# Imbalance

① Split ⟶ Stratify = y

② Set Suitable Metric : Recall, precision, f1-score

$$\underline{\underline{\quad}}$$

Auc Score

⟶ ③ use suitable handling (Resampling)

↳ SMOTE / oversample / undersample

④ use Decision Tree or Ensemble Methods

(less Sensitive for imbalance)

⑤ use class.weight (if applicable)

Majority

0 : 500

1 : 5

} Imbalance

1 : 100

1 : 3   (light)

Minority

① underSampling → 

0 : 5   (Random Sampling)

1 : 5

② overSampling → 

0 : 500

1 : 500   (Resampling with Replacement)

③ SMOTE → 

0 : 500

1 : 500   (Synthetic Oversampling)

\* **Ensemble Methods :**

Random forest → Bagging
      ↳ Base : DT

- LR, SVR, Knn, DT → Single learner (weak)

- Ensemble → Strong learner

① Voting   Same Dataset    ↪ Same Model   ② Bagging   Diff Subset ← ③ Boosting
      Diff Models         of Data

[SVR, LR, DT]

→Reg :   Avg

→Clf : ⟶ hard voting ⟶ voting
    ↳ Soft ~ → Avg (prob)

↪ parallel

↪ Random forest

↪ Sequential

⇒ Xgboost (learning-rate)

⇒ Gradient Boosting

\* Dimension Reduction

$(\#Rows, \#Cols)$

↑ observations

↑ features

- Curse of Dimensionality

① Stonge

② Time / Response

③ Visualization

feature Selection

feature Extraction

Compression

indicator for information

Age   inCome

original Space  $\boxed{f_1 \; f_2 \; f_3 - f_{100}}$

⇓ Linear Transformation

New Space  $\boxed{PC_1 \; PC_2 \; PC_3 \; \dots \; PC_{100}}$

ordered by variance

Example   50%   20%   15% → Select  $\boxed{PC1, PC2, PC3}$  85%