

Machine learning e modelli di classificazione di dati biomedici

Christian Salvatore
Scuola Universitaria Superiore IUSS Pavia

VALIDATION & PERFORMANCE EVALUATION

How to test a machine-learning classifier?

A good validation process allows to obtain a **minimally biased estimate** of the true diagnostic performance of the classifier



- Correct quantification of the discriminatory power of a given model (model evaluation)
- Possibility to compare classification techniques based on different approaches (model selection)

How to test a machine-learning classifier?

For example, if parameter selection, training of the predictive model and validation are performed using the same dataset, the generated classifier will show limited generalization ability when classifying unseen samples



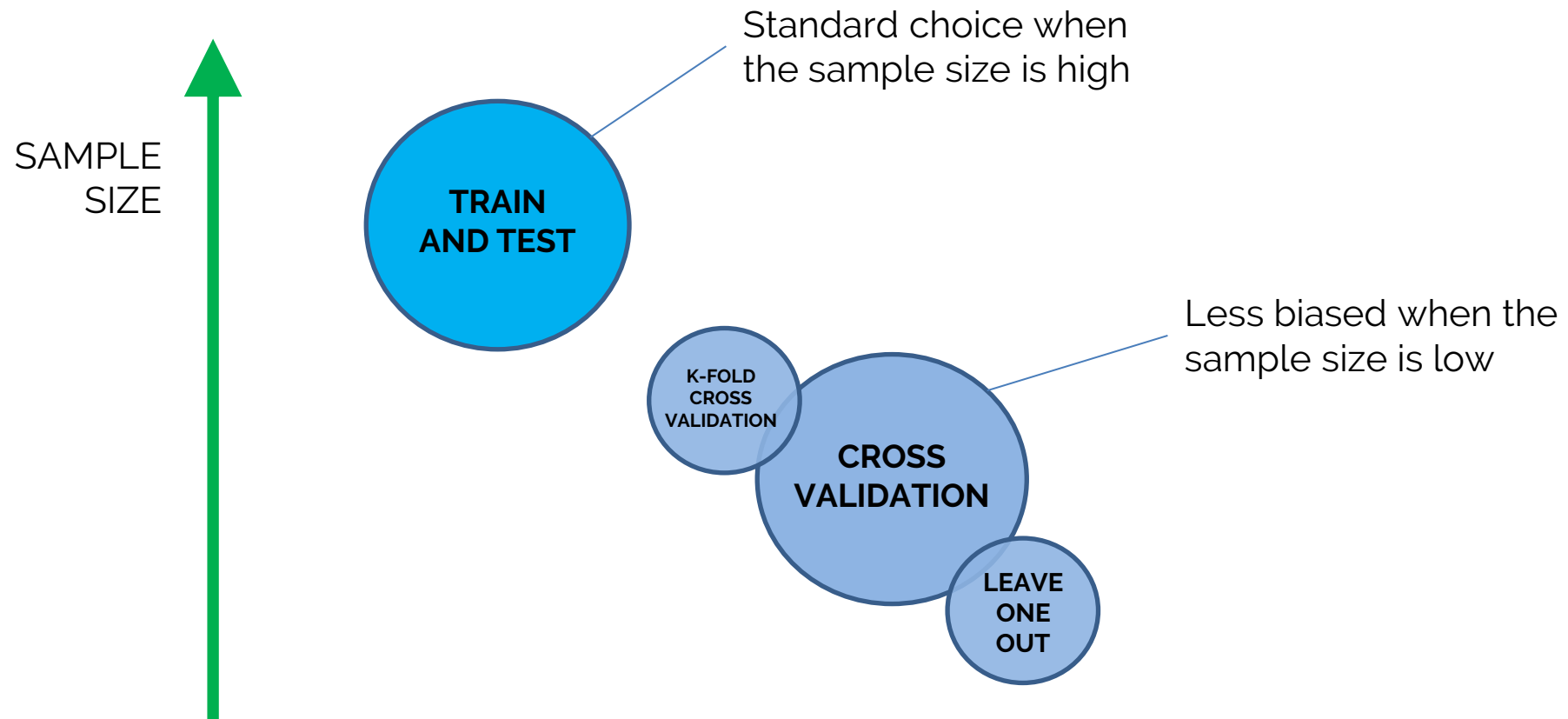
LOW
TRAINING ERROR



HIGH
TESTING ERROR
(low generalization
ability)

OVERFITTING

Which validation approach?

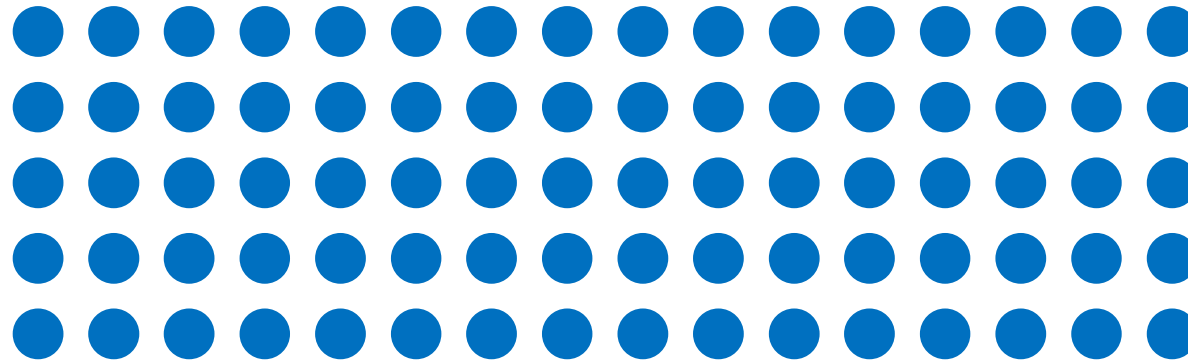


Train-and-test

This kind of procedure is used when the number of samples in the original dataset is high enough to allow its splitting into two subsets including different samples, which can be used to train and test the classifier.

Train-and-test

This kind of procedure is used when the number of samples in the original dataset is high enough to allow its splitting into two subsets including different samples, which can be used to train and test the classifier.

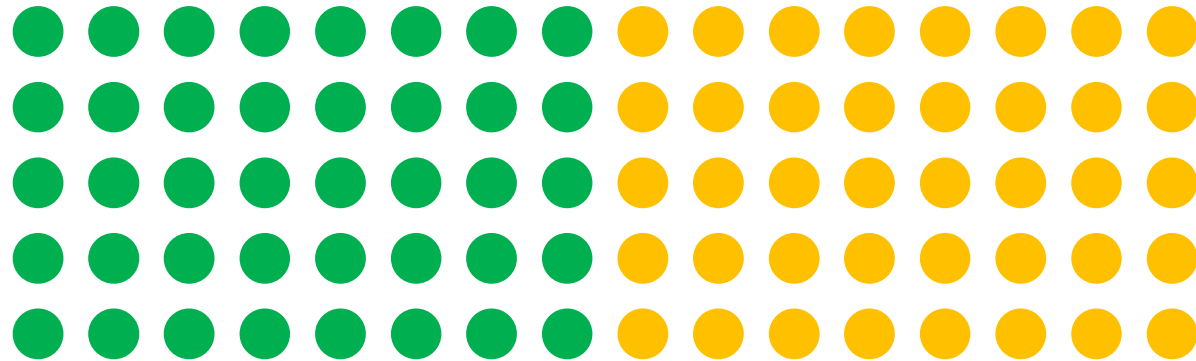


The original dataset is partitioned into 2 complementary subsets, the TRAINING set and the TESTING set.

The TRAINING set is used to train the classifier

The TESTING set is used for validation

Train-and-test



● Training

● Testing

Advantages:

- Over-training problems are reduced, because the training and testing sets are completely independent

Drawbacks:

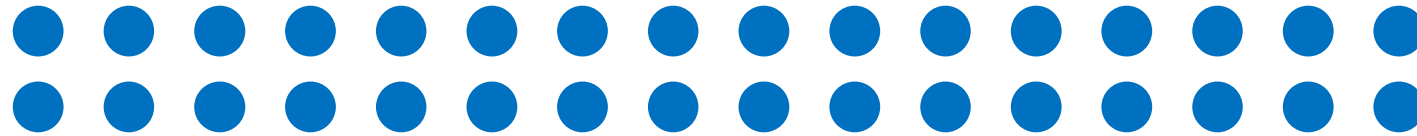
- Results could be related to the particular choice of the partition subsets

Leave-one-out cross validation

Leave-One-Out (LOO) CV can be considered a particular form of k-fold CV in which

k = number of samples in the original dataset

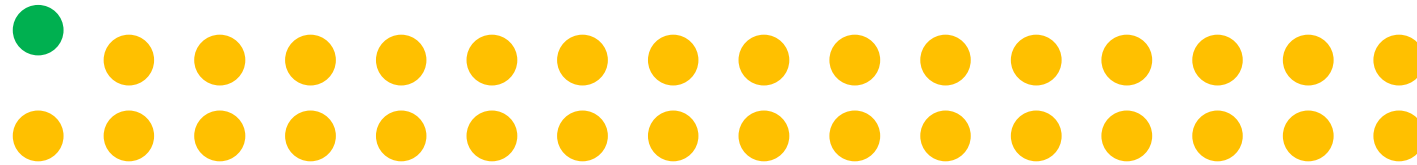
Leave-one-out cross validation



TRAINING of the classifier is performed using $n-1$ samples of the original dataset

TESTING is performed using the remaining sample (n being the total number of samples in the original dataset)

Leave-one-out cross validation

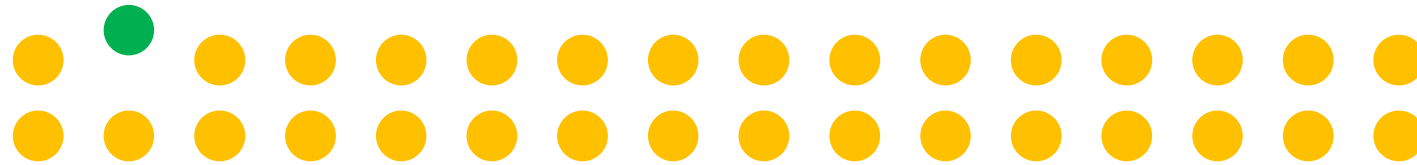


● Training

● Testing

The procedure is then repeated n times, until all samples are used once for validation.

Leave-one-out cross validation

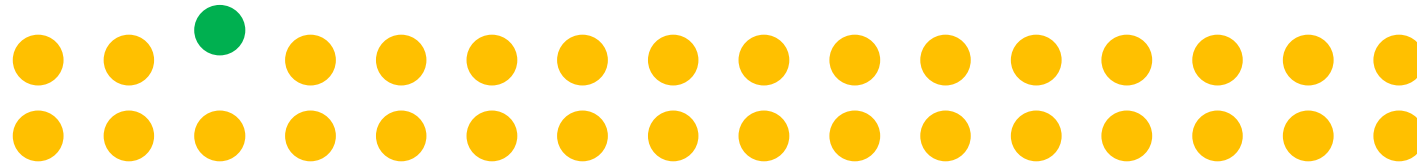


● Training

● Testing

The procedure is then repeated n times, until all samples are used once for validation.

Leave-one-out cross validation

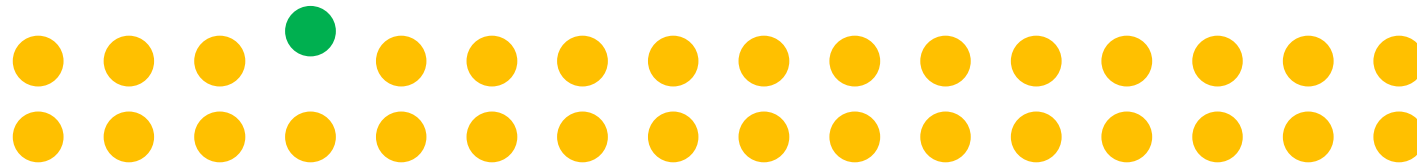


● Training

● Testing

The procedure is then repeated n times, until all samples are used once for validation.

Leave-one-out cross validation



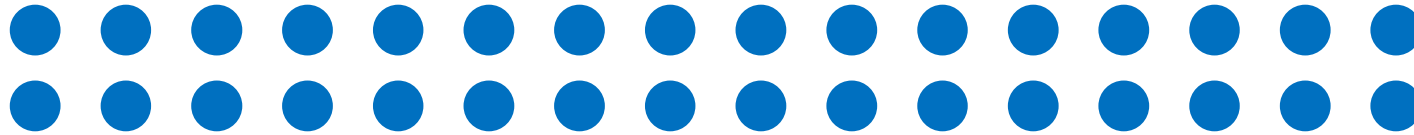
● Training

● Testing

and so on...

Cross validation

Quantification of the discriminatory power of a predictive model even if the size of the dataset is small



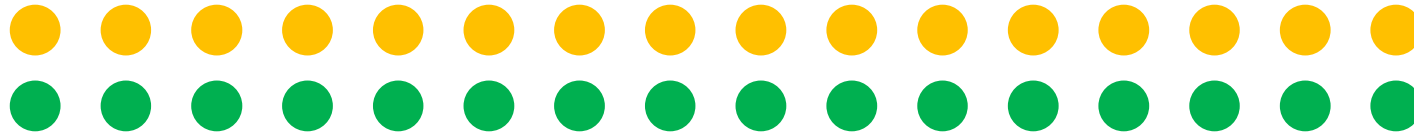
CV involves partitioning the original dataset into complementary subsets, the training set and the testing set

The TRAINING set is used to train the classifier

The TESTING set is used to validate the generated predictive model

Cross validation

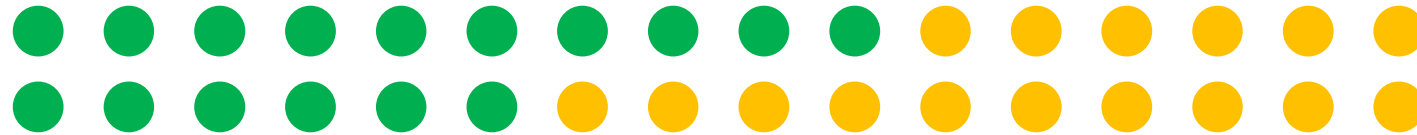
● Training
● Testing



By using different partitions of the original dataset, multiple rounds of CV can be performed, which can aid reducing the variability of the partitioned subsets.

Cross validation

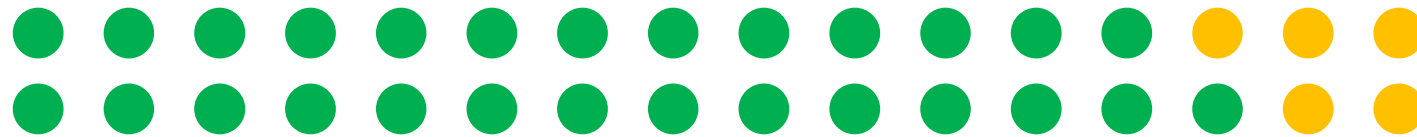
● Training
● Testing



By using different partitions of the original dataset, multiple rounds of CV can be performed, which can aid reducing the variability of the partitioned subsets.

Cross validation

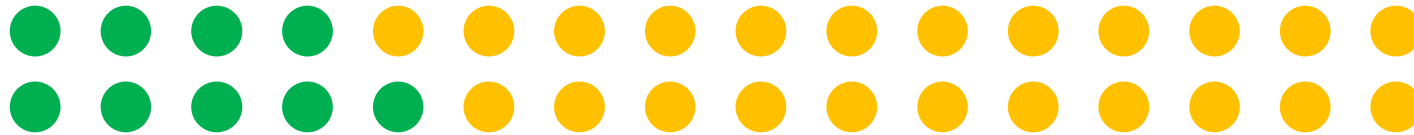
● Training
● Testing



By using different partitions of the original dataset, multiple rounds of CV can be performed, which can aid reducing the variability of the partitioned subsets.

Cross validation

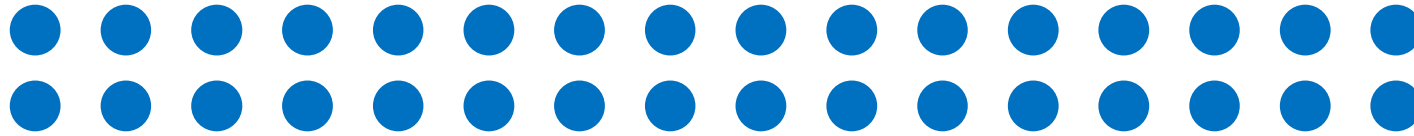
● Training
● Testing



and so on...

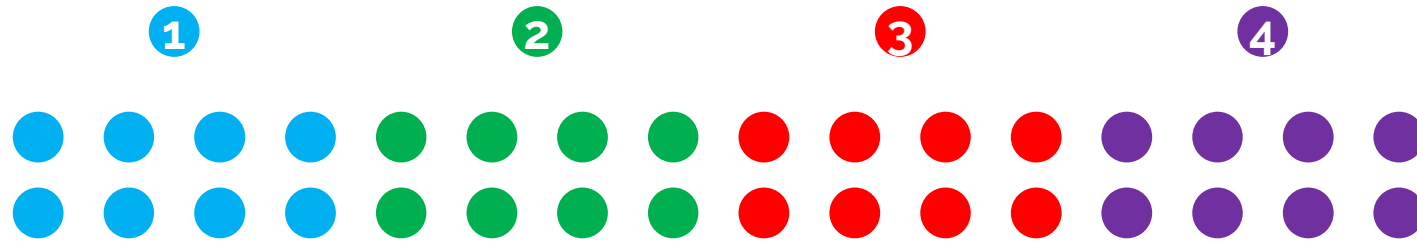
Results obtained from multiple rounds can be averaged in order to obtain a quantification of the performance of the classifier.

K-fold cross validation



The original dataset is randomly partitioned into k subsets of equal size.

K-fold cross validation

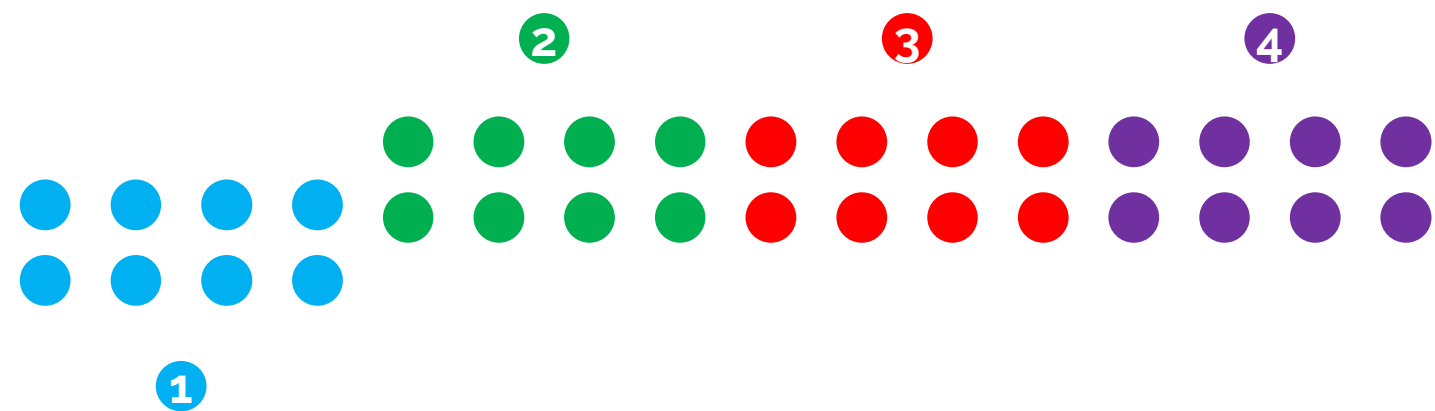


The original dataset is randomly partitioned into k subsets of equal size

TRAINING of the classifier is performed using $k-1$ subsets

TESTING is performed using the remaining subset

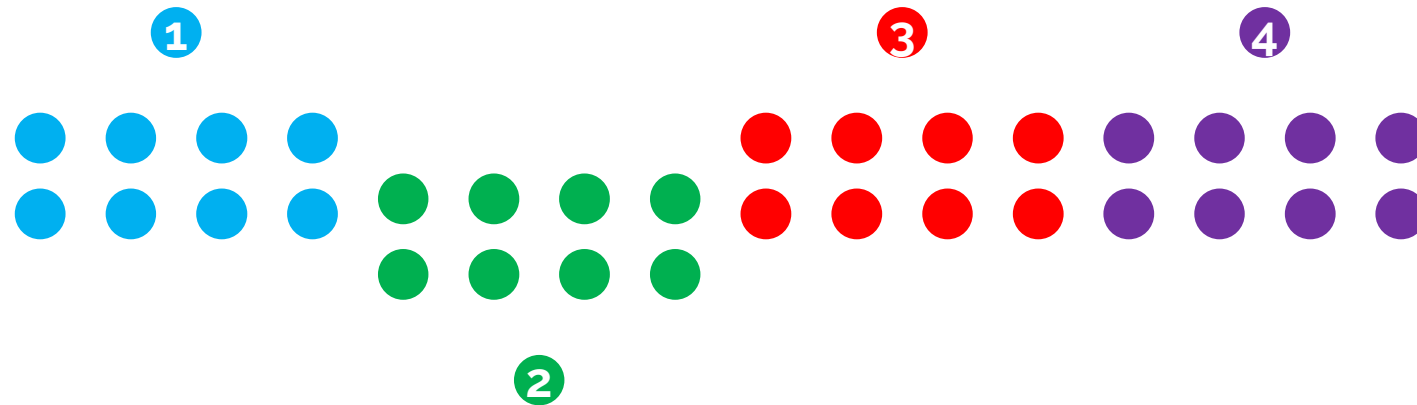
K-fold cross validation



2 3 4 Training
1 Testing

The procedure is then repeated k times, until all subsets are used once as testing set.

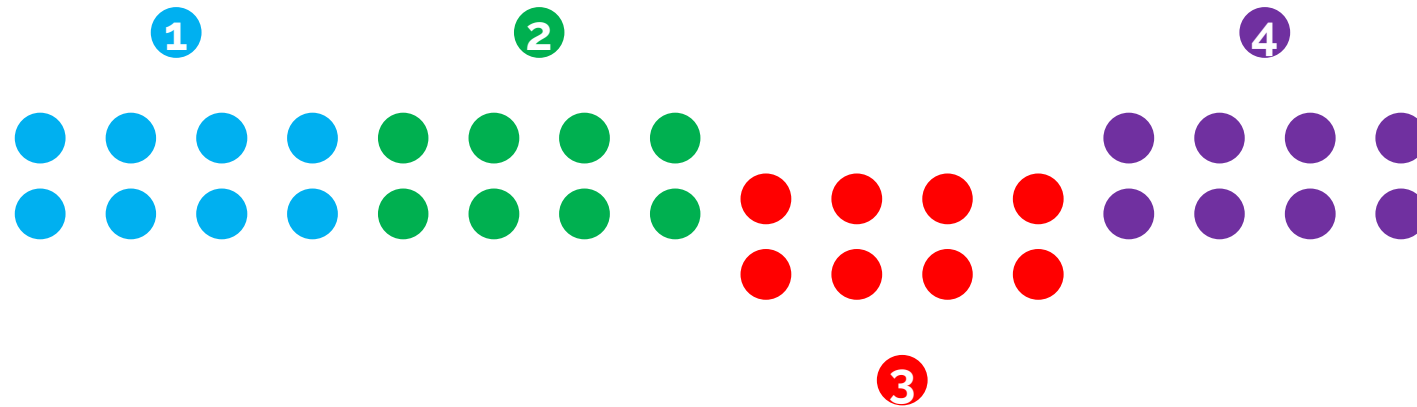
K-fold cross validation



① ③ ④ Training
② Testing

The procedure is then repeated k times, until all subsets are used once as testing set.

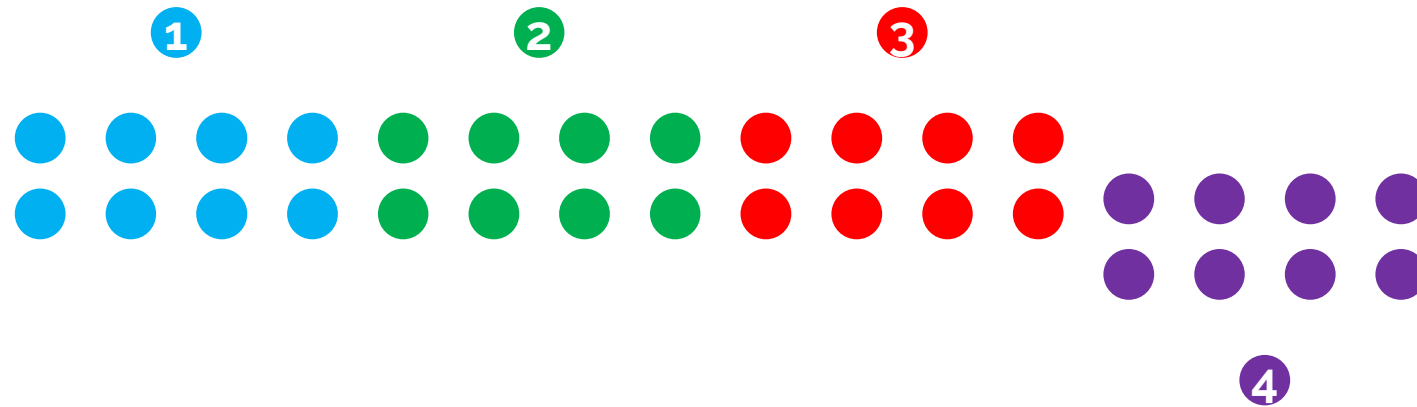
K-fold cross validation



1 **2** **4** Training
3 Testing

The procedure is then repeated k times, until all subsets are used once as testing set.

K-fold cross validation



1 2 3

Training

4

Testing

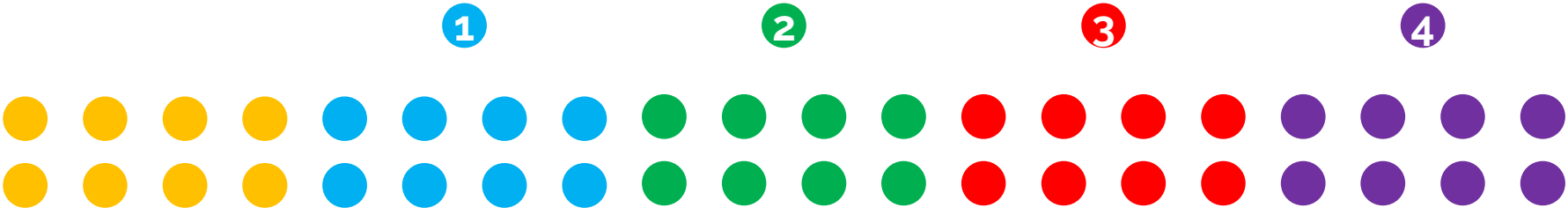
The procedure is then repeated k times, until all subsets are used once as testing set.

K-fold cross validation

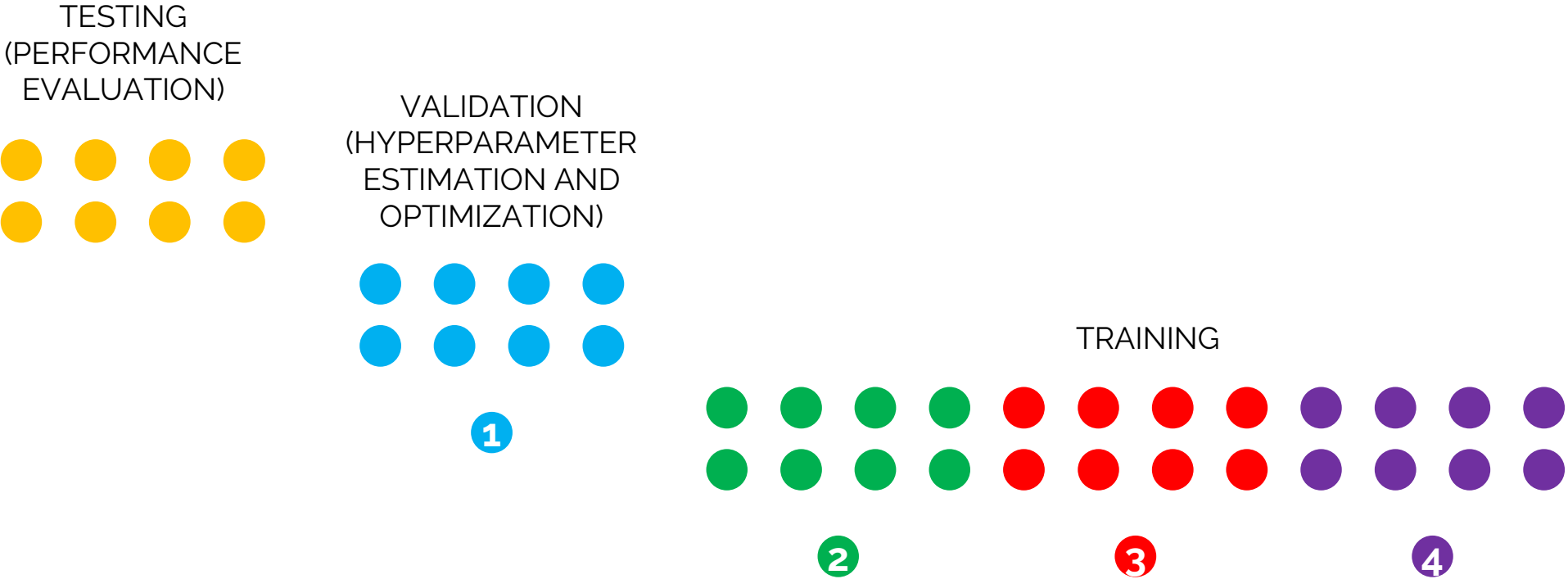
Advantages:

- each sample of the original dataset is used once for validation
- all samples being used for both training and testing phases

Nested K-fold cross validation



Nested K-fold cross validation



Nested K-fold cross validation

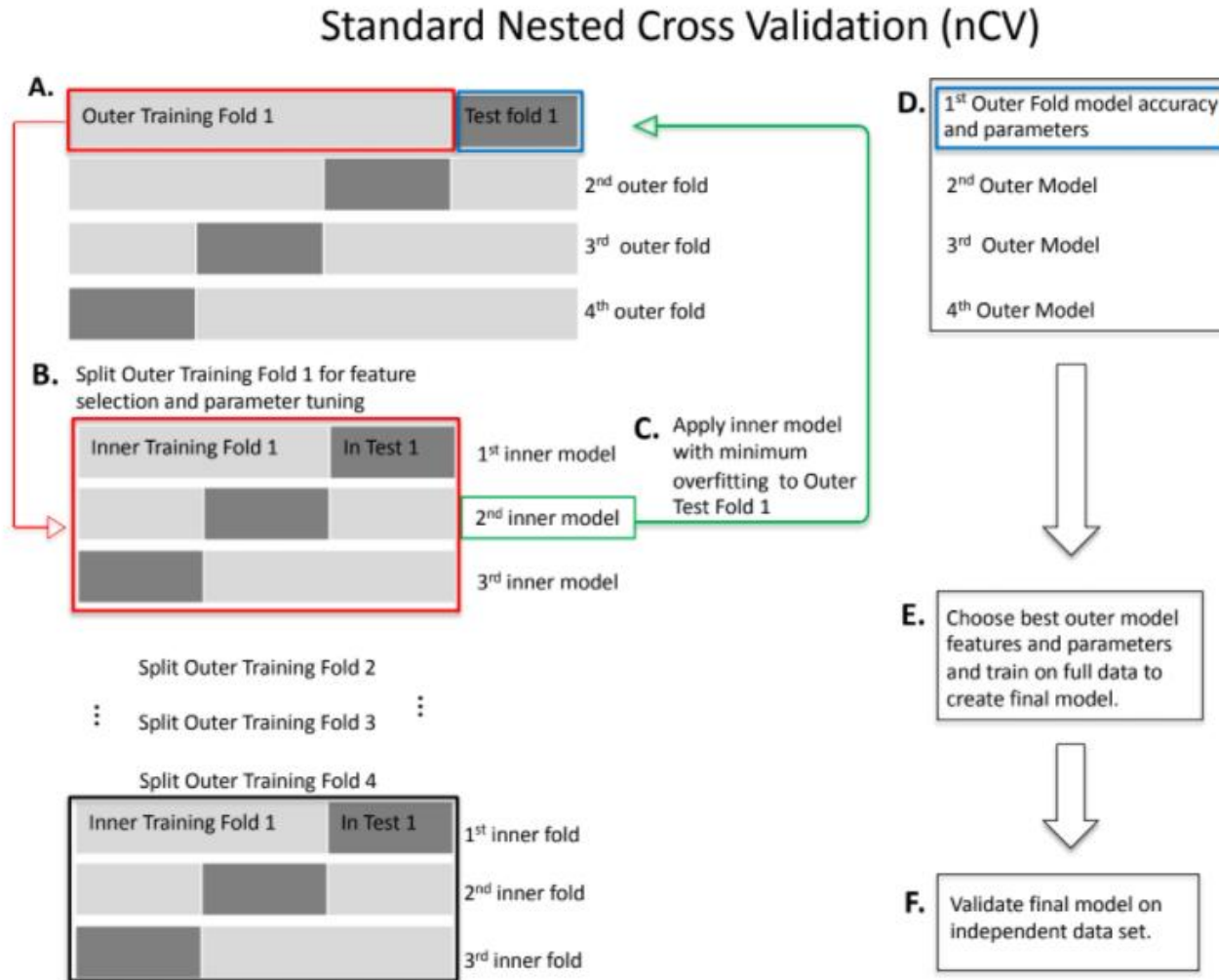


Fig. 1. Standard nested Cross-Validation (nCV). **A.** Split the data into outer folds of training and testing data pairs (4 outer folds in this illustration). Then do the following for each outer training fold (illustration starting with Outer Training Fold 1 (red box, A)). **B.** Split outer training fold into inner folds for feature selection and possible hyperparameter tuning by grid search. **C.** Use the best inner training model including features and parameters (2nd inner model, green box, for illustration) based on minimum overfitting (difference between training and test accuracies) in the inner folds to test on the outer test fold (green arrow to blue box, Test Fold 1). **D.** Save the best model for this outer fold including the features and test accuracies. Repeat B-D for the remaining outer folds. **E.** Choose the best outer model with its features based on minimum overfitting. Train on the full data to create the final model. **F.** Validate the final model on independent data.

Nested K-fold cross validation

Consensus Nested Cross Validation (cnCV)

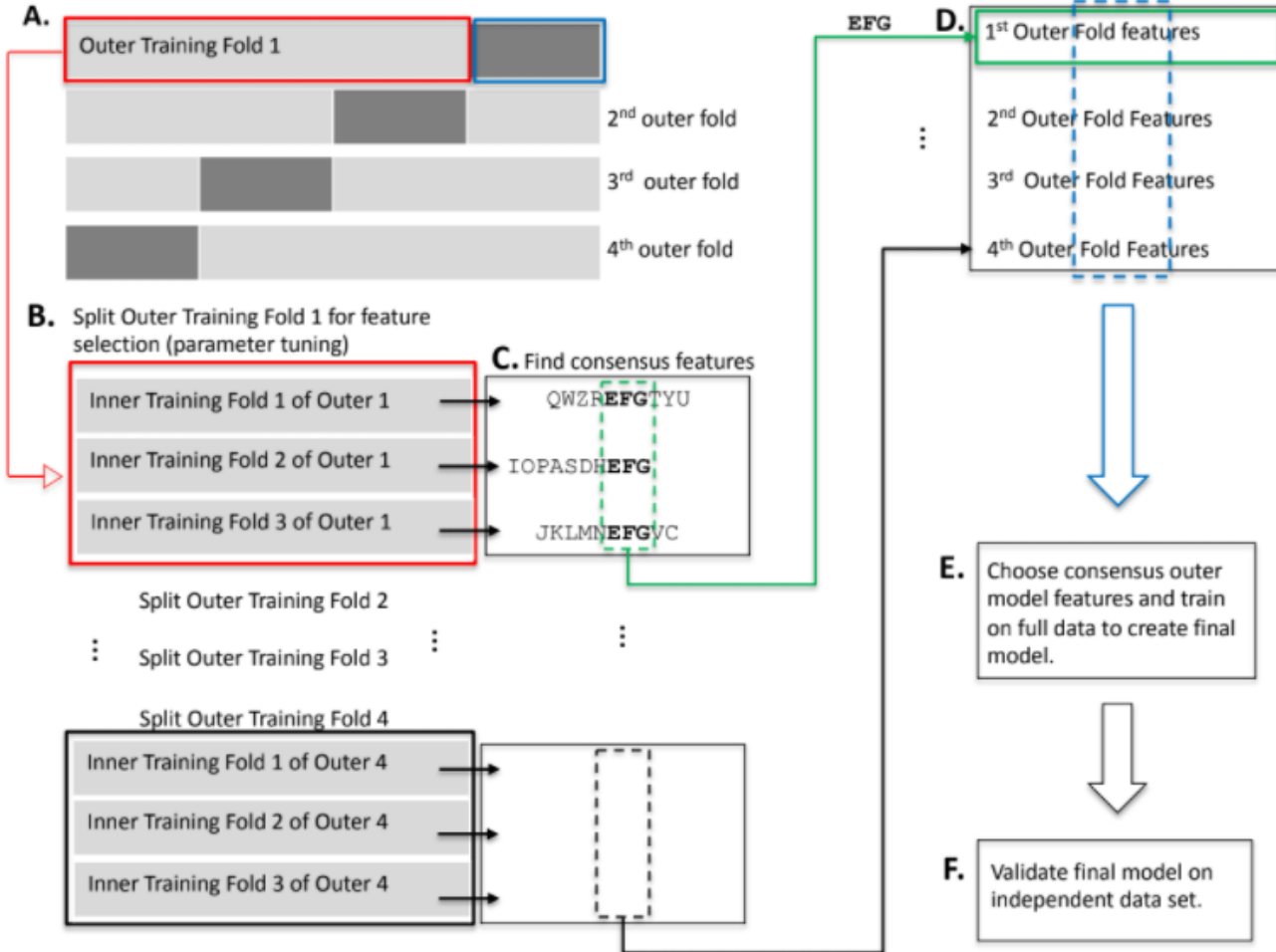


Fig. 2. Consensus Nested Cross-Validation (cnCV). **A.** Split the data into outer folds (4 outer folds in this illustration). Then do the following for each outer training fold (illustration starting with Outer Training Fold 1 (red box, A)). **B.** Split outer training fold into inner folds for feature selection and optional hyperparameter tuning by grid search. **C.** Find consensus features. For each fold, features with positive Relief scores are collected (e.g., “QWZREFGTYU” for fold 1). Negative Relief scores have high probability of being irrelevant to classification. The implementation allows for different feature importance methods and tuning the number of input features. Consensus features (found in all folds) are used as the best features in the corresponding outer fold. For example, features “EFG” are shared across the three inner folds. This procedure is used in the inner and outer folds of cnCV. Classification is not needed to select consensus features. **D.** The best outer fold features (green arrow to green box) are found for each fold (i.e., Repeat B-D for all outer folds). **E.** Choose the consensus features across all the outer folds to train the final model on full data. Consensus features are selected based on training data only. Classification is not performed until the outer consensus features are selected (A-D). **F.** Validate the final model on independent data.

Evaluation metrics

A brief overview. . .
(see next lessons)

1. Accuracy
2. Sensitivity and specificity
3. AUC

Accuracy

Accuracy is the most used metric in classification problems

Accuracy of classification =

correctly classified
samples (for both
classes)



classified samples

If the error rate is defined as the number of misclassified samples (both classes) divided by the total number of classified samples, it is evident that accuracy and error rate are complementary measures.

Sensitivity and specificity

Two metrics of great importance in medicine are sensitivity and specificity, as they measure the rate of correctly classified samples in the positive (pathological) and negative (normal) class, respectively.

Sensitivity (also known as True Positive Rate or Recall) is given by the number of correctly classified samples belonging to the positive class (true positives) divided by the total number of samples belonging to the positive class (true positives plus false negatives).

Sensitivity =

correctly classified
samples in the positive
class



positive samples

Sensitivity and specificity

Specificity (also known as True Negative Rate) is given by the number of correctly classified samples belonging to the negative class (true negatives) divided by the total number of samples belonging to the negative class (true negatives plus false positives).

Specificity =

correctly classified
samples in the negative
class

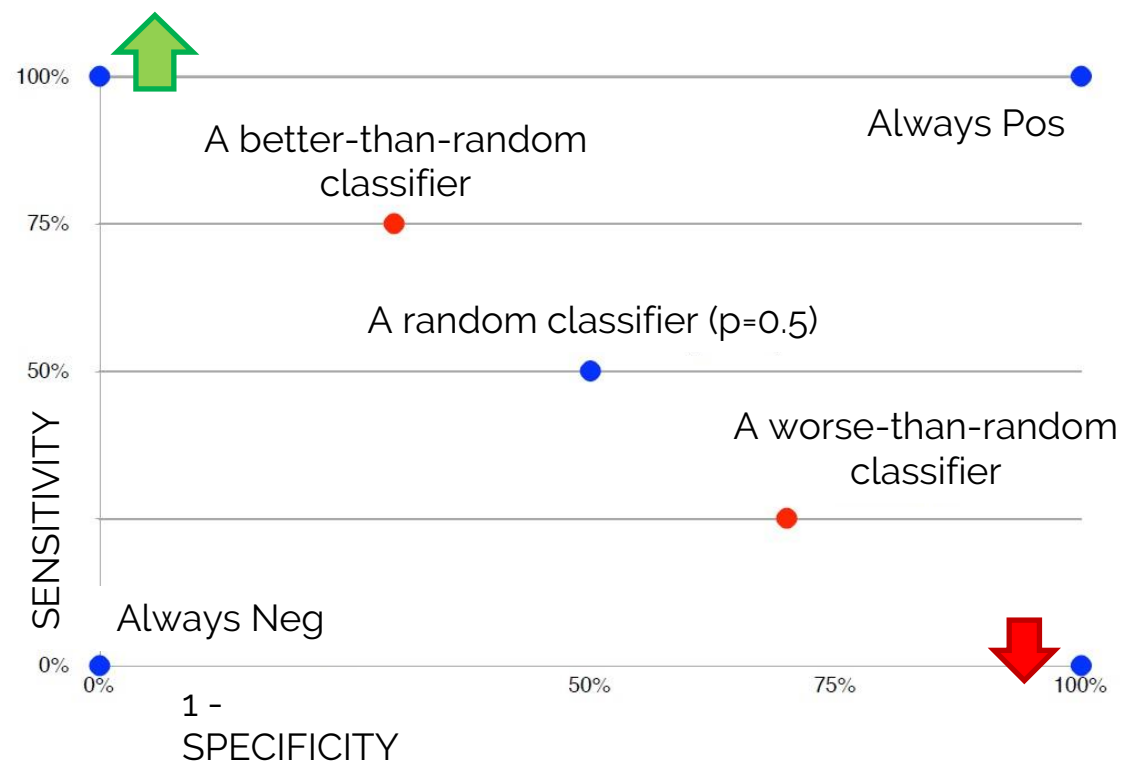
negative samples

Here, true positive (negative) gives the number of correctly classified samples belonging to the positive (negative) class, while false positive (negative) gives the number of misclassified samples belonging to the negative (positive) class.

ROC analysis and Area Under the (ROC) Curve

Another important metric in classification problems is given by the study of the Receiver Operating Characteristic (ROC) curve.

For a binary classifier, A ROC curve is a plot of the TPR (sensitivity) against the FPR (1 – specificity), which can be obtained at different setting thresholds.

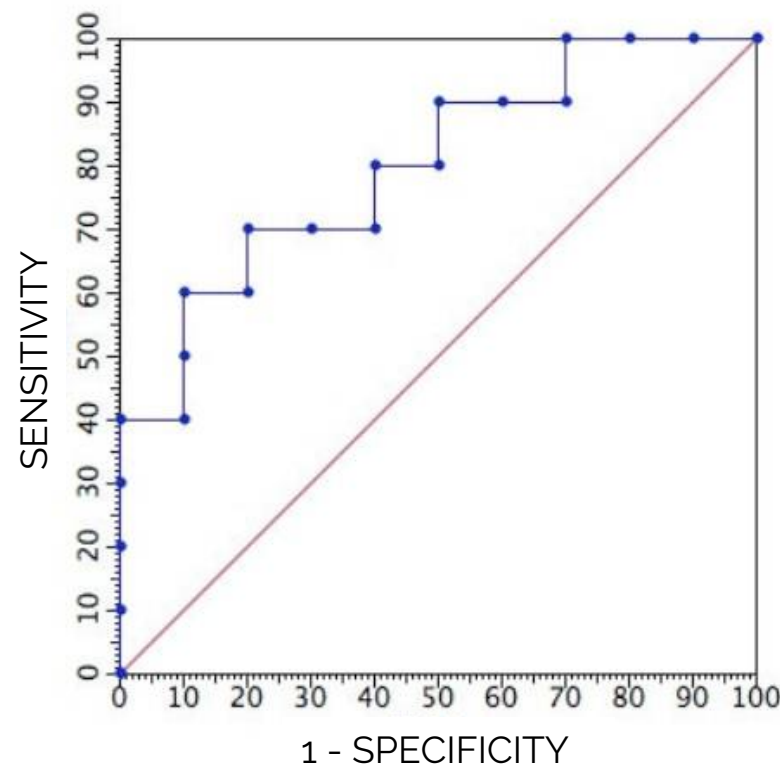


The Area Under the ROC Curve (AUC) gives a quantification of the classifier performance, with a higher statistical consistency than accuracy.

ROC analysis and Area Under the (ROC) Curve

Another important metric in classification problems is given by the study of the Receiver Operating Characteristic (ROC) curve.

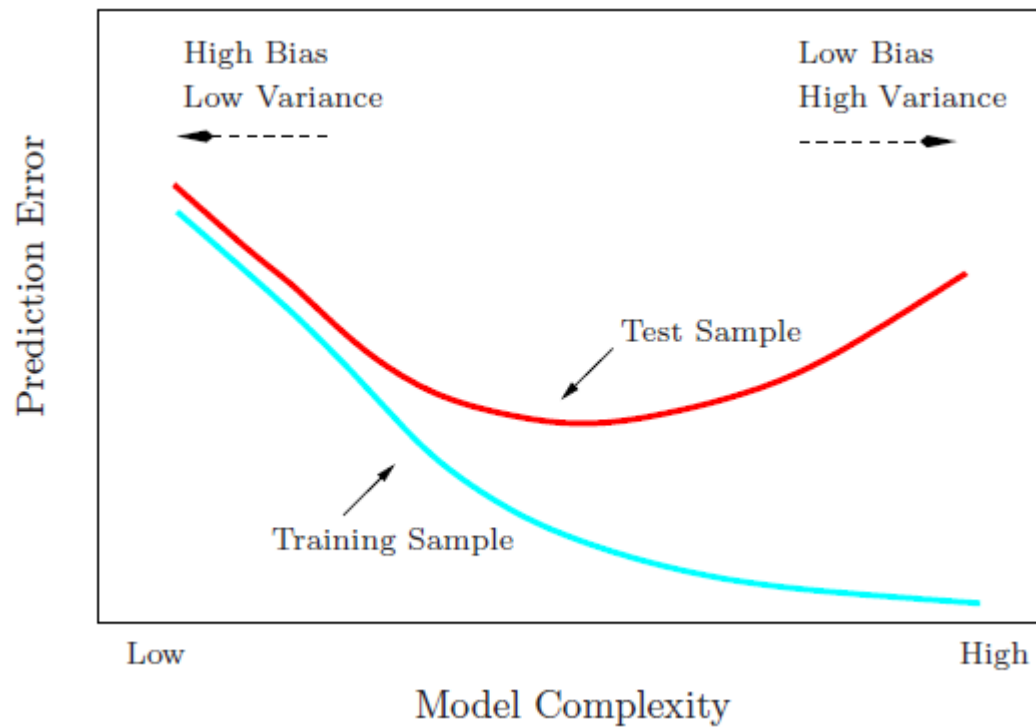
For a binary classifier, A ROC curve is a plot of the TPR (sensitivity) against the FPR ($1 - \text{specificity}$), which can be obtained at different setting thresholds.



The Area Under the ROC Curve (AUC) gives a quantification of the classifier performance, with a higher statistical consistency than accuracy.

UNDERFITTING, OVERFITTING AND...

... Best Fitting

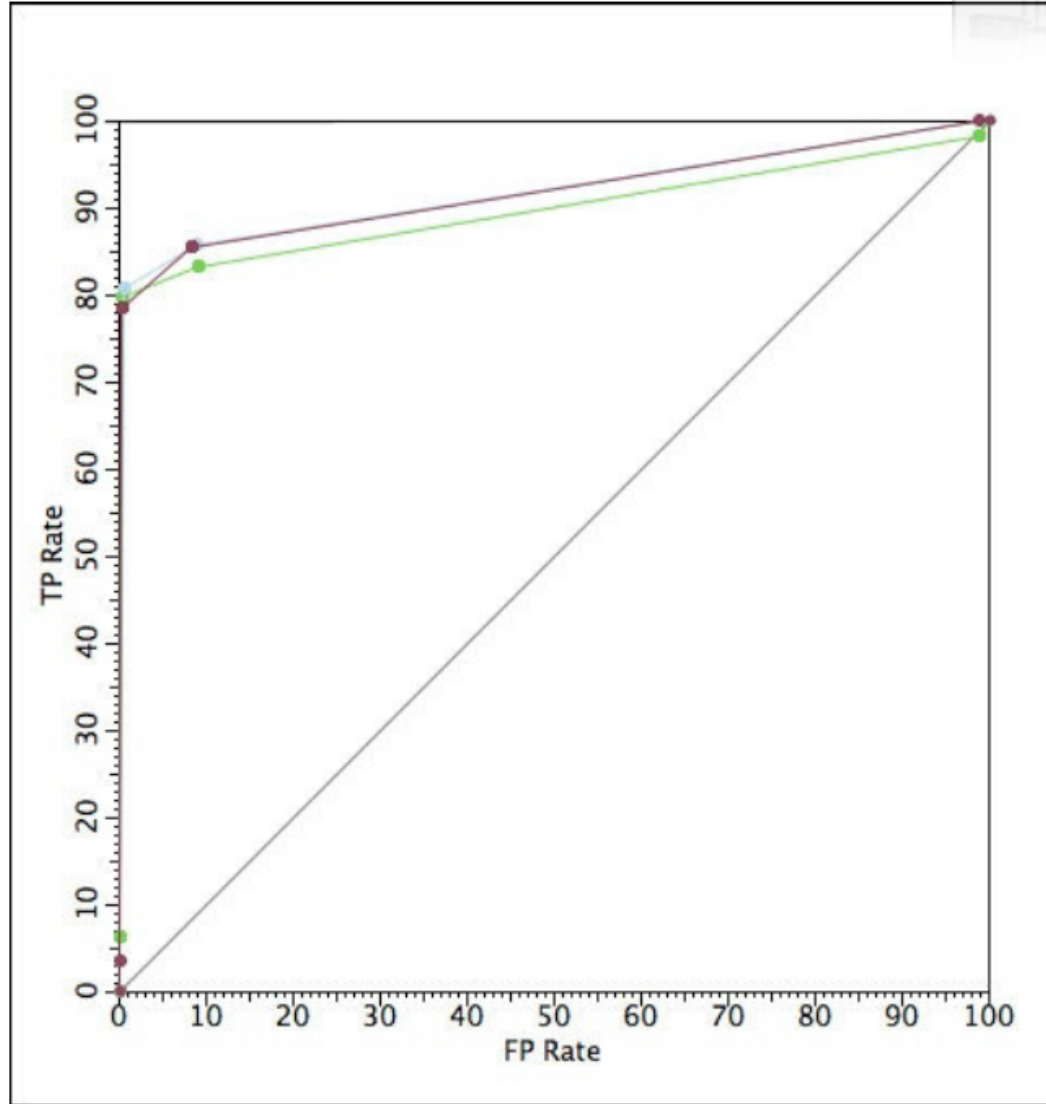


	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none"> • High training error • Training error close to test error • High bias 	<ul style="list-style-type: none"> • Training error slightly lower than test error 	<ul style="list-style-type: none"> • Very low training error • Training error much lower than test error • High variance
Regression illustration			
Classification illustration			
Deep learning illustration			
Possible remedies	<ul style="list-style-type: none"> • Complexity model • Add more features • Train longer 		<ul style="list-style-type: none"> • Perform regularization • Get more data

STANFORD.EDU

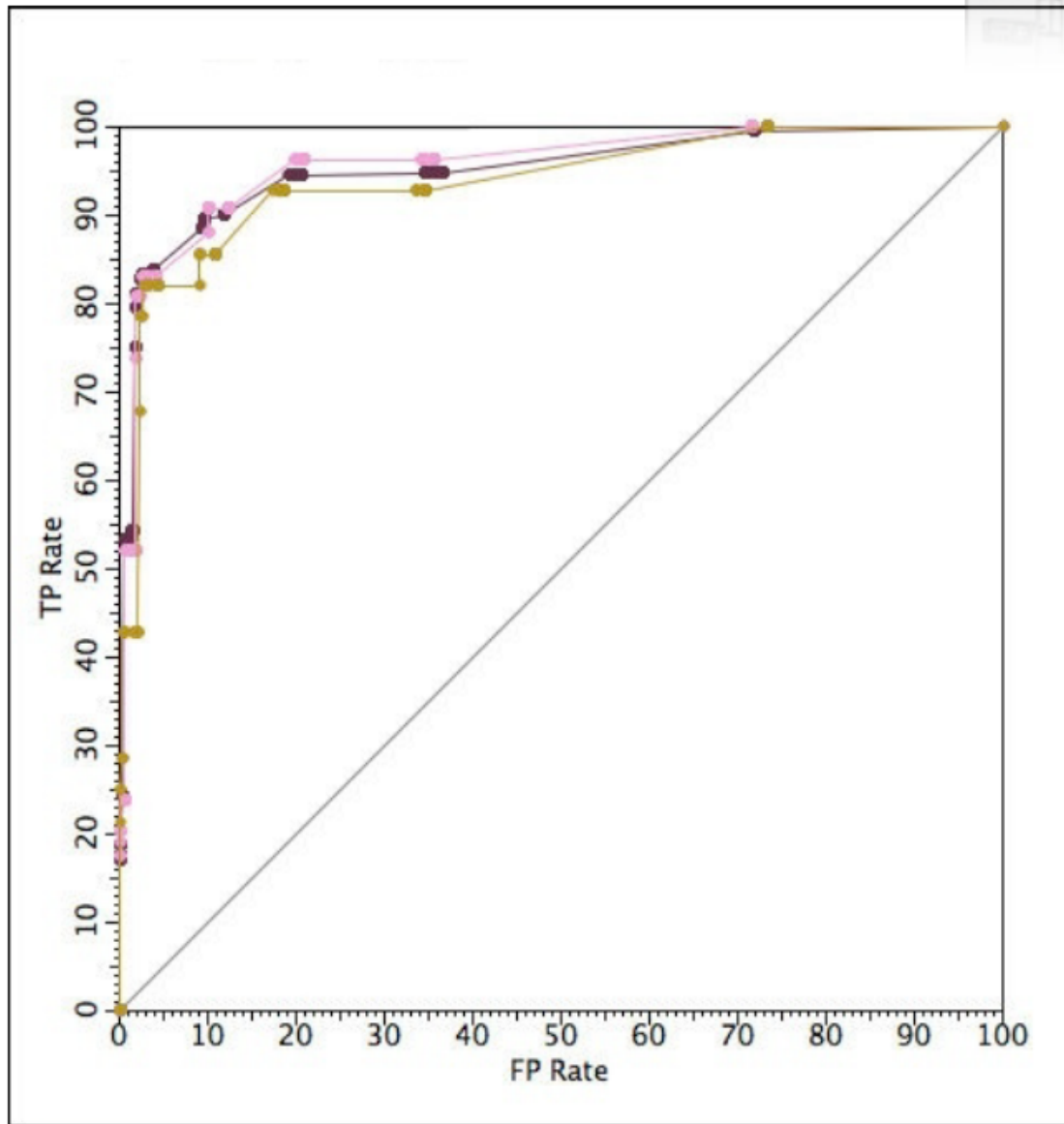
QUIZ

Which is which?



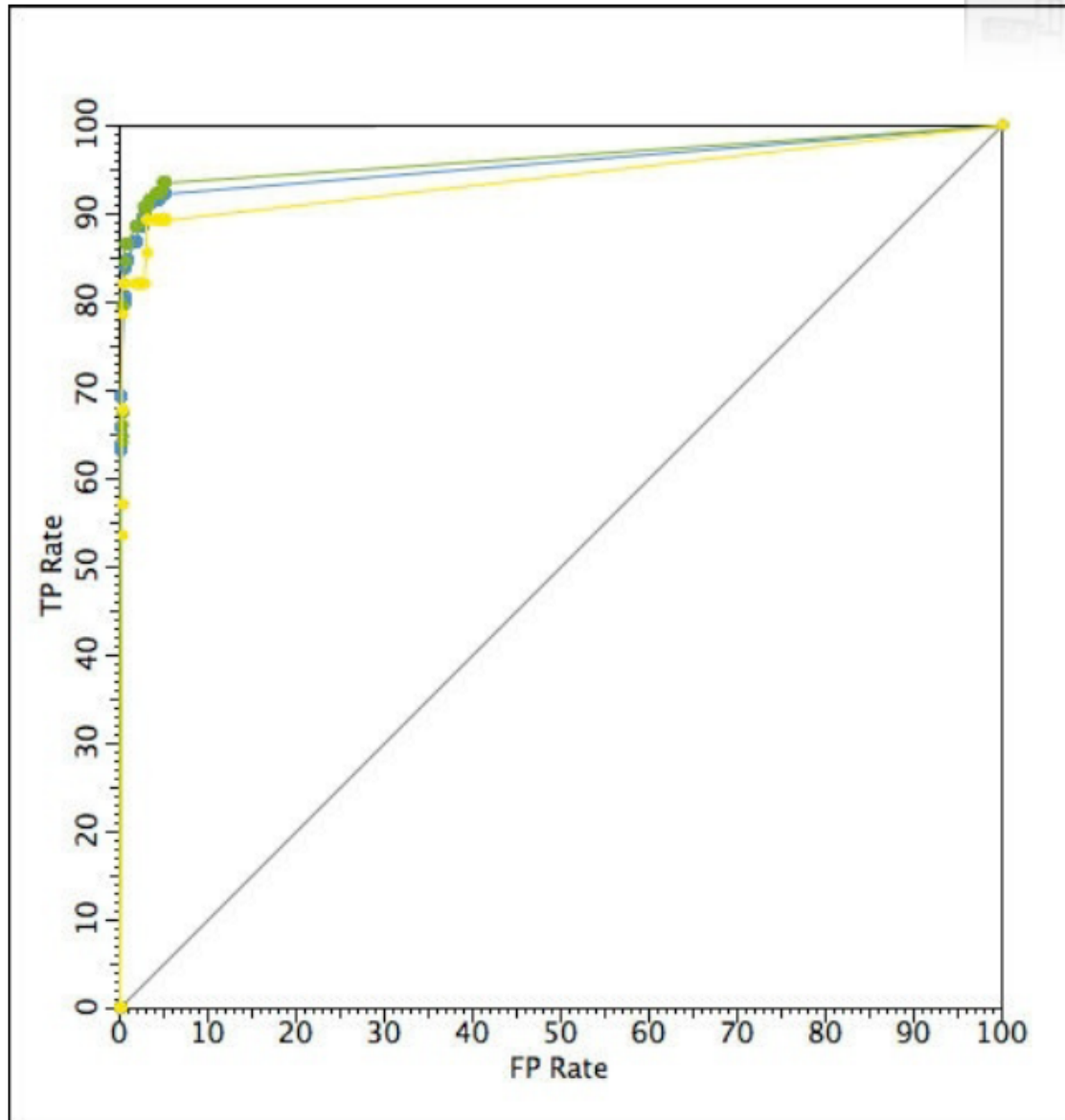
- ❖ Four models:
 - ❖ decision tree
 - ❖ k-nearest neighbour
 - ❖ linear classifier
 - ❖ naive Bayes
- ❖ trained on 2,000 examples and evaluated on
 - ❖ 18,000 test examples
 - ❖ 3,600 of those (20%)
 - ❖ 720 of those (4%)

Which is which?



- ❖ Four models:
 - ❖ decision tree
 - ❖ k-nearest neighbour
 - ❖ linear classifier
 - ❖ naive Bayes
- ❖ trained on 2,000 examples and evaluated on
 - ❖ 18,000 test examples
 - ❖ 3,600 of those (20%)
 - ❖ 720 of those (4%)

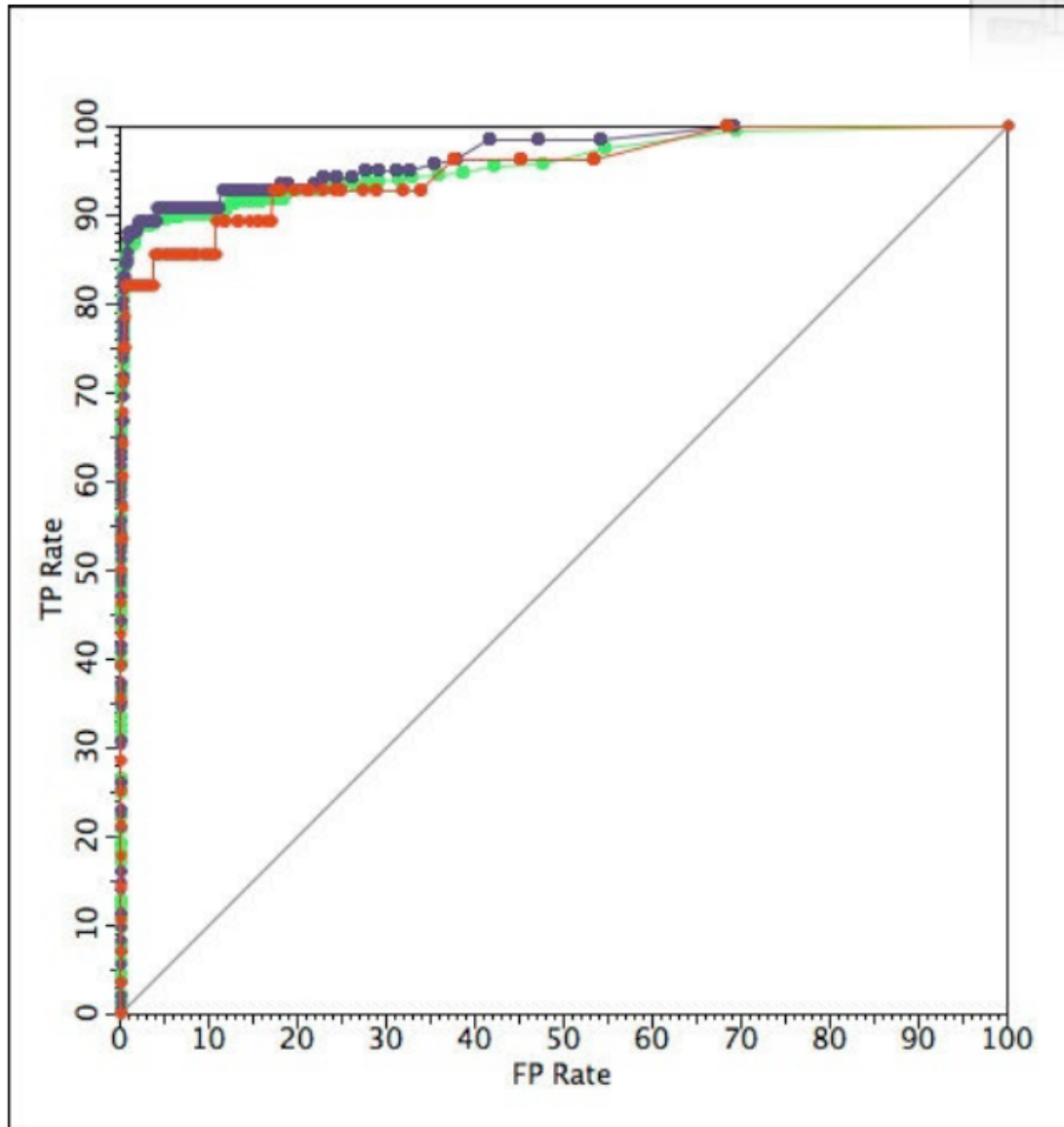
Which is which?



❖ Four models:

- ❖ decision tree
 - ❖ k-nearest neighbour
 - ❖ linear classifier
 - ❖ naive Bayes
- ❖ trained on 2,000 examples and evaluated on
- ❖ 18,000 test examples
 - ❖ 3,600 of those (20%)
 - ❖ 720 of those (4%)

Which is which?



❖ Four models:

❖ decision tree

❖ k-nearest neighbour

❖ linear classifier

❖ naive Bayes

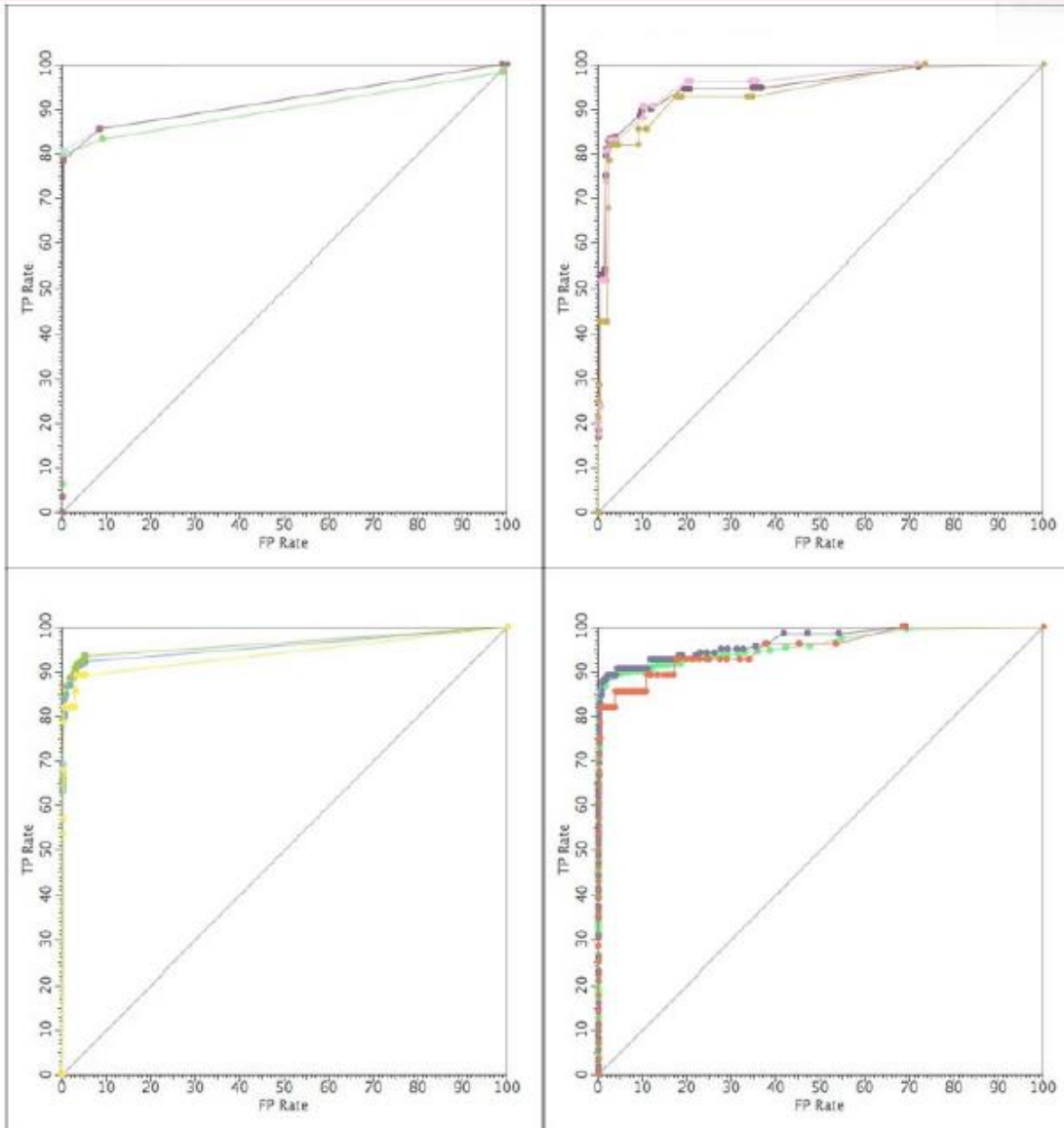
❖ trained on 2,000 examples
and evaluated on

❖ 18,000 test examples

❖ 3,600 of those (20%)

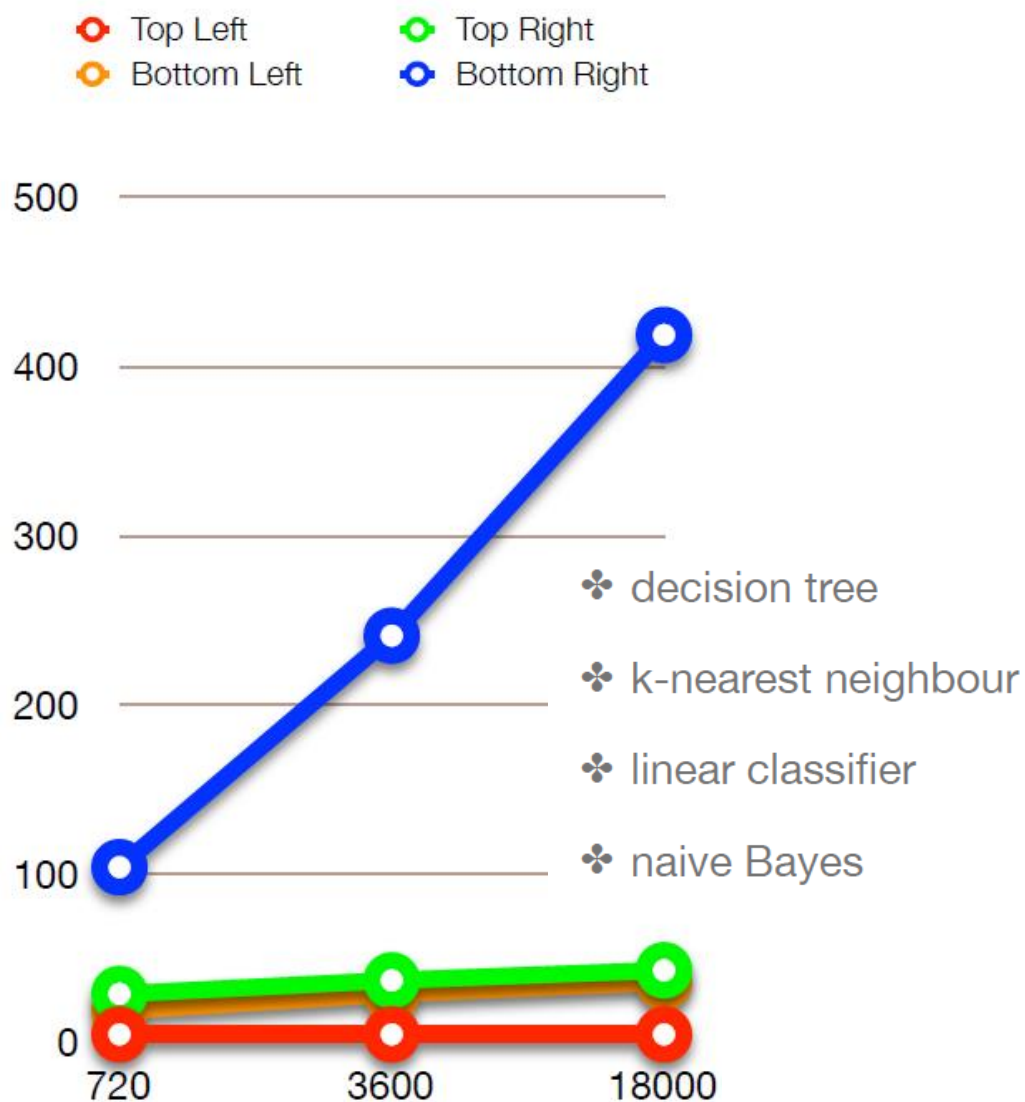
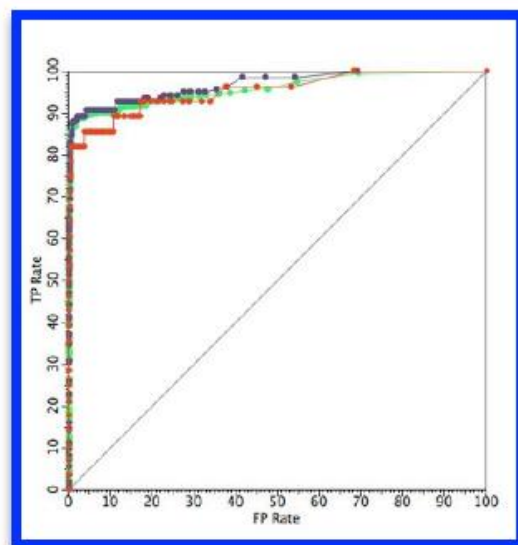
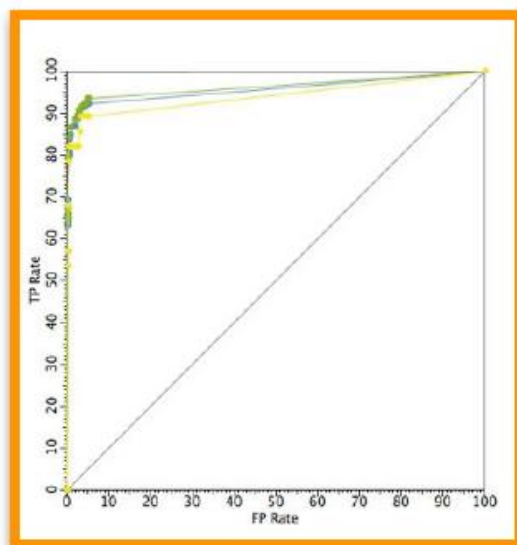
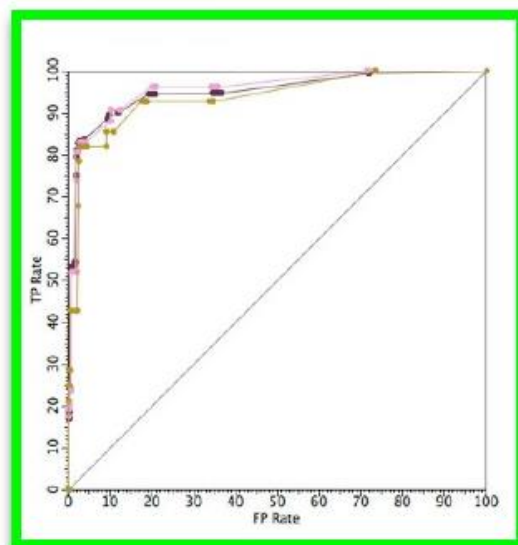
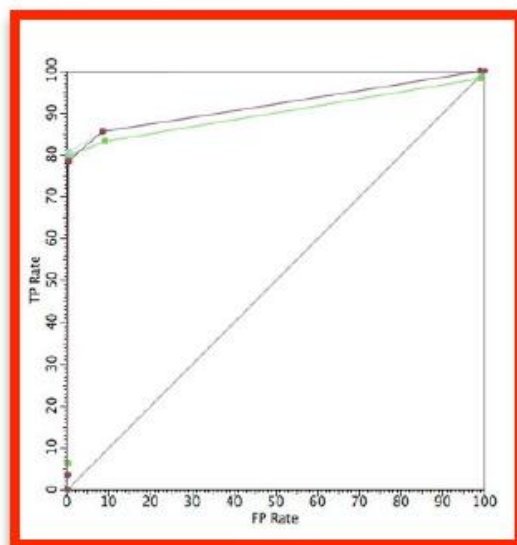
❖ 720 of those (4%)

Which is which?



- ❖ Four models:
 - ❖ decision tree
 - ❖ k-nearest neighbour
 - ❖ linear classifier
 - ❖ naive Bayes
- ❖ trained on 2,000 examples and evaluated on
 - ❖ 18,000 test examples
 - ❖ 3,600 of those (20%)
 - ❖ 720 of those (4%)

Which is which?



christian.salvatore@iusspavia.it

<https://christiansalvatore.github.io/machinelearning-iuss/>