

EVALUATION MEASURES

(24) Bootstrapping and Cross Validation

Bootstrapping: resampling technique used to estimate stats (like mean, variance, accuracy) or build models when you have limited data.

We repeatedly sample {with replacement} from the dataset to create multiple bootstrap samples.

We train a model on each $D_i \rightarrow$ we get model $M_1, M_2, M_3 \dots$
Finally we aggregate them.

→ Out of Bag samples: On avg. 63.2% of the original data appear in each bootstrap sample.

- That means roughly 36.8% are not used - These are called out of bag (OOB) samples.
- Used to estimate model accuracy

Why

Bootstrap?

1. Reduces Variance
2. Works well with small dataset.
3. Basis for Bagging and Boosting.

Output: Many models (ensemble)

Cross-Validation:

Estimate how well your model generalizes on unseen data by training and testing on different subsets.

Why

Cross

Validation?

1. Reliable estimate of true performance
2. Helps detect overfitting
3. For model selection and hyperparameter tuning.

2-Class Evaluation Metrics

1. Confusion Matrix

	Pred Positive	Pred Negative
Actual Positive	True Pos	False Neg
Actual Negative	False Pos	True Neg

$$FP = \text{Type I Error} \quad | \quad FN = \text{Type II error} \quad (\text{Issue})$$

2. Accuracy : Measures overall fraction of correct prediction

- Can be misleading if classes are imbalanced

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

3. Precision: Of all predicted positives , how many are correct?

$$\text{Precision} = \frac{TP}{TP + FP}$$

Medical Diagnosis
4. Recall (Sensitivity) : Of all actual positives , how many did we correctly predict

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

Fraud Detection
4. Specificity : Of actual negatives ; how many did we correctly predict as negative

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

5. F1-Score : Harmonic mean of precision and recall

- Good for imbalance datasets

$$F_1 = \frac{2 \times \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

26 ROC and AUC

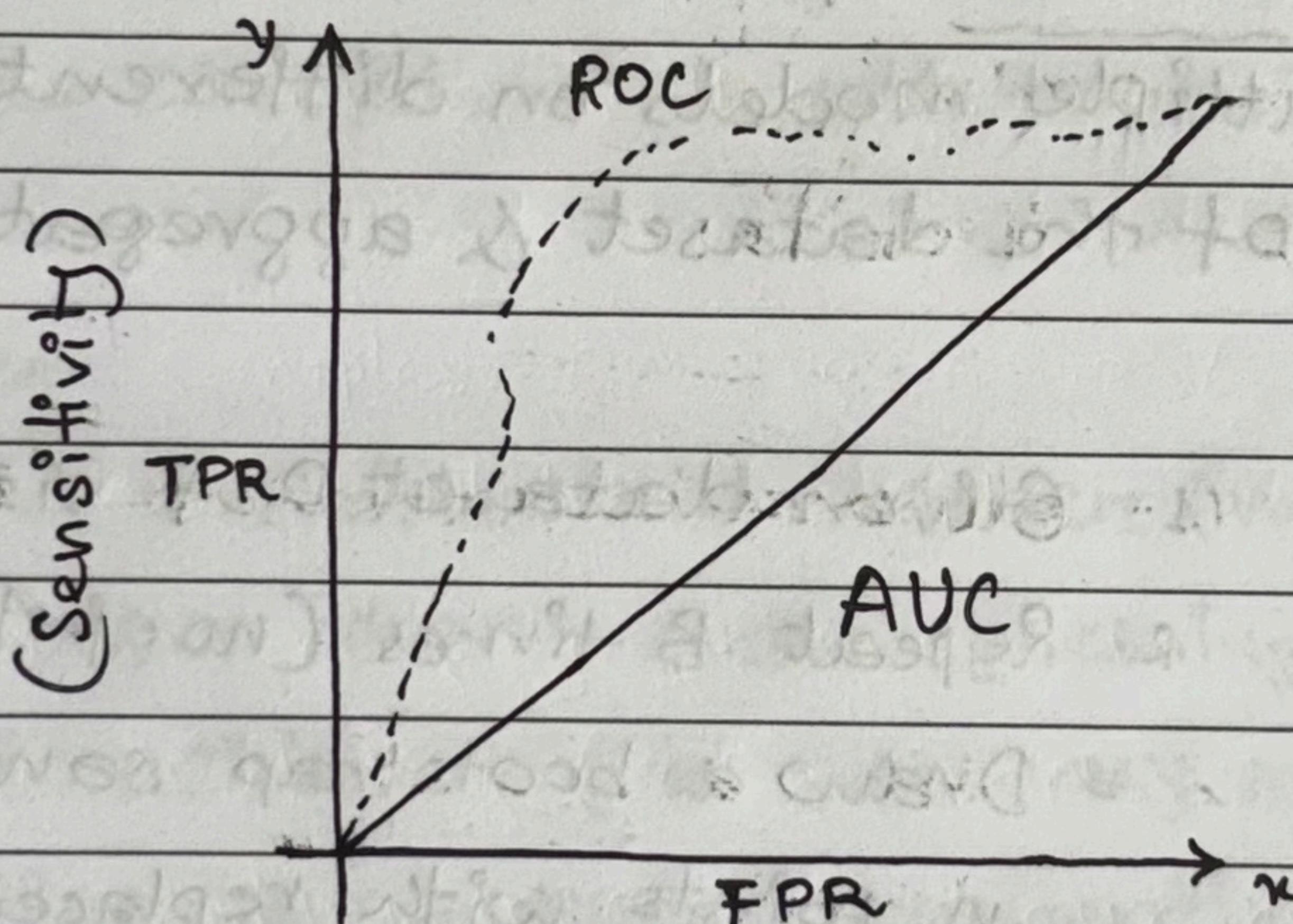
ROC Curve - Receiver Operating Characteristic Curve

- Evaluates a classifiers performance across all thresholds rather than a single threshold.

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (\text{sensitivity against FPR})$$

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

Threshold: The cutoff probability at which we classify as true



The closer the ROC curve is to top left, the better the classifier.

AUC Curve: A single-number metric that summarizes ROC Curve.

- AUC = 1 → perfect classifier

- AUC = 0.5 → random classifier

- $0.5 < \text{AUC} < 1$ → Better than random

Probability that a randomly chosen positive sample is ranked higher than a randomly chosen negative sample.

→ Precision-Recall curve works well when classes are imbalanced.

→ We require both ROC and AUC to give statement on AUC

AUC = 0.5 Random Model? → ROC curve can flip labels