# Q1-1: Consider following statements and choose the correct option (True/False for all the statements A/B/C/D).

A. *ROC curve summarize the trade-off between the true positive rate and the positive predictive value for a model*

B. *Precision-Recall curve summarize the trade-off between the true positive rate and false positive rate for a model*

C. *In both imbalanced and balanced datasets, the area under the curve (AUC) can be used as a summary of the model performance.*

D. *If we decrease the false negative (select more positives), recall always increases, but precision may increase or decrease.*

1. A: True, B: False, C: True, D: False

2. A: False, B: False, C: True, D: True

3. A: True, B: True, C: True, D: True

4. A: False, B: True, C: False, D: True

**Q1-1: Consider following statements and choose the correct option (True/False for all the statements A/B/C/D).**

A. *ROC curve summarize the trade-off between the true positive rate and the positive predictive value for a model*

B. *Precision-Recall curve summarize the trade-off between the true positive rate and false positive rate for a model*

C. *In both imbalanced and balanced datasets, the area under the curve (AUC) can be used as a summary of the model performance.*

D. *If we decrease the false negative (select more positives), recall always increases, but precision may increase or decrease.*

1. A: True, B: False, C: True, D: False

2. A: False, B: False, C: True, D: True ⬅

3. A: True, B: True, C: True, D: True

4. A: False, B: True, C: False, D: True

# Q1-2: Which of the following metrics is NOT specifically tailored to help with evaluating highly imbalanced data?

1. Precision and Recall

2. Area Under the ROC curve

3. Accuracy

4. All of the above helps in evaluating highly imbalanced data

# Q1-2: Which of the following metrics is NOT specifically tailored to help with evaluating highly imbalanced data?

1. Precision and Recall

2. Area Under the ROC curve

3. Accuracy ⬅

4. All of the above helps in evaluating highly imbalanced data

If you have an imbalanced dataset **accuracy** can give you false assumptions regarding the classifier's performance, it's better to rely on precision and recall.

**Q2-1:** A learned model $h$ makes 10 errors over the 100 instances. Calculate the 95% confidence interval i.e. With approximately 95% probability, the true error lies in the interval _____ . Take $z_C = 2$

1. $0.1 \pm 0.02$

2. $0.1 \pm 0.04$

3. $0.1 \pm 0.06$

4. $0.1 \pm 0.08$

$$error_S(h) \pm z_C \sqrt{\frac{error_S(h)(1-error_S(h))}{n}}$$

**Q2-1:** A learned model $h$ makes 10 errors over the 100 instances. Calculate the 95% confidence interval i.e. With approximately 95% probability, the true error lies in the interval _____ . Take $z_C = 2$

1. $0.1 \pm 0.02$
2. $0.1 \pm 0.04$
3. $0.1 \pm 0.06$ ⬅
4. $0.1 \pm 0.08$

$$error_S(h) \pm z_C \sqrt{\frac{error_S(h)(1 - error_S(h))}{n}}$$

Here, r = 10, n = 100.
error(h) = r/n = 0.1
Substitute all the values in the formula
$0.1 \pm 2 \times$ sqrt(0.1 x 0.9 / 100)
$= 0.1 \pm 2 \times 0.3/10 = \mathbf{0.1 \pm 0.06}$

# Q2-2: Which of the following statements is FALSE?

1. The null hypothesis states that the 2 learning systems have the same accuracy
2. Alternative hypothesis states that one of the systems is more accurate than the other
3. If p is sufficiently small, then reject the alternative hypothesis
4. A two tailed test asks if the accuracy of the two systems are different.

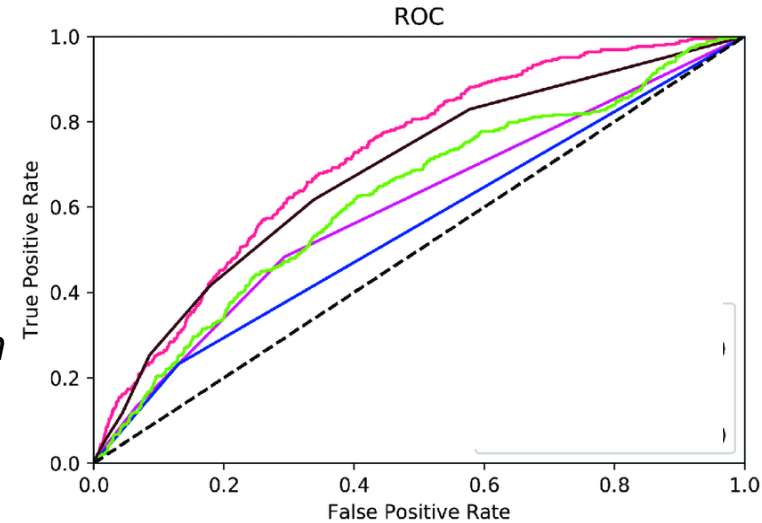# Q2-2: Which of the following statements is FALSE?

1. The null hypothesis states that the 2 learning systems have the same accuracy
2. Alternative hypothesis states that one of the systems is more accurate than the other
3. If p is sufficiently small, then reject the alternative hypothesis ⬅
4. A two tailed test asks if the accuracy of the two systems are different.

If p is sufficiently small, then reject the **null** hypothesis

# Q3-1: The figure shows ROC curve for different models. Select the correct option.
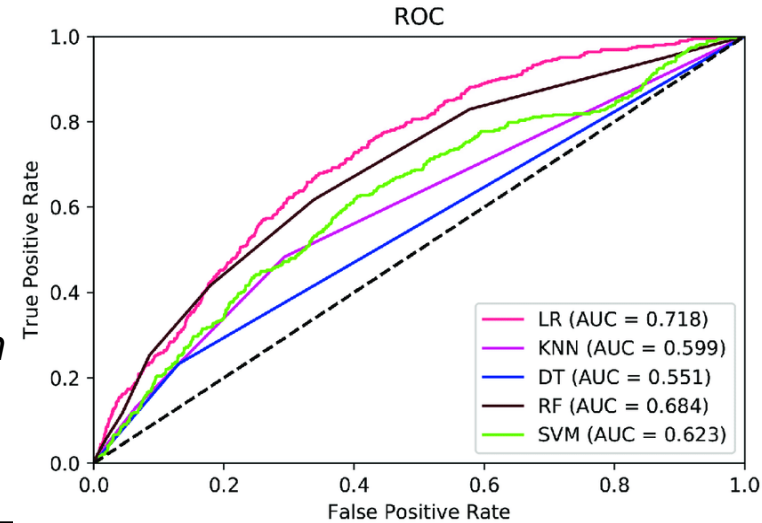
A. *Dashed black line represents random classification.*

B. *ROC curve for any model can't fall below the dashed black line.*

C. *The model represented by solid blue line is better than that represent by solid lime.*



1. Statement A is true. Statement B, C are false.

2. Statement A, B are true. Statement C is false.

3. Statement B, C are true. Statement A is false.

4. All Statements are true.

# Q3-1: The figure shows ROC curve for different models. Select the correct option.

A. *Dashed black line represents random classification.*

B. *ROC curve for any model can't fall below the dashed black line.*

C. *The model represented by solid blue line is better than that represent by solid lime.*



1. Statement A is true. Statement B, C are false.

2. Statement A, B are true. Statement C is false.

3. Statement B, C are true. Statement A is false.

4. All Statements are true.

# Q3-2: How to avoid pitfalls while training a model?

1. Collect test data that is true representation of real world.

2. Don't access the label of a test instance while training.

3. Avoid excessive preprocessing/training on a particular dataset.

4. All of the above.

# Q3-2: How to avoid pitfalls while training a model?

1. Collect test data that is true representation of real world.

2. Don't access the label of a test instance while training.

3. Avoid excessive preprocessing/training on a particular dataset.

4. All of the above.