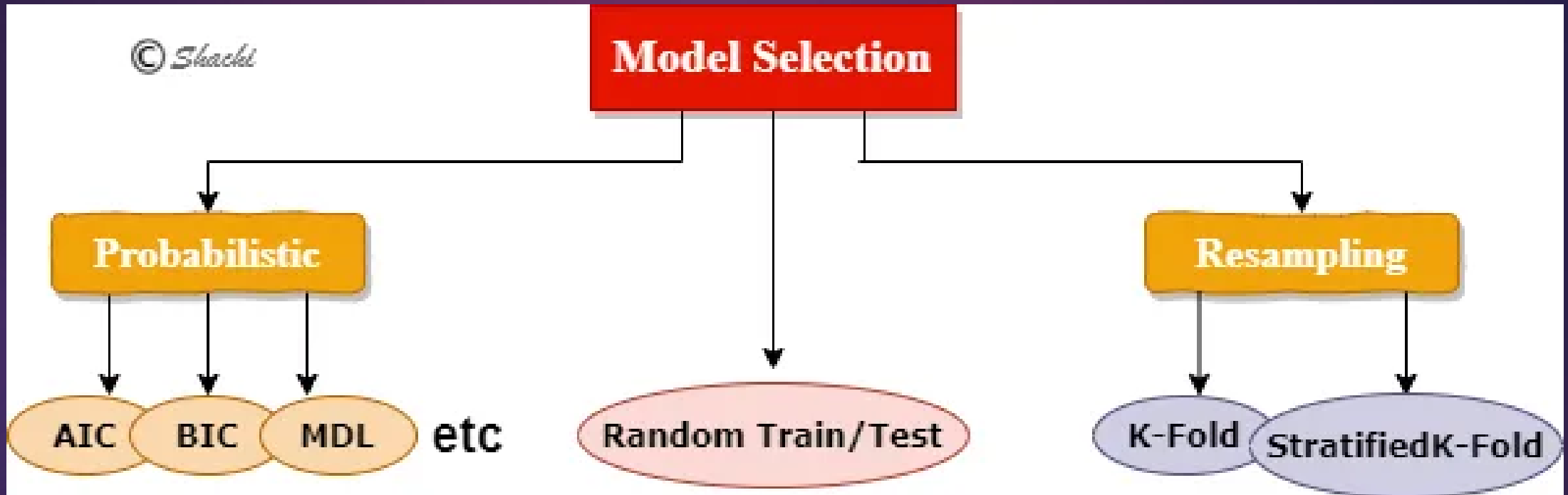Stage 1: Selecting the regression model forms

Stage 2: Selecting the regression model and the independent variables

Stage 3: Fitting the model

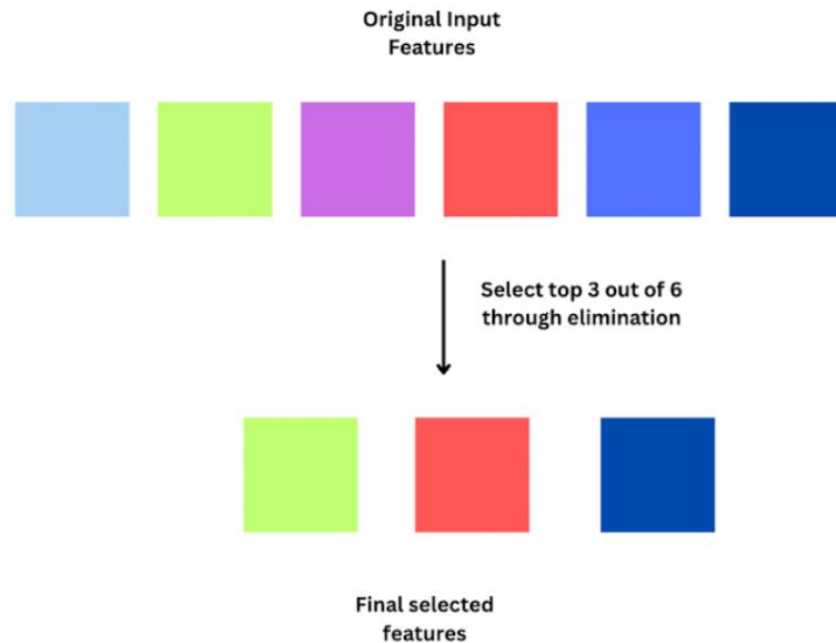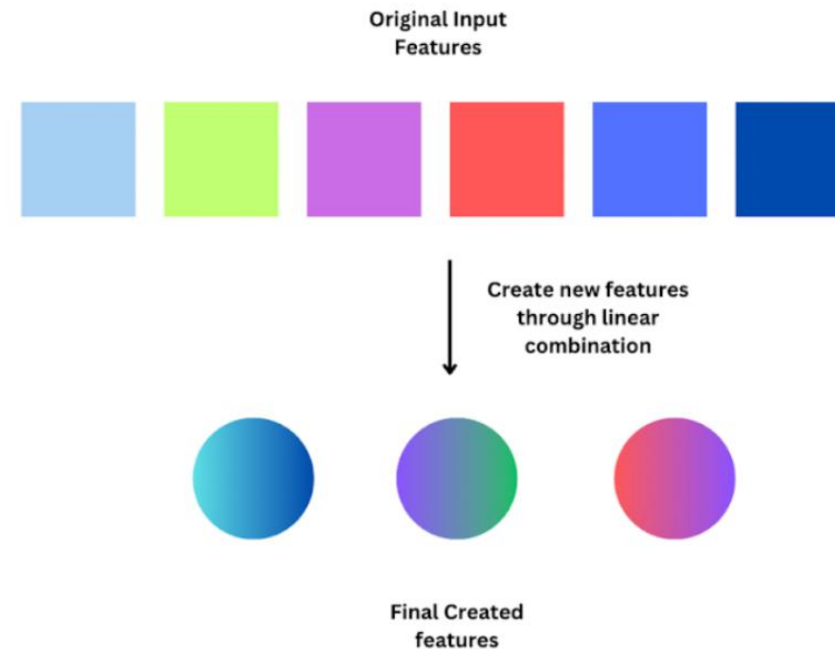Stage 4: Examining or validation of the applied model

# Model selection: methods

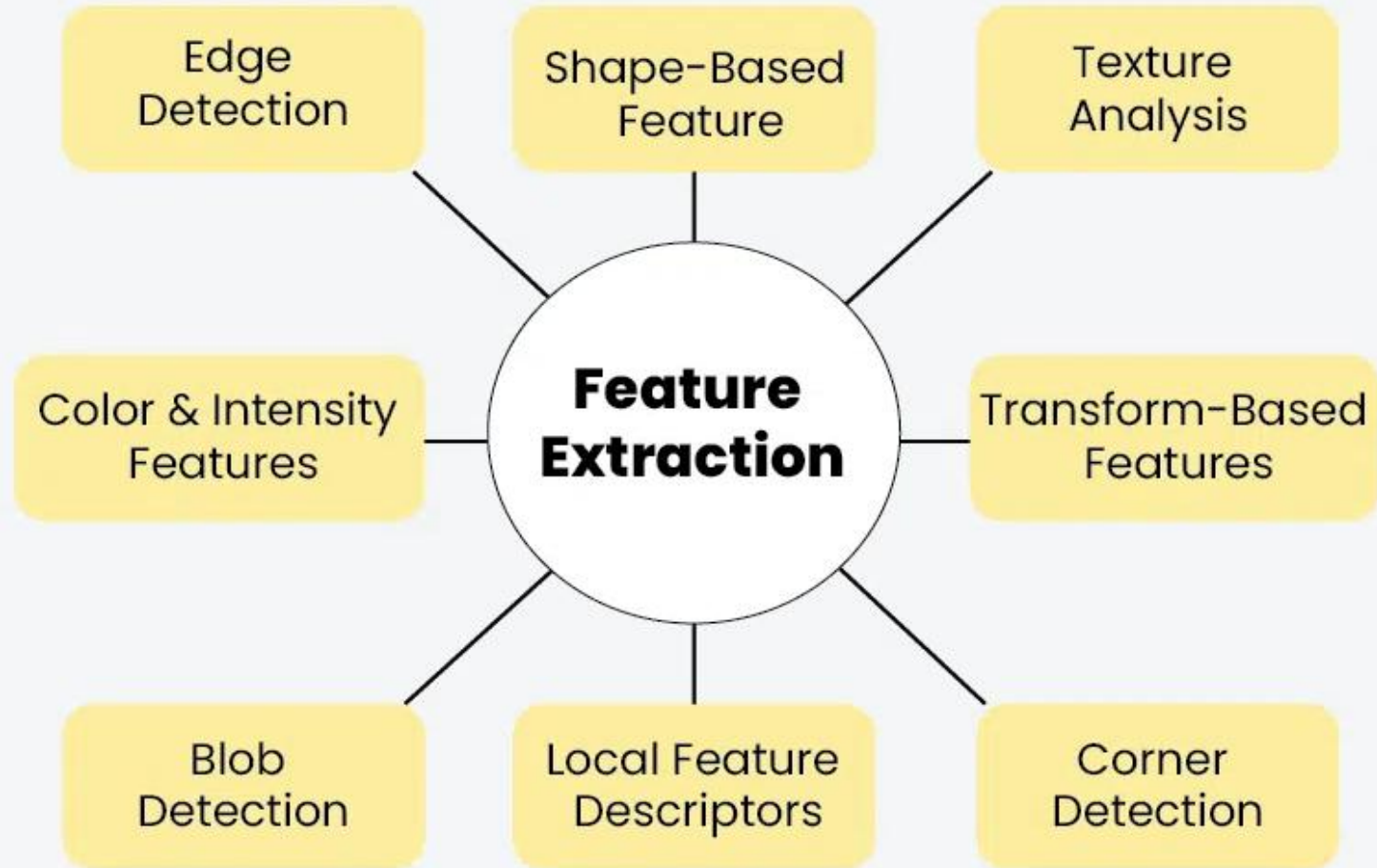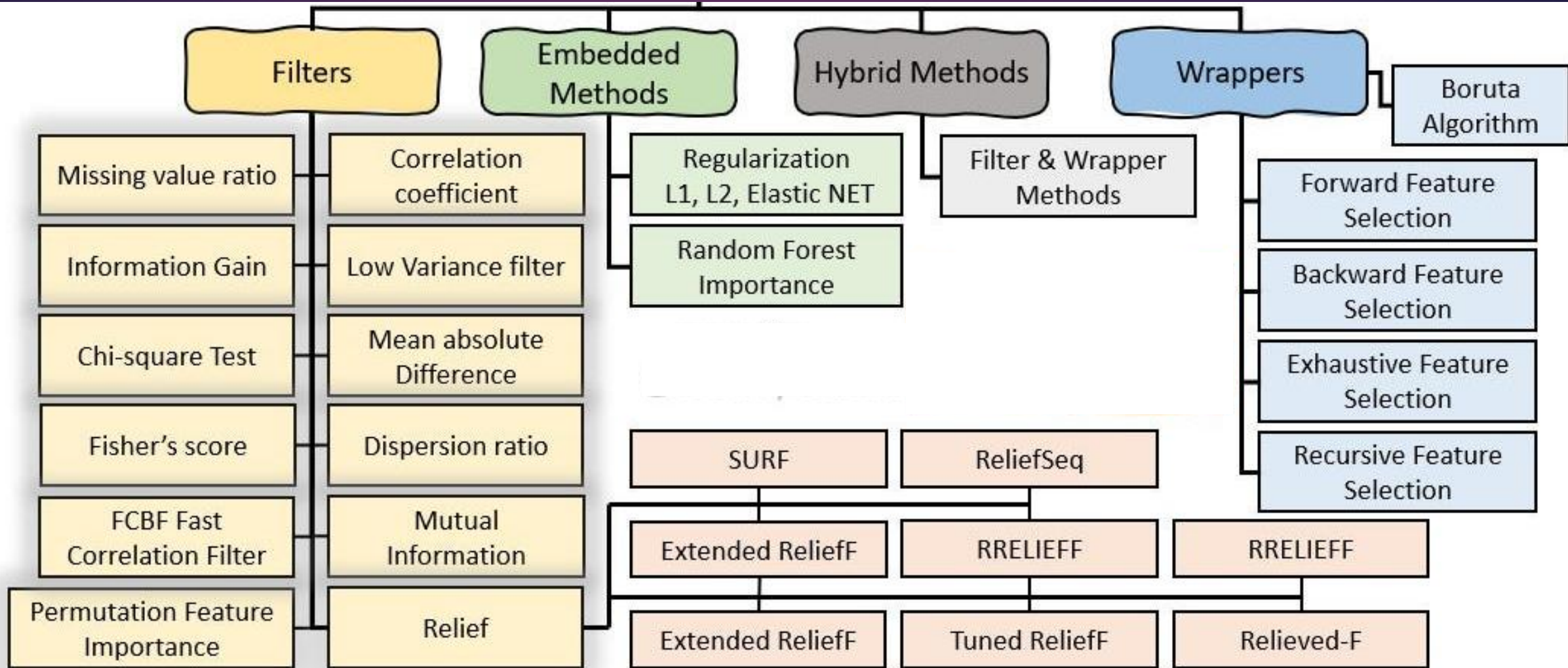# Feature Selection vs Feature Extraction/Engineering

# Feature Extraction/Engineering

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

# Xin chân thành cảm ơn!

LUU PHUC LOI, PHD

ZALO: 0901802182

LUU.P.LOI@GOOGLEMAIL.COM