# MACHINE LEARNING

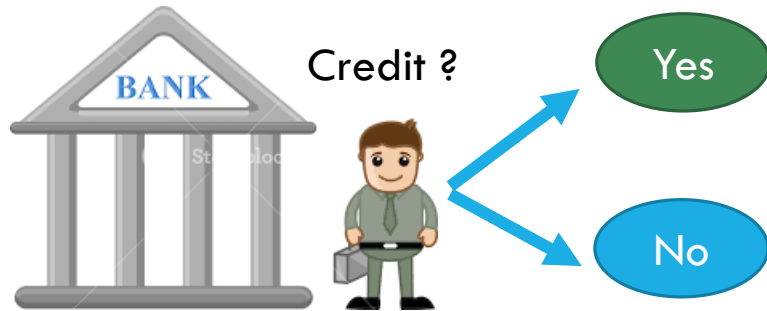# INTRODUCTION TO PYTHON FOR STATISTICAL LEARNING

Fathi Abdul Muhyi

# MACHINE LEARNING

"Sebuah Algoritma yang dapat mempelajari karakteristik data, membuat prediksi, menggali informasi berdasarkan data"
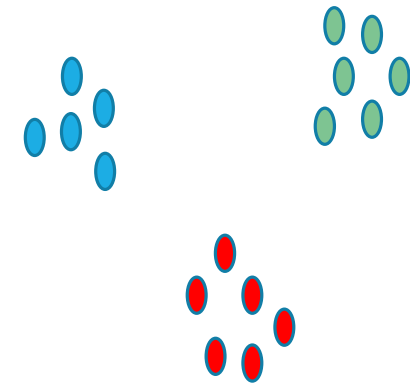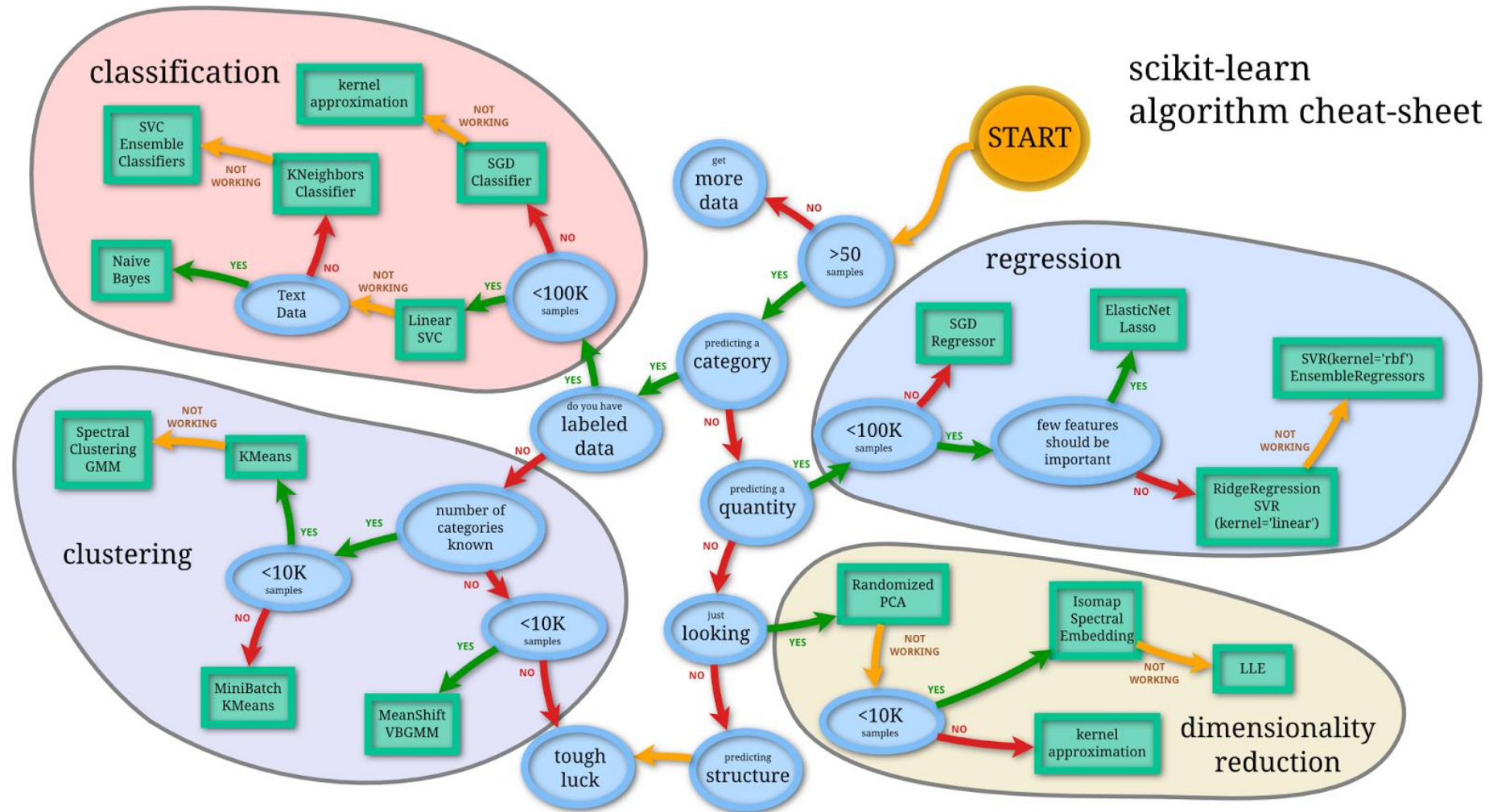
# PENERAPAN MACHINE LEARNING

**Classification**

**Regression**

**Clustering**

BANK

Credit ?

Yes

No

DIJUAL!!!

*How Much ?*

# ML FLOWCHART



scikit-learn
algorithm cheat-sheet

# METODE KLASIFIKASI

Classification is a method to build a model to classify the response variable by certain charcteristic (feature variable).

There are many problem in case of classification :

1. Credit scoring, calon penerima kredit mampu bayar atau tidak

2. Churn, customer yang berpotensi meninggalkan jasa/produk

3. Direct Marketting, identifikasi prospective customer

# CONTOH KASUS SEDERHANA

Toko X mengumpulkan data pelanggan untuk mengetahui apakah pelanggan mereka tertarik untuk membeli produk jenis baru

Toko tersebut memiliki sejumlah pelanggan baru yang belum diketahui tertarik atau tidak

Peubah Penjelas

Peubah Respon

| | No | Jenis Kelamin | Single | Tinggal di Kota | usia | Perokok | Budget | Kesukaan | Tertarik Beli? |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 1 | 32 | 0 | low | Tekno | 0 |
| 1 | 2 | 0 | 1 | 0 | 38 | 0 | medium | Tekno | 0 |
| 2 | 3 | 0 | 0 | 0 | 33 | 0 | low | Tekno | 0 |
| 3 | 4 | 0 | 1 | 0 | 27 | 0 | medium | Lainnya | 0 |
| 4 | 5 | 1 | 1 | 1 | 30 | 0 | medium | Busana | 1 |
| 5 | 6 | 0 | 1 | 0 | 44 | 0 | medium | Tekno | 0 |
| 6 | 7 | 1 | 1 | 1 | 36 | 0 | medium | Seni | 1 |
| 7 | 8 | 1 | 0 | 0 | 32 | 0 | low | Seni | 0 |
| 8 | 9 | 0 | 1 | 0 | 31 | 0 | medium | Seni | 0 |
| 9 | 10 | 1 | 1 | 0 | 40 | 0 | high | Tekno | 0 |
| 10 | 11 | 1 | 0 | 1 | 34 | 0 | low | Lainnya | 0 |

# REGRESI LOGISTIK BINER

Memprediksi peluang ketertarikan untuk membeli

$$P(Y = 1) = \frac{\exp(\beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k)}{1 + \exp(\beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k)}$$
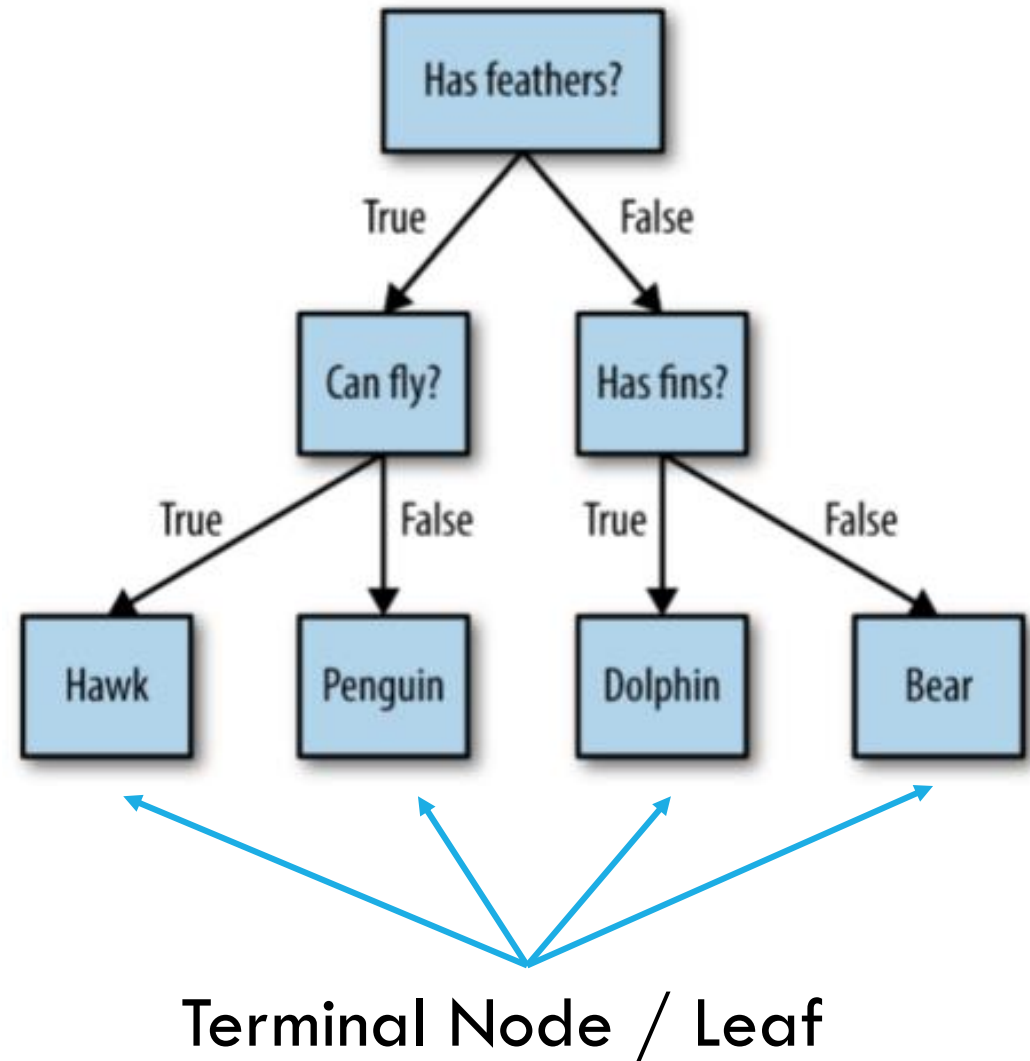
$Y$      : peubah kategorik biner

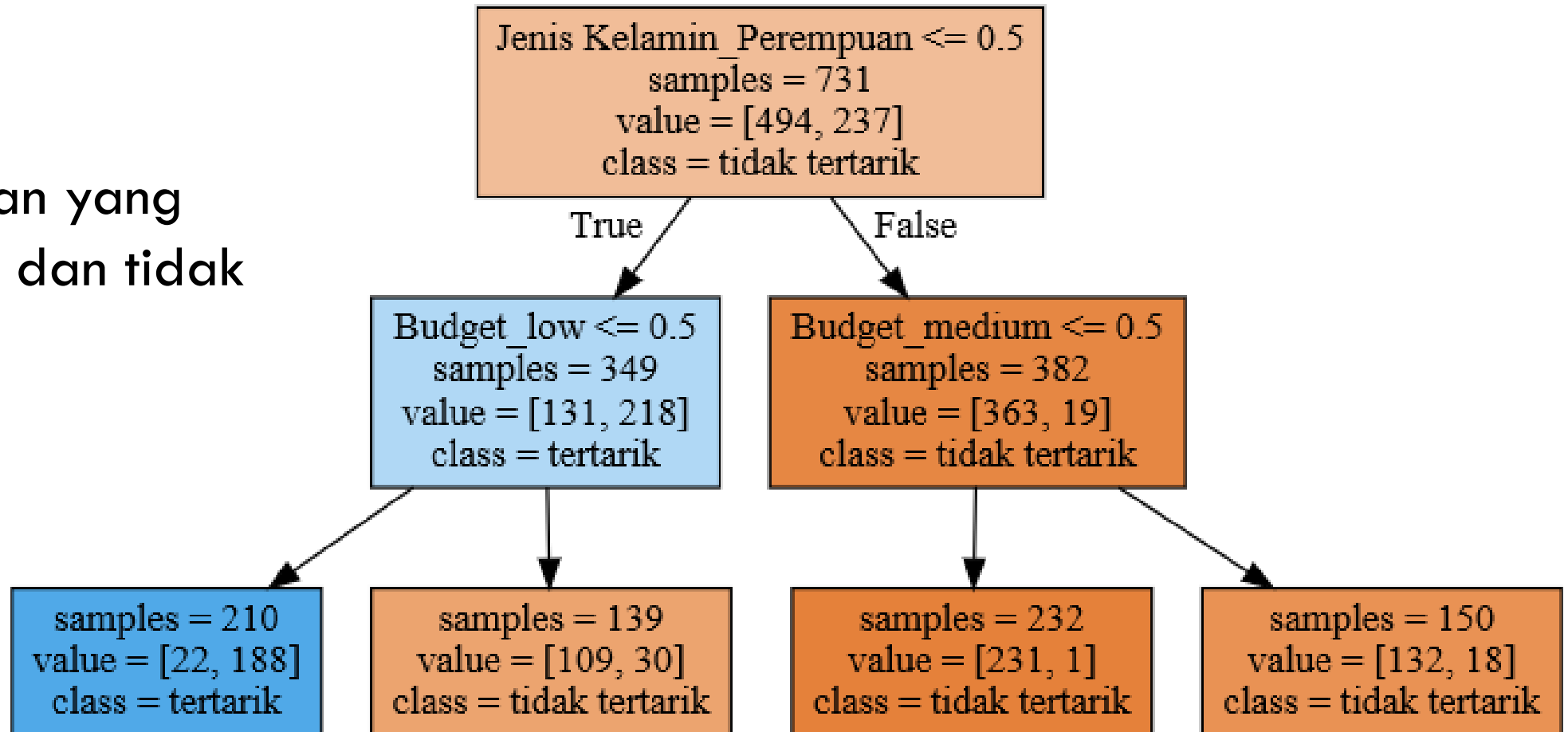$x_i$      : peubah penjelas

$\beta_i$      : Parameter Regresi

# POHON KEPUTUSAN

Hierarchi of if/else question, leading to a decision (Muller and Guido, 2016)
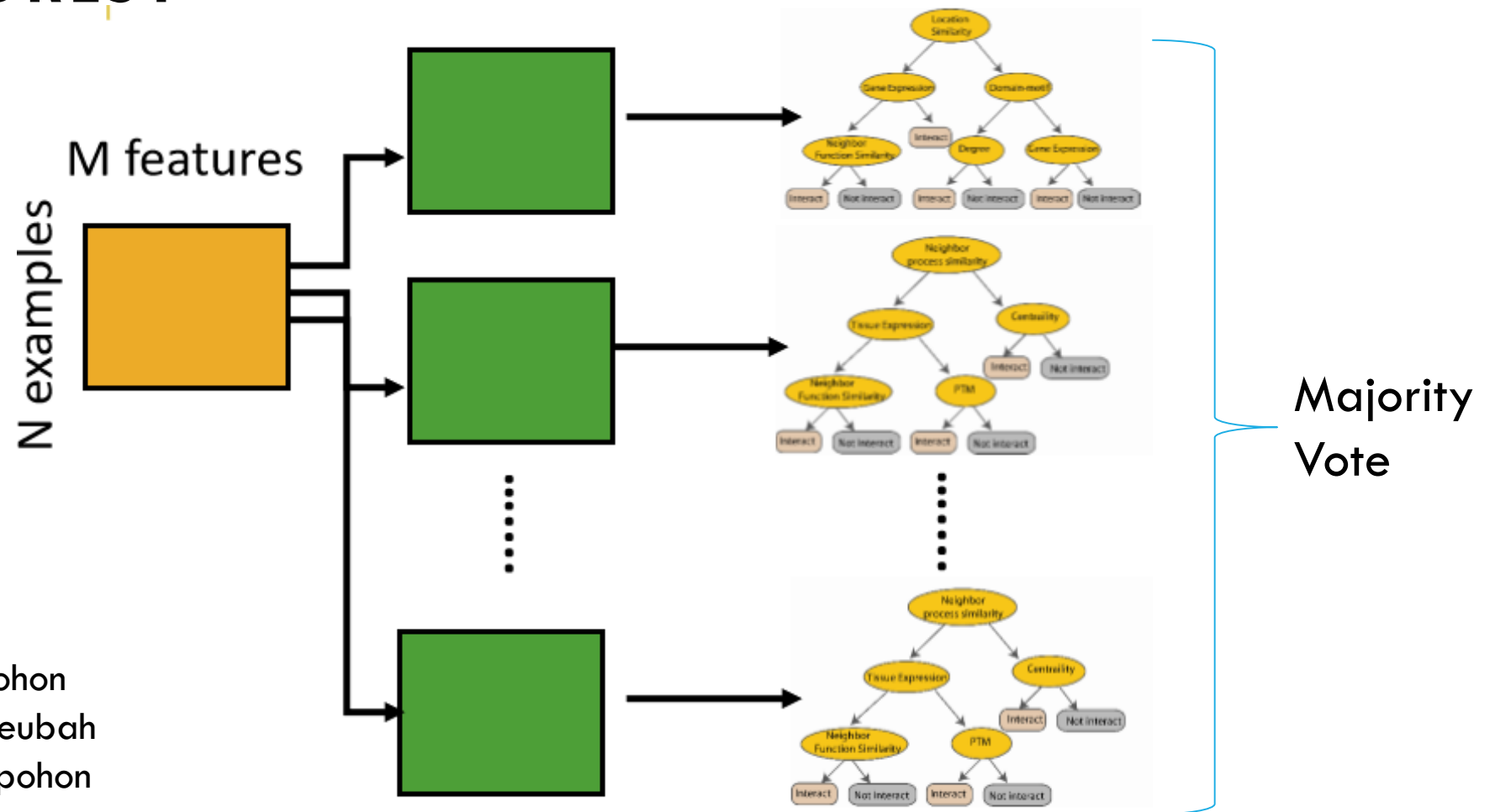
How to distinguish :
Hawk, Penguin, Dolphin, Bear ?



Terminal Node / Leaf

# POHON KEPUTUSAN

Membedakan yang tertarik beli dan tidak

# RANDOM FOREST



**M features**

**N examples**

**Majority Vote**

Random Forest :
1. Tree cenderung overfit dan keragamannya tinggi
2. Menggabungkan banyak pohon
3. Pengambilan contoh acak peubah dan observasi untuk setiap pohon

# EVALUASI

Ukuran Ketepatan Prediksi :

 1. Akurasi, persentase yang terprediksi benar

 2. Sensitivitas/Recall, persentase yang terprediksi benar pada kelas positif

 3. Spesifisitas, persentase yang terprediksi benar pada kelas negatif

 4. dan lain-lain

Metode Evaluasi :

1. Validasi (Train - Test)

2. Validasi Silang

**Table 1** Confusion Matrix

| Prediction | Aktual | |
|---|---|---|
| | + | - |
| + | True Positive (A) | False Positive(B) |
| - | False Negative (C) | True Negative(D) |

**Table 2** Metrics of the classification

| Metrics | Information |
|---|---|
| Accuracy | (A + D) / (A+B+C+D) |
| Recall (postive) | A / (A+C) |
| Recall (negative) | D / (B+D) |
| Precision (Positive) | A / (A+B) |
| F1 score | Weighted Average between precision and recall |
| Geometric mean score | $(Accuracy*Recall(+)*Recall(-))^{1/3}$ |

# METRICS

# CONTOH

| No | Perdiksi | Aktual |
|----|----------|--------|
| 1 | 1 | 1 |
| 2 | 1 | 0 |
| … | … | … |
| … | … | … |
| 553 | 0 | 0 |
| 554 | 0 | 1 |

Matriks Klasifikasi

| Prediksi | Aktual | |
|----------|--------|--------|
| | 1 | 0 |
| 1 | 121 | 23 |
| 0 | 10 | 400 |

Akurasi  :
94.04 %
Sensitivitas :
92.36 %
Spesifisitas :
94/56 %

# VALIDASI

# VALIDASI SILANG



Model Fitting

Model Fitting

Model Fitting

Model Fitting

CV Train set

CV Test set

# TASK

Prediksi variabel default pada data titanic dan upload di kaggel.com