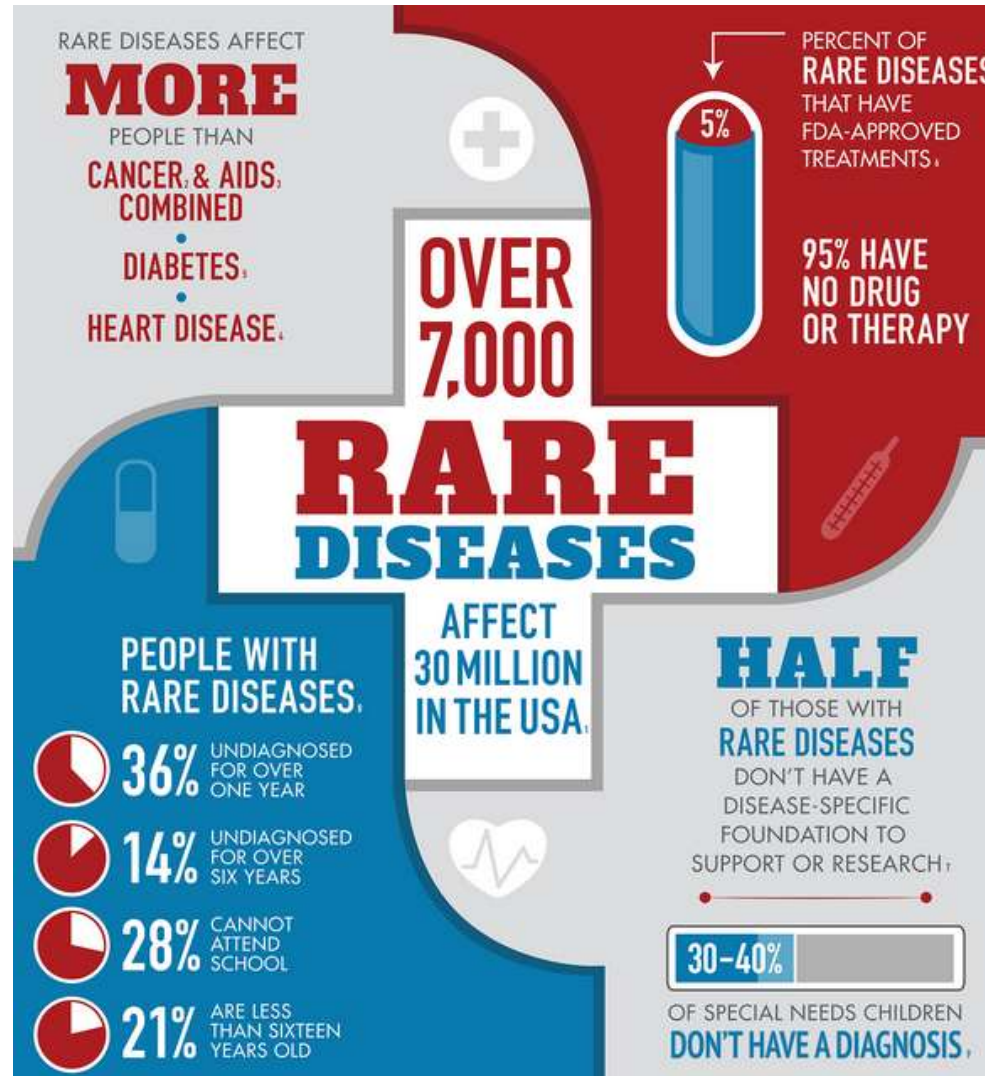# Precision-Recall (PR) vs Receiver Operating Characteristics (ROC) curves.

# Which one to use if data is imbalanced?

*dr. ir. Sultan K. Imangaliyev*
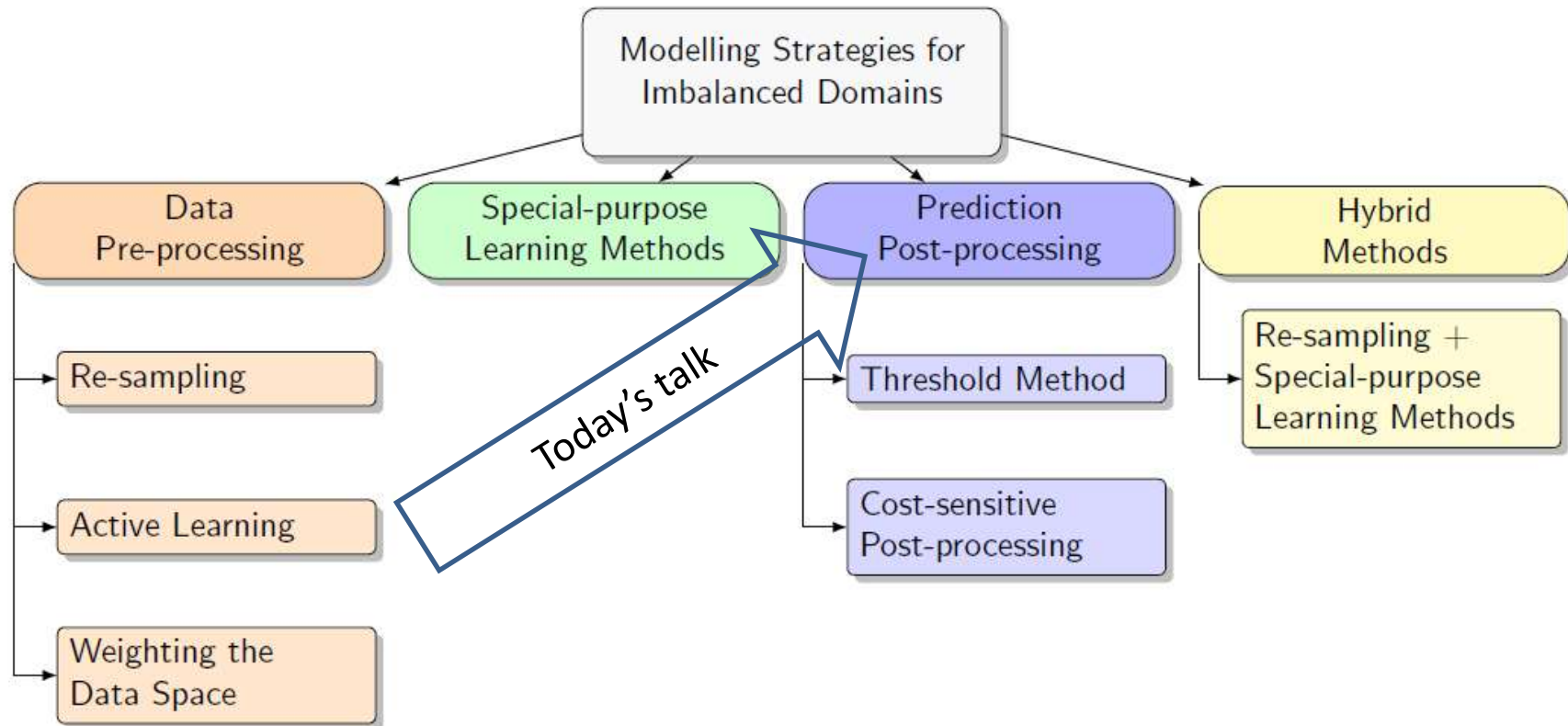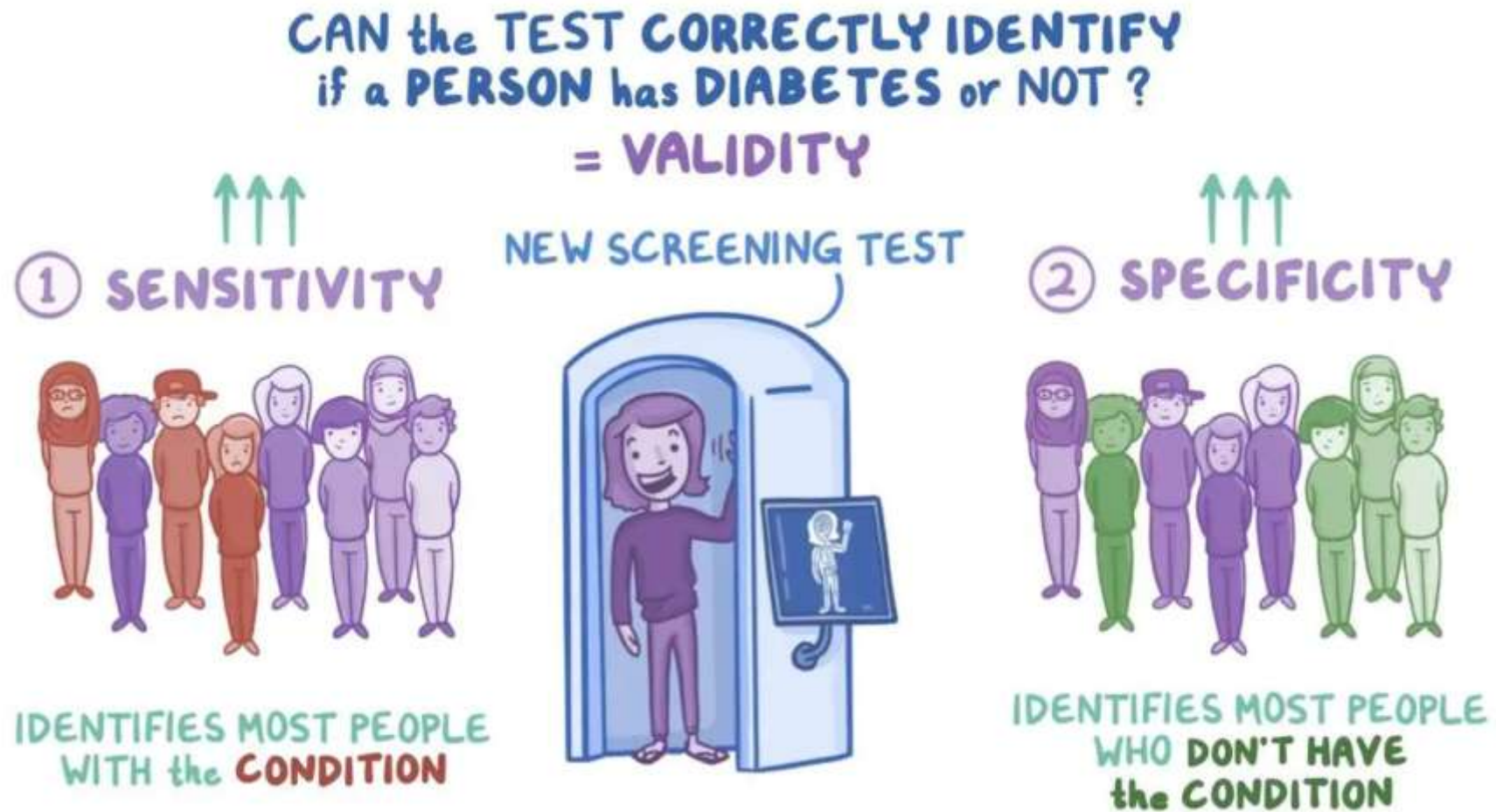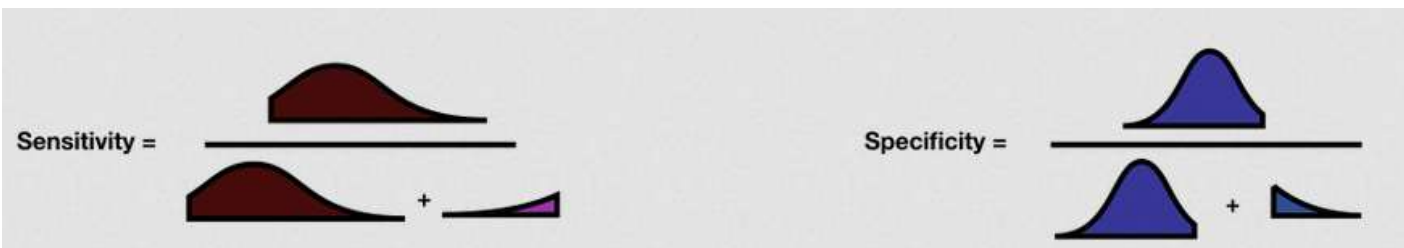
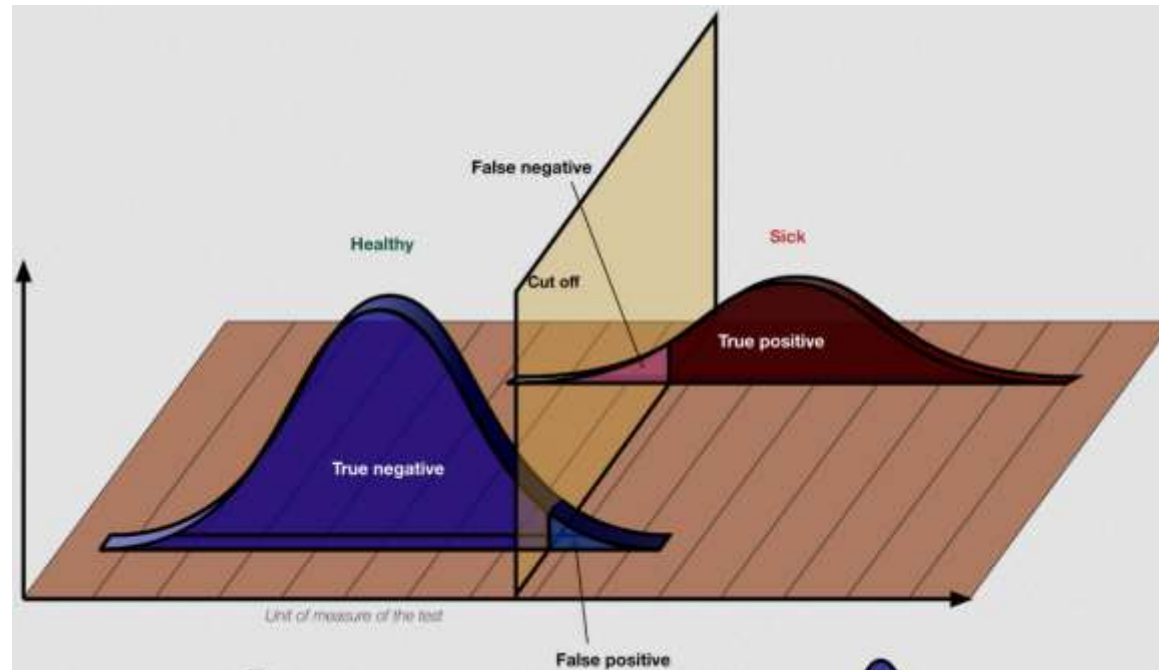*Leiden, 22/01/26*

# Imbalanced data (1)

# Imbalanced data (2)



Figure 7: Main modelling strategies for imbalanced domains.

Branco, Paula, Luís Torgo, and Rita P. Ribeiro. "A survey of predictive modeling on imbalanced domains." *ACM Computing Surveys (CSUR)* 49.2 (2016): 1-50.

# Sensitivity vs Specificity (1)



CAN the TEST CORRECTLY IDENTIFY
if a PERSON has DIABETES or NOT ?
= VALIDITY

NEW SCREENING TEST

① SENSITIVITY

② SPECIFICITY

IDENTIFIES MOST PEOPLE
WITH the CONDITION

IDENTIFIES MOST PEOPLE
WHO DON'T HAVE
the CONDITION

# Sensitivity vs Specificity (2)

# Precision vs Recall (1)



**RECALL**
How many relevant items are returned.

9 returned
18 relevant
9/18 = 50%

Relevant or Responsive Docs

Non-Responsive or Irrelevant Docs

**PRECISION**
How many selected items are relevant.

9 relevant
12 returned
9/12 = 60%

# Precision vs Recall (2)

# Precision vs Recall (3)



Low recall, low precision



High recall, low precision



High recall, high precision
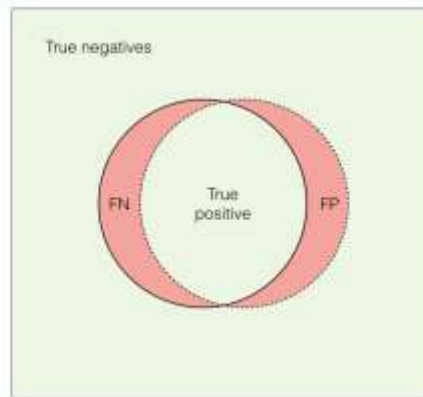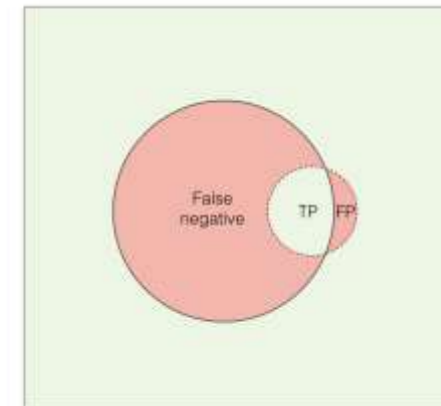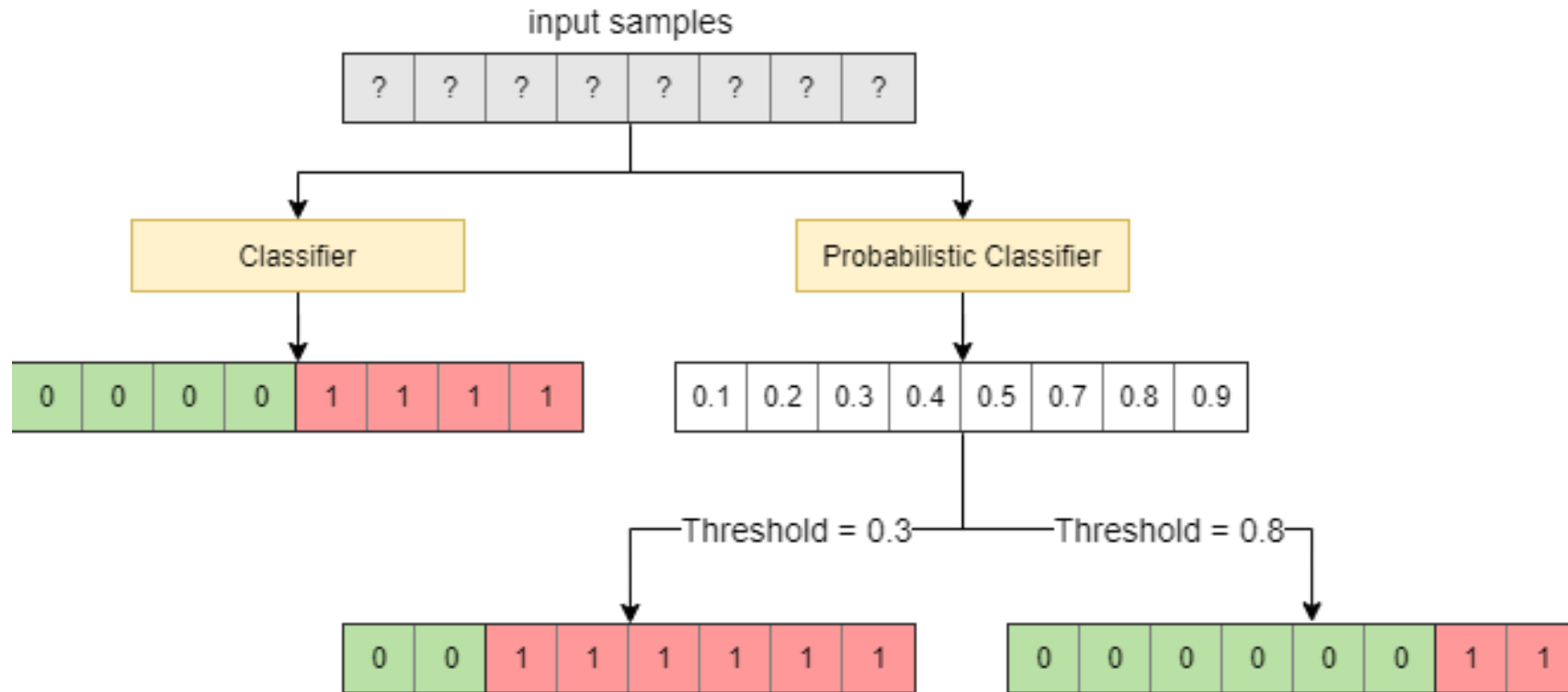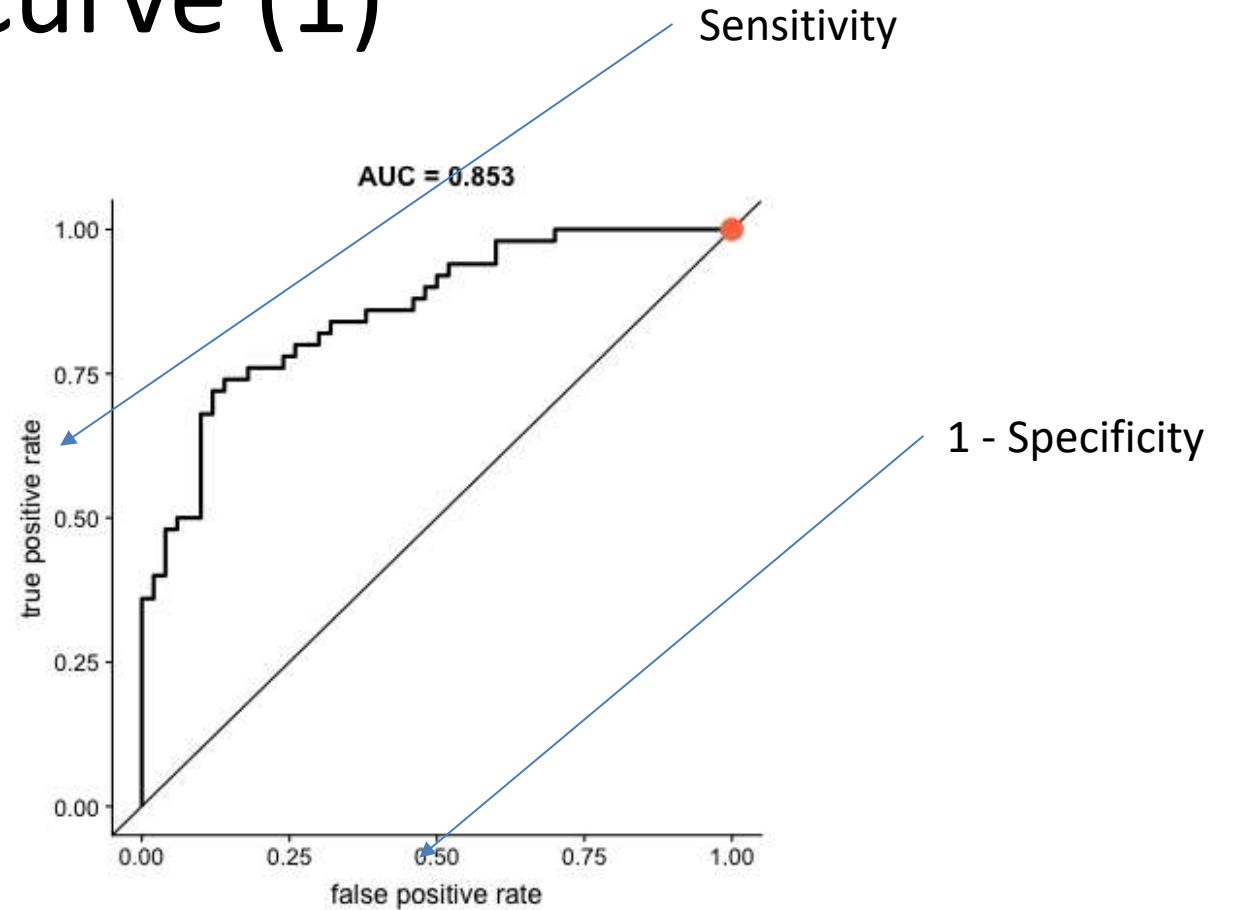


Low recall, high precision

# Probabilistic classifiers

# ROC curve (1)



*AUCROC can be interpreted as the probability that the scores given by a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one.*

# ROC curve (2)



AUC = 0.496

# PR curve

*The PR curves plots the following parameters:*

*Precision = TP/(TP+FP)*

*Recall = TP/(TP+FN)*

*Notice how True Negatives (TN) are absent from the equation?*

*PR curves are useful when positive examples are rare.*

# Comparing curves, balanced simulated data

# PR curve, balanced data

# Misleading AUC values

# Comparing curves, imbalanced simulated data
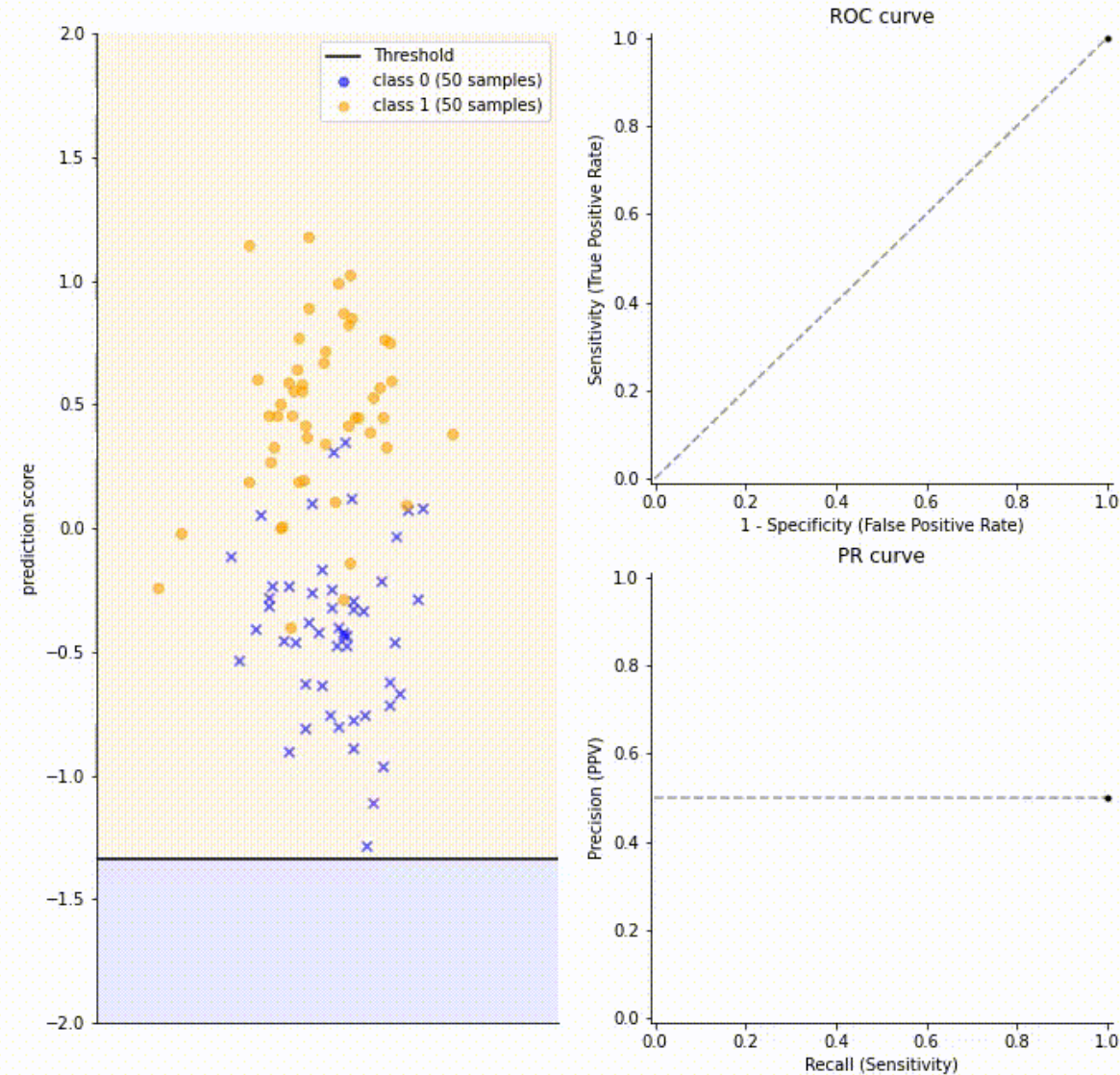


ROC curve with its optimal threshold (green) and the PR optimal threshold (grey)

optimal threshold for ROC (purple) and PR (green).

PR curve with its optimal threshold (purple) and the ROC optimal threshold (grey)

*https://towardsdatascience.com/on-roc-and-precision-recall-curves-c23e9b63820c*

# PR curve, imbalanced data
## (*positive examples are rare*)



*... precision and recall make it possible to assess the performance of a classifier on the minority class*

# PR curve, imbalanced data
# (*negative examples are rare*)

# Practical tip (simplified)

- When to use PR AUC or AUROC?

  - When two classes are equally important
    - AUROC if the goal is to perform equally well on both classes.
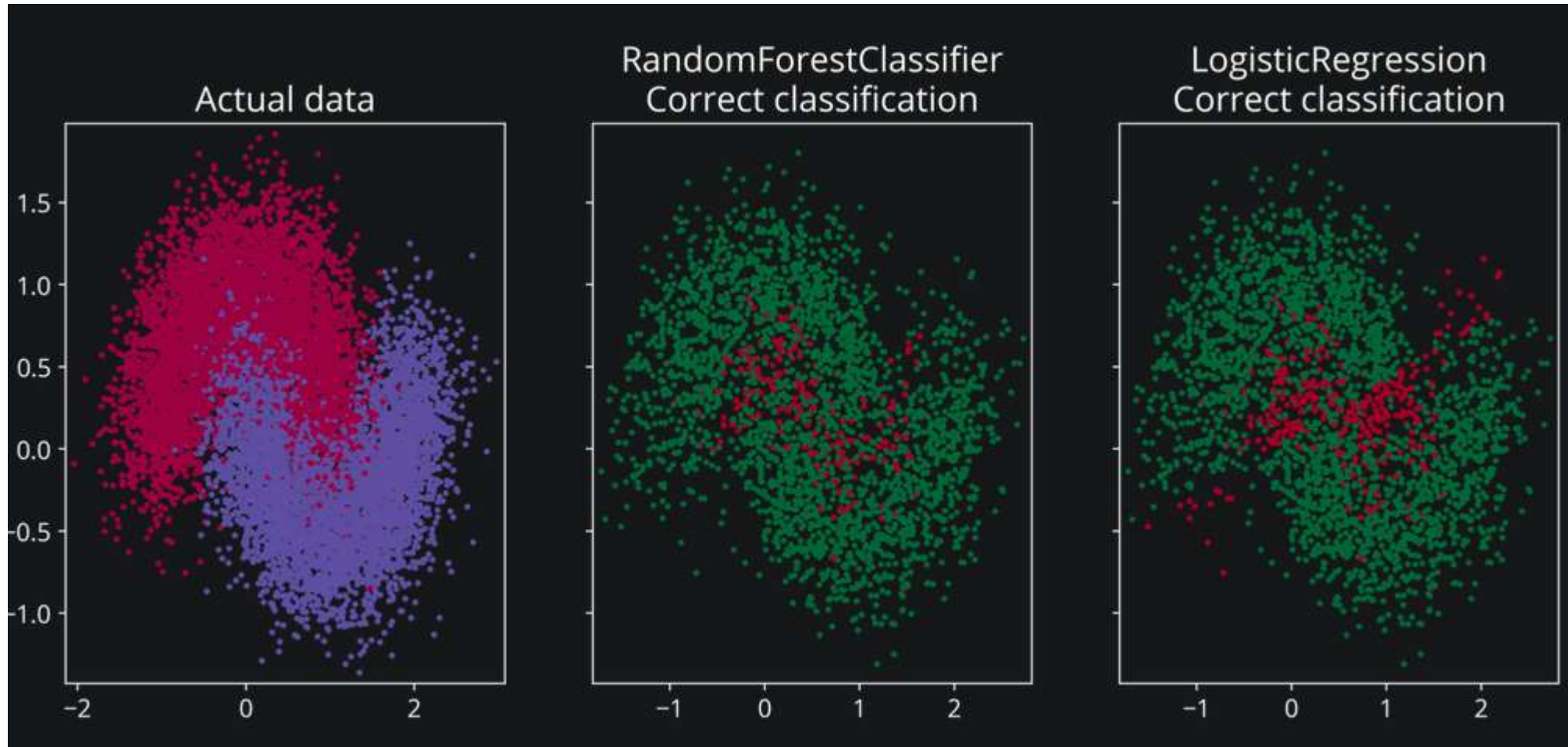      - Image classification between cats & dogs; the performance on cats is equally important on dogs.

  - When minority class is more important
    - PR AUC if the focus of the model is to identify correctly as many positive samples as possible.
      - Spam detectors; regular emails are not of interest at all — they overshadow the number of positives.
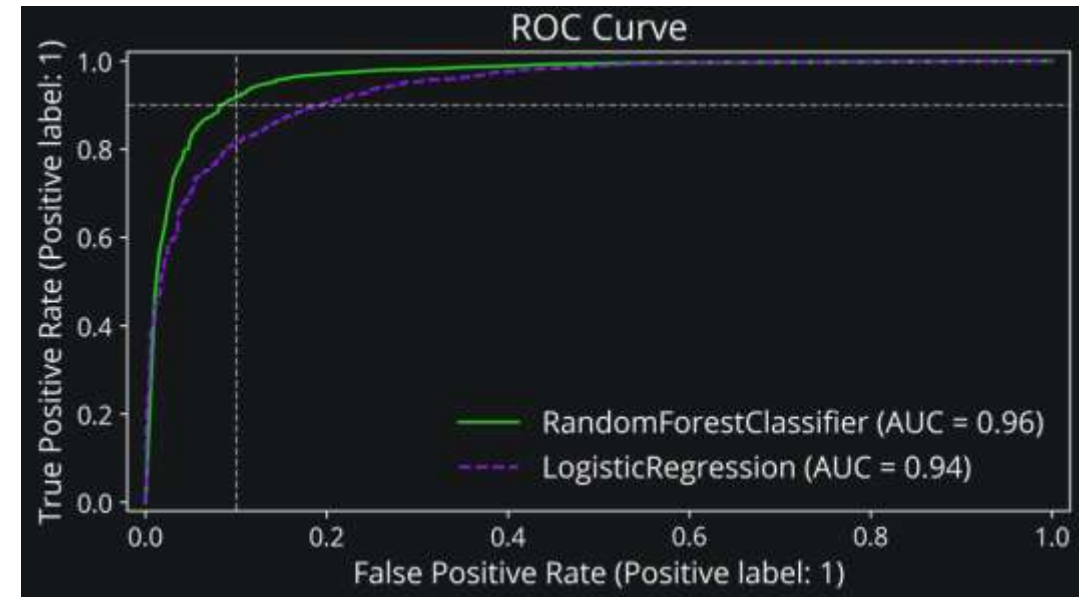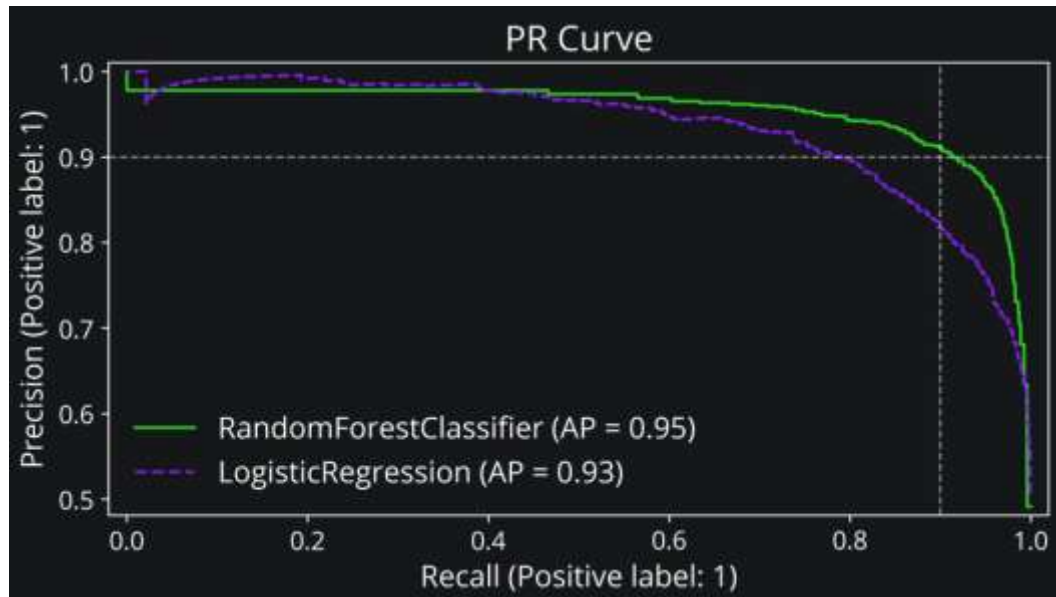
# Words of caution

- ROC curves have common advantages
  - Universal baselines
    - A random classifier is represented by a major diagonal
  - Linear interpolation
    - Any two points on an ROC curve can be linearly interpolated
  - Interpretable area:
    - AUROC can be used to calculate an expected accuracy of the model:

- PR curves have pitfalls
  - No universal baseline
    - Performance of a random classifier depends on the prevalence in the data set
  - No linear interpolation
  - Uninterpretable area
    - Area under curve, commonly calculated by a trapezodial rule that performs linear interpolation, might be an overly optimistic measure
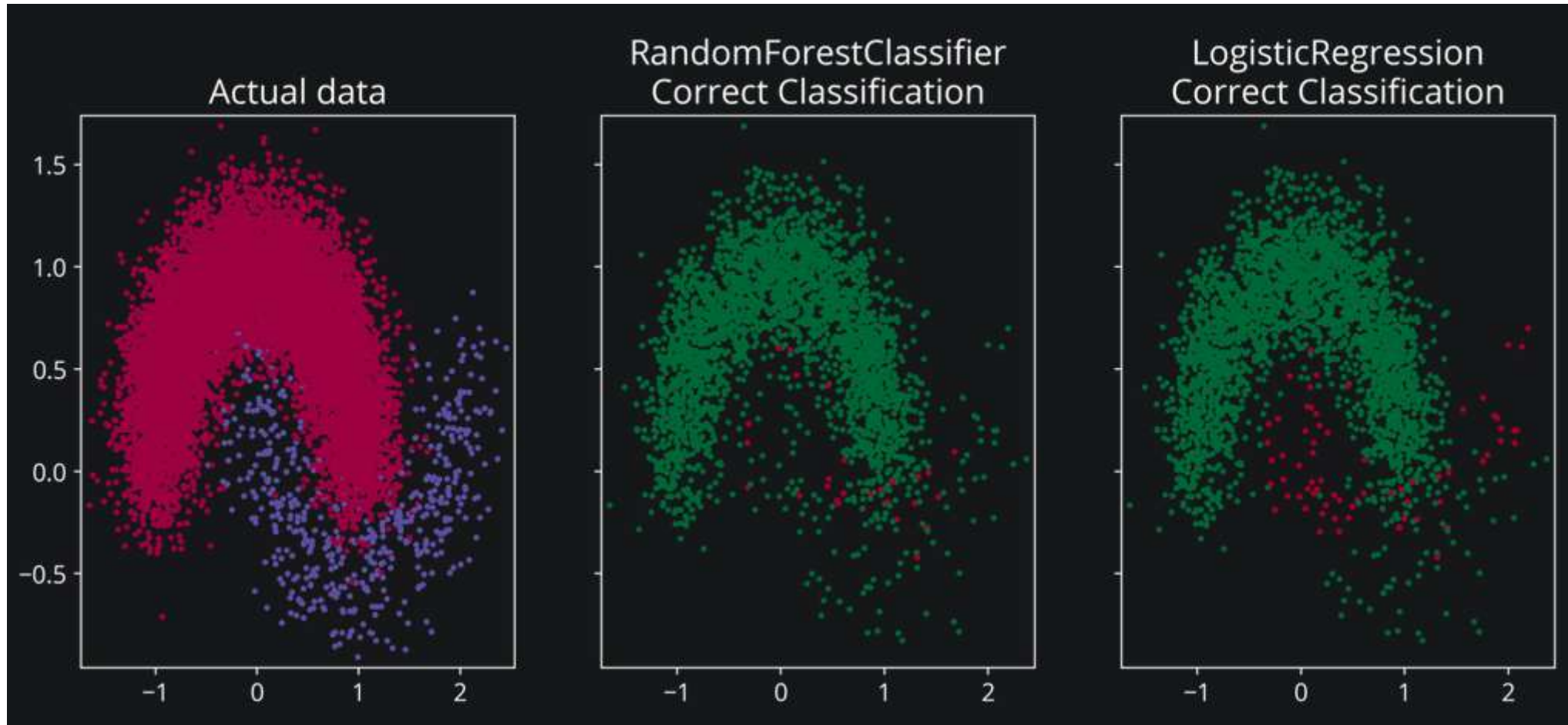
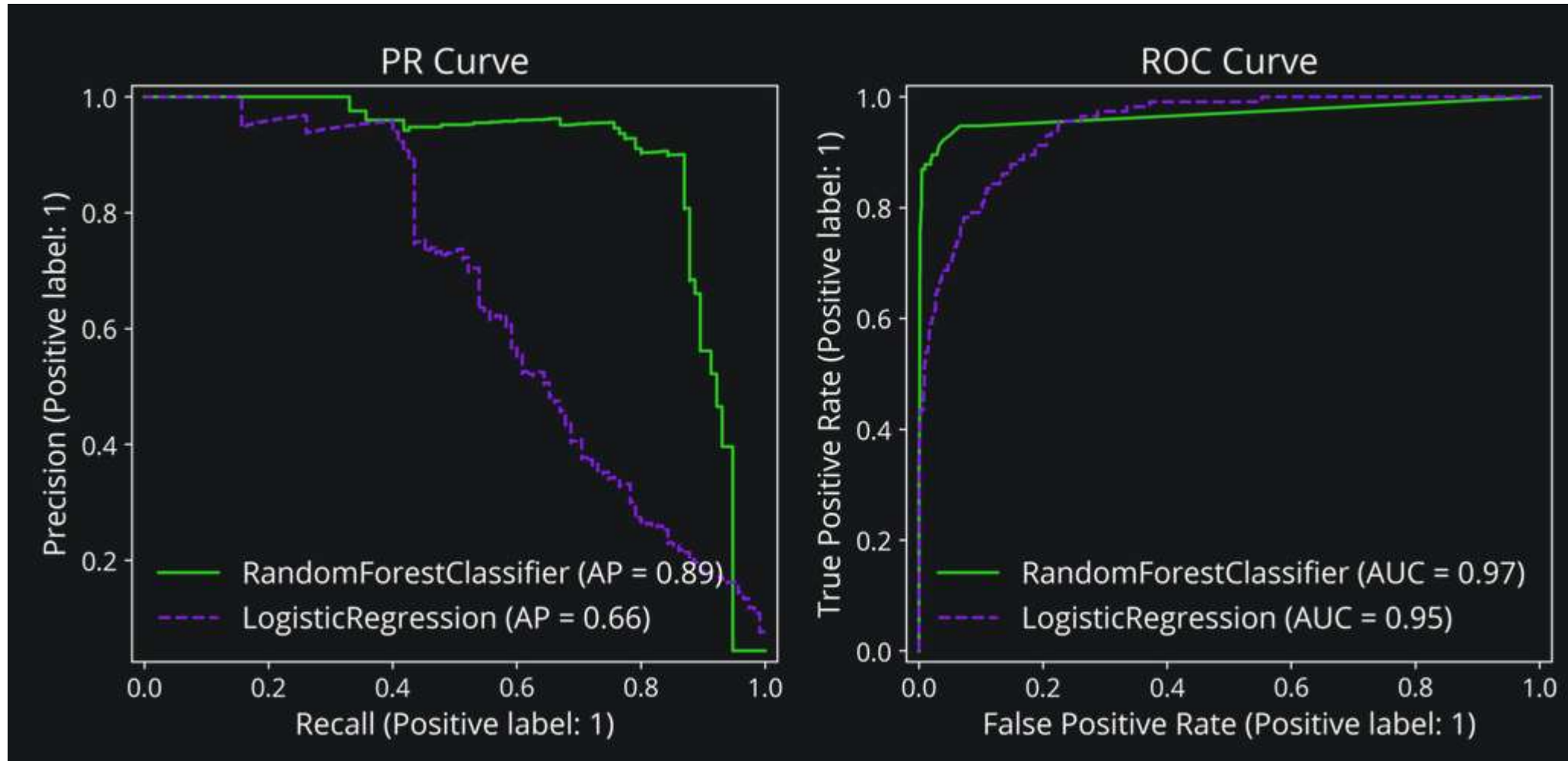# Comparing classifiers, balanced simulated data (1)

# Comparing classifiers, balanced simulated data (2)

# Comparing classifiers, imbalanced simulated data (1)

23

# Comparing classifiers, imbalanced simulated data (2)

# Conclusions

- There is no universal way of treating imbalanced data
  - Choices are application-specific

- In case of balanced data AUROC and AP are comparable

- For imbalanced data
  - Both classes are important – use AUROC
  - Minority positive class is more important – use PR curve

- Extremely high imbalance?
  - Consider switching to anomaly detection methods

# Q&A