

PS: Pre-Learning .

<I> Motivation: Prediction & Inference.

- **Prediction** → Mainly in the predicted valuable as result of input.
<Result> But not each input.
- **Inference** → Mainly in the way each one of the input affect.
<Process>

<II> Types:

- **Regression** → Continuous. ➤ all is to predict, but the type of results are different. And Both are supervised problem.
- **Classification** → Discrete.
- **Clustering** → A set of inputs is to be divided into groups, which has the same properties.
- **Density Estimation** → The distribution of some space. Generally, it supposes that the samples are obey a special distribution D in some hypothetic space.

<III> Kind. (Type).

- **Parametric** → Based on the inputs. I can select or design a model which has been consisted by several unknown parameters. And I need to tune it to find which params. are most suited.
- **Non-Parametric** → It is not given by params. but it doesn't mean that it doesn't need params. It just need a large number of observations to obtain an accurate estimate for f .

<IV> Categories:

- **Supervised** → The goal is to learn a general rules that $\text{input} \rightarrow \text{output}$
ex: classification, regression
- **Un-supervised** → No label are given to the learning algorithm.
ex: clustering

ex: clustering.

- Reinforcement Learning → The program is provided feedback in terms of rewards and punishments as it navigates its problem space.

(V). Performance Analysis

• Confusion Matrix :

		True condition		Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Total population	Condition positive	Condition negative			
Predicted condition	Predicted condition positive	True positive	False positive (Type I error)	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative (Type II error)	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
	True positive rate (TPR), Recall, Sensitivity, probability of detection $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) $= \frac{\text{LR}^+}{\text{LR}^-}$	F ₁ score = $\frac{2}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}}$
	False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	True negative rate (TNR), Specificity (SPC) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$		

- Error Rate: The percentage which the samples classified false in the whole samples.
- Accuracy: The percentage which the ones classified true in the all.

		TRUE CONDITION	
		Condition Positive	Condition Negative
Predicted Positive	TRUE POSITIVE	FALSE POSITIVE	TRUE NEGATIVE
	FALSE NEGATIVE	TRUE NEGATIVE	

• Accuracy (AC): $\frac{\sum \text{TP} + \sum \text{TN}}{\sum \text{Total Samples}}$. present the model's Accuracy.

• Precision : $P = \frac{\text{TP}}{\text{TP} + \text{FP}}$ "查准率"

• Recall : $R = \frac{\text{TP}}{\text{TP} + \text{FN}}$ "查全率"

• F₁ score : $\frac{2 \times P \times R}{P + R} = \frac{2 \times \text{TP}}{\sum \text{Total} + \text{TP} - \text{TN}} = \frac{2}{\frac{1}{P} + \frac{1}{R}}$

→ For important one of it, the f₁ score can be described as :

$$F_1 \text{ score} = \frac{(1 + \beta^2) \times P \times R}{(\beta^2 \times P) + R}$$

• ROC & AUC

ROC is for research the ability of generalization.

• ROC & AUC

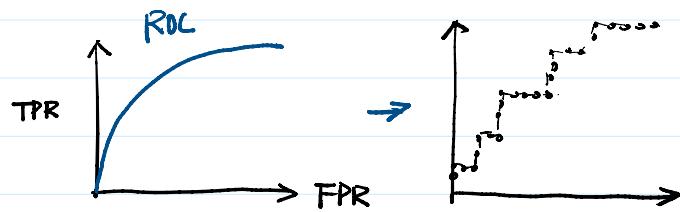
ROC is for research the ability of generalization.

□ TPR (TRUE POSITIVE RATE)

$$= \frac{TP}{TP + FN}$$

□ FPR (FALSE POSITIVE RATE)

$$= \frac{FP}{TN + FP}$$



Problem: Realization. - the ROC can't be so smooth like this. \rightarrow AUC

该图(ROC)过程:

STEP ①: 给定 m^+ Positive and m^- NEGATIVE

STEP ②: 把步类阈值设到最大, \rightarrow 把所有样例 \rightarrow 反例. $\rightarrow (0,0)$.

STEP ③: 把步类阈值依次设为每个样例的预测值。

若为真例, $\rightarrow (x, y + \frac{1}{m^+})$ 若为反例, $\rightarrow (x + \frac{1}{m^-}, y)$

AUC (Area of Under ROC Curve).

$$AUC = \frac{1}{2} \cdot \sum_{i=1}^{m-1} (x_{i+1} - x_i) \cdot (y_i + y_{i+1}) \quad (\text{由来: 梯形面积的表示})$$

• MSE (Mean Squared Error)

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

Measures the average of the squares of the error or deviations — that is, the difference between the estimator and what is estimated.