① Accuracy

| Actual | Predicted | |
|--------|-----------|---|
| 1 | 0 | ✗ |
| 0 | 1 | ✗ |
| 1 | 1 | ✓ |
| 1 | 0 | ✗ |
| 0 | 1 | ✗ |
| 1 | 0 | ✗ |
| 0 | 0 | ✓ |
| 0 | 0 | ✓ |

Same logic for binary & multi-class classification problem

Accuracy score = $\frac{\text{correctly predicted}}{\text{total } \# \text{ of preds}}$   [0,1] ↑ better

$= \frac{3}{9} = \frac{1}{3}$

$= 33\%$

✳ Same logic for multi-classification

┌─────────────────────────────────────────┐
│ ✳ Always ask 3 questions b4 selection of │
│ measure in classification problems.      │
│                                          │
│                        ⟋ Balanced        │
│        1) Dataset ⟨                       │
│                        ⟍ Imbalanced      │
│                                          │
│                        ⟋ Binary          │
│        2) Problem ⟨                       │
│                        ⟍ Mult-class      │
│                                          │
│        3) Which error are we more concerned │
│                   about?                 │
└─────────────────────────────────────────┘

✳ How much accuracy is good?

→ Depends on the problem we are dealing with & the implications / consequences the inaccuracy will lead to.

90% accuracy might sound good while predicting if a customer will renew subscription based on the current month's activity but the same accuracy when predicting the presence of cancer by looking at x-ray is very bad as it means for every 100 patients there a chance 10 people may die.

✳ Accuracy doesn't tell us the nature of the mistake.

i.e.   actual  predicted
        1        0      → How many
        0        1      → How many

② Confusion matrix

|  | Predicted 1 (+) | Predicted 0 (−) |
|--------|--------|--------|
| (+) 1 | TP ✓ | FN [Type II] |
| (−) 0 | FP [Type I] | TN ✓ |

Actual

✳ We can extract the accuracy score from the confusion matrix

accuracy score = $\frac{TP+TN}{TP+TN+FP+FN}$

MCR = 1 − accuracy score

```
from sklearn.metrics import confusion_matrix, classification_report
```

⤷ confusion_matrix (true, observed)

Note: Not choosing the right evaluation metrics can mis lead your judgement.

Accuracy vs Imbalanced dataset

→ Assume we are building a terrorist classification system

|  | Predicted |  |
|--------|--------|--------|
| Actual 1 | 0 | 1 |
| 0 | 0 | 99999 |

accuracy = $\frac{99999}{99999+1} = 99.99\%$

✳ Accuracy can be misleading when dealing with imbalanced dataset.

It is obvious if a model is trained on Balanced Dataset, most chances predictions will also be balanced.
Actual # of Positives ≈ Predicted # of Positives
Actual # of Negatives ≈ Predicted # of Negatives
TP ≈ TN
FP ≈ FN

③ Precision [what proportion of predicted positives is truly positive]

|  | Predicted 1 | 0 |
|--------|--------|--------|
| Actual 1 | TP | FN |
| 0 | FP | TN |

precision = $\frac{TP}{TP+FP}$   [0,1] ↑ better

```
from sklearn.metrics import recall_score, precision_score, f1_score
```
average = None

✳ There is a tradeoff b/w Precision & Recall.

④ Recall [what proportion of actual positives is correctly classified] [True Positive Rate] [Sensitivity]

|  | Predicted 1 | 0 |
|--------|--------|--------|
| Actual 1 | TP | FN |
| 0 | FP | TN |

Recall = $\frac{TP}{TP+FN}$   [0,1] ↑ better

$\frac{1060}{1200}$  $\frac{1030}{1200}$

* eg cancer

| high threshold | low threshold |
|--------|--------|
| TP↓ FP↓↓ | TP↑ FP↑↑ |
| TN↑ FN↑↑ | TN↓ FN↓↓ |
| P↑ | P↓ |
| R↓ | R↑ |

Choosing threshold — w.r.t. Precision } → Precision & Recall vs Threshold (P)
— w.r.t. recall } → Precision vs Recall curve (P) (R)
— w.r.t. f1-score } → F1 score vs threshold (f) (t)

we can use AUPRC for model comparison also. ↑ better

✳ Whether to use precision or recall depends on if the Type I or Type II error has more implications / consequences.

[FP]        [FN]
[Type I]    [Type II]

⑤ F1 score [not sure if Type I or Type II error is more problematic? use F1 score]

F1 score = $\frac{2PR}{P+R}$ → Penalizes the lower value [Harmonic mean]

✳ In Binary classification for Precision & Recall, we discuss in terms of 1/positive

Multi-class Precision & Recall

|  | Predicted |  |  |  |
|--------|--------|--------|--------|--------|
|  | Setosa | versicolor | Virginica | Actual total |
| Setosa | 15 | 5 | 10 | 40 |
| versicolor | 0 | 20 | 4 | 24 |
| virginica | 6 | 10 | 20 | 34 |
| predicted total | 24 | 15 | 34 | 108 |

Actual

$F1_s = 2P_sR_s / (P_s + R_s)$
$F1_{ver} = 2P_{ver}R_{ver} / (P_{ver} + R_{ver})$ ⟍ 'macro'
$F1_{vir} = 2P_{vir}R_{vir} / (P_{vir} + R_{vir})$ ⟍ 'weighted'

$P_s = 15/24$   $R_{ver} = 20/15$   $P_{vir} = 20/34$     $R_s = 15/40$   $R_{ver} = 20/24$   $R_{vir} = 20/34$

[balanced]          [imbalanced]           'macro recall'        'weighted recall'
'macro precision'    'weighted precision'    [balanced]           [imbalanced]

$\frac{P_s + P_{ver} + P_{vir}}{3}$     $\frac{40}{108}P_s + \frac{24}{108}P_{ver} + \frac{34}{108}P_{vir}$     $\frac{R_s + R_{ver} + R_{vir}}{3}$     $\frac{40}{108}R_s + \frac{24}{108}R_{ver} + \frac{34}{108}R_{vir}$

Recommendation :- Always label ⇒ majority-class → 0-class ✓
                                  minority-class → 1-class ✓
                          ↓
So that, will not confuse in selecting the performance measure

Predictions
|  | 0 | 1 |
|---|---|---|
| Actual 0 | TN | FP |
| 1 | FN | TP |

$$TNR \ (specificity) = \frac{TN}{TN + FP}$$

Predictions
|  | 0 | 1 |
|---|---|---|
| Actual 0 | TN | FP |
| 1 | FN | TP |

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

Matthews correlation coefficient
→ Symmetric metric for imbalanced dataset
→ measure of the quality of the classification

→ MCC is a measure of association for two binary variables

Actual      Prediction
  0            0
  0  ⟨——⟩    0
  ⋮            ⋮

Range $[-1, +1]$

→ $MCC = \dfrac{TP \times TN - FP \times FN}{\sqrt{(TP+FN)(TP+FP)(TN+FP)(TN+FN)}}$

} → cares all counts
   ie. TP, TN, FP, FN
   more symmetric than F1-score

} → just to make output within $[-1,+1]$
   Normalization constant, always +ve

$\left. \begin{array}{l} FP=0 \\ FN=0 \end{array} \right\}$ → perfect +1 ⎤ Perfect prediction

$(TP+TN) = (FP+FN)$ } → very bad 0

$\left. \begin{array}{l} TP=0 \\ TN=0 \end{array} \right\}$ → fault in model -1 ⎦ Inverse prediction

→ Use cases ⟨ Imbalance Dataset
             all classes are equally important { like in F1-score }

→ MCC is more informative than F1-score in evaluating binary classification problems.
   ↳ → B/c it takes into account the balance ratios of 4 confusion matrix cells.

$$MCC = (Pos_{precision} + Neg_{precision} - 1) \times PosNegRatio$$

Where each element of the formula is:

$$Pos_{precision} = \frac{TP}{TP + FP}$$

$$Neg_{precision} = \frac{TN}{TN + FN}$$

$$PosNegRatio = \frac{PosPredictionCount \times NegPredictionCount}{PosLabelCount \times NegLabelCount}$$

$$PosPredictionCount = TP + FP$$

$$NegPredictionCount = TN + FN$$

The reformulation helps to clarify, in a more intelligible form than the original, that you can get higher performance from improving both positive and negative class precision, but that's not enough: you also have to have positive and negative predictions in proportion to the ground truth, or your submission will be greatly penalized.

* MCC doesn't depend on which class is the positive one

Multi-class classification Performance measures

1 ⟨ 0 as 1
     1 as 0

In Multi-class : (class) Choose which class we are more concerned about → (error) Choose which error we are more concerned about ⟨ Type 1 — Precision / Type II — Recall / Both — F1 score
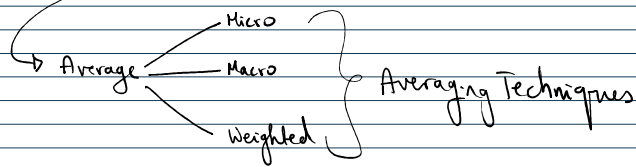
→ In multi-class problems, we extend the formulas of Binary class.

→ We are dealing with multiple classes, so calc concerned score for each class & then take the Average.

→ Average
  - Micro
  - Macro          } Averaging Techniques
  - Weighted

Macro → Compute performance for each class & then average



Micro   [not suitable for imbalanced dataset]





| Micro Recall = Micro Precision = Micro F1 score = Accuracy |

  └→ Micro Averaging holds same properties as Accuracy

  └→ Global Metric

  └→ Not good measure when classes are not balanced

  → Micro avg gives equal weight to each sample

  → Macro avg gives equal weight to each class

  → If we have same # of samples for each class, both
     macro & micro will provide the same score.



  → The metric is very intuitive & easy to understand.
     Both in binary cases & multi-class cases the Accuracy
     assumes values b/w 0 & 1, while the quantity
     missing to reach 1 is called misclassification rate.

- summarizing the two main steps of Balanced Recall, first we compute a measure of recall for the algorithm on each class, then we apply the arithmetic mean of these values to find the final Balanced Recall score.

- All in all, Balanced Recall consists in the arithmetic mean of the recall of each class, so it is "balanced" because every class has the same weight and the same importance.

- A consequence is that smaller classes eventually have a more than proportional influence on the formula, although their size is reduced in terms of number of units. This also means that Balanced Score is insensitive to imbalanced class distribution and it gives more weight to the instances coming from minority classes. On the other hand, Accuracy treats all instances alike and usually favours the majority class.

- Balanced Score: This may be a perk if interested in having good prediction also for under-represented classes, or a drawback if we care more about good prediction on the entire dataset.

- The smallest classes when misclassified, are able to drop down the value of Balanced score, since they have the same importance as largest classes have in the equation.

- The class presents a high number of individuals, its bad performance is caught up also by the Accuracy.

- Instead, the class has just few individuals, the model's bad performance on this class cannot be caught up by Accuracy.

- If we are interested in achieving good predictions also for rare classes, the information of Balanced Score guarantees to spot possible predictive problems also for the under-represented classes.
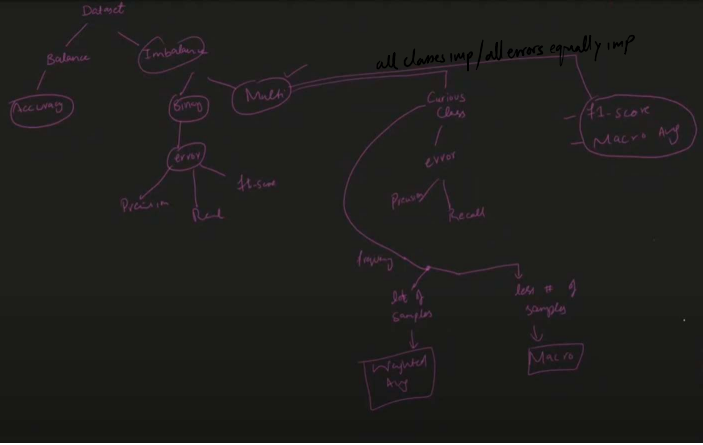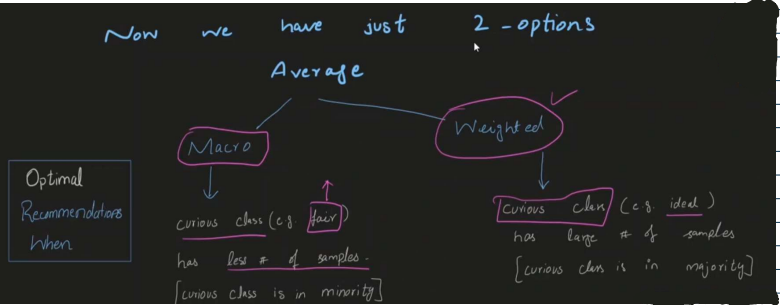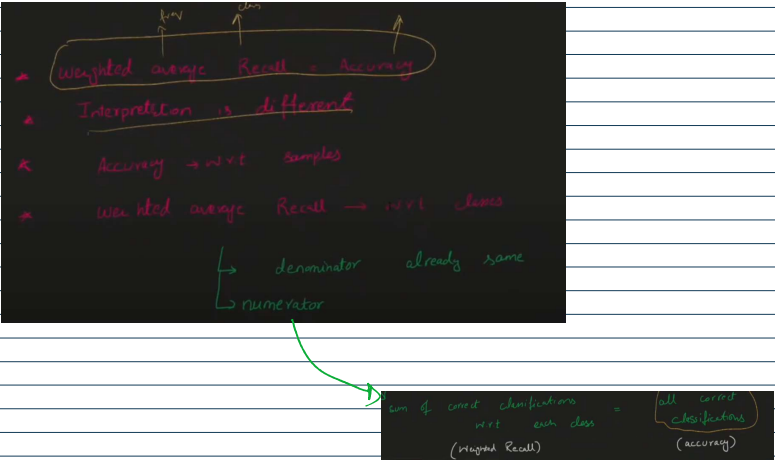
**Weighted :** compute performance for each class & then take weighted average

```
y_test.value_counts()

Ideal     5388
Good      4247
Premium   3488
Fair       402
Name: cut, dtype: int64
```

weighted Recall = 0.57

- The **Balanced Score Weighted** takes advantage of the Balanced Score formula multiplying each performance(i.e., recall) by the **weight of its class,** namely the frequency of the class on the entire dataset. We add also the sum of the weights at the denominator, with respect to the Balanced Score.

- Balanced score weighted holds the **goodness of both accuracy and balanced score metrics.**

- Once recalls have been weighted by the frequency of each class, **the average of recall is no longer dirtied by low frequency classes(unlike balanced score):** large classes will have a proportional weight to their size, and small ones will have a resized effect, compared with the Balanced Score formula.

- Since every recall is weighted by the class frequency of the initial dataset, Balanced Recall Weighted could be a good performance indicator when the aim is to train a classification algorithm on a wide number of classes.

- In fact, this metric **allows to keep separate algorithm performances on the different classes,** so that we may track down which class causes poor performance(like balanced score). At the same time, it keeps track of the importance of each class thanks to the frequency(unlike balanced score).

- This ensures to obtain a reliable value of the overall performance on the dataset: we may interpret this metric as the **probability to correctly predict a given unit,** even if the formula is slightly different from the Accuracy.

- Weighted average Recall = Accuracy
- Interpretation is different
- Accuracy → w.r.t. samples
- weighted average Recall → w.r.t. classes
  ↳ denominator already same
  ↳ numerator

sum of correct classifications w.r.t. each class = all correct classifications
(weighted Recall)                                    (accuracy)

Now we have just 2 - options

Average

Macro                    Weighted

curious class (e.g. Fair) has less # of samples.
[curious class is in minority]

curious class (e.g. ideal) has large # of samples
[curious class is in majority]

Optimal Recommendations when

all classes imp/all errors equally imp?

F1-score Macro Avg

```
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=, stratify=y)
```