

Random Forest and Applications

Sambit Panda

Question: What to do with all this data?

Mouse	Group	Locomotor Activity	Grip Strength
10	WT	24	0.1
7.2	Mutant	18	0.12
34	Mutant	17	0.22
20	WT	33	0.05
8.8	Mutant	15	0.3



Mutant



WT

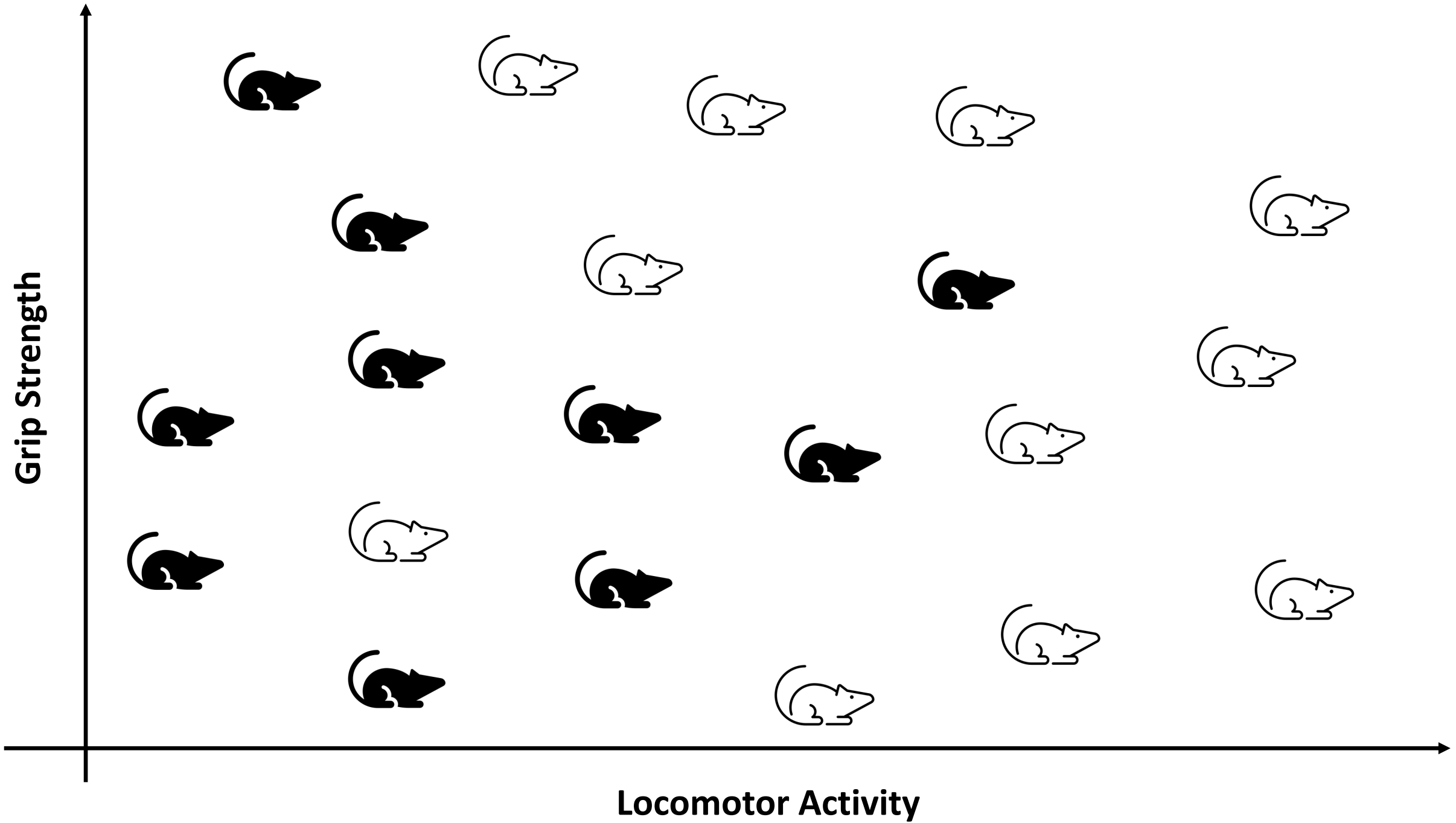
Columns = Features = Dimensions (p or d)

Rows = Subjects (n)

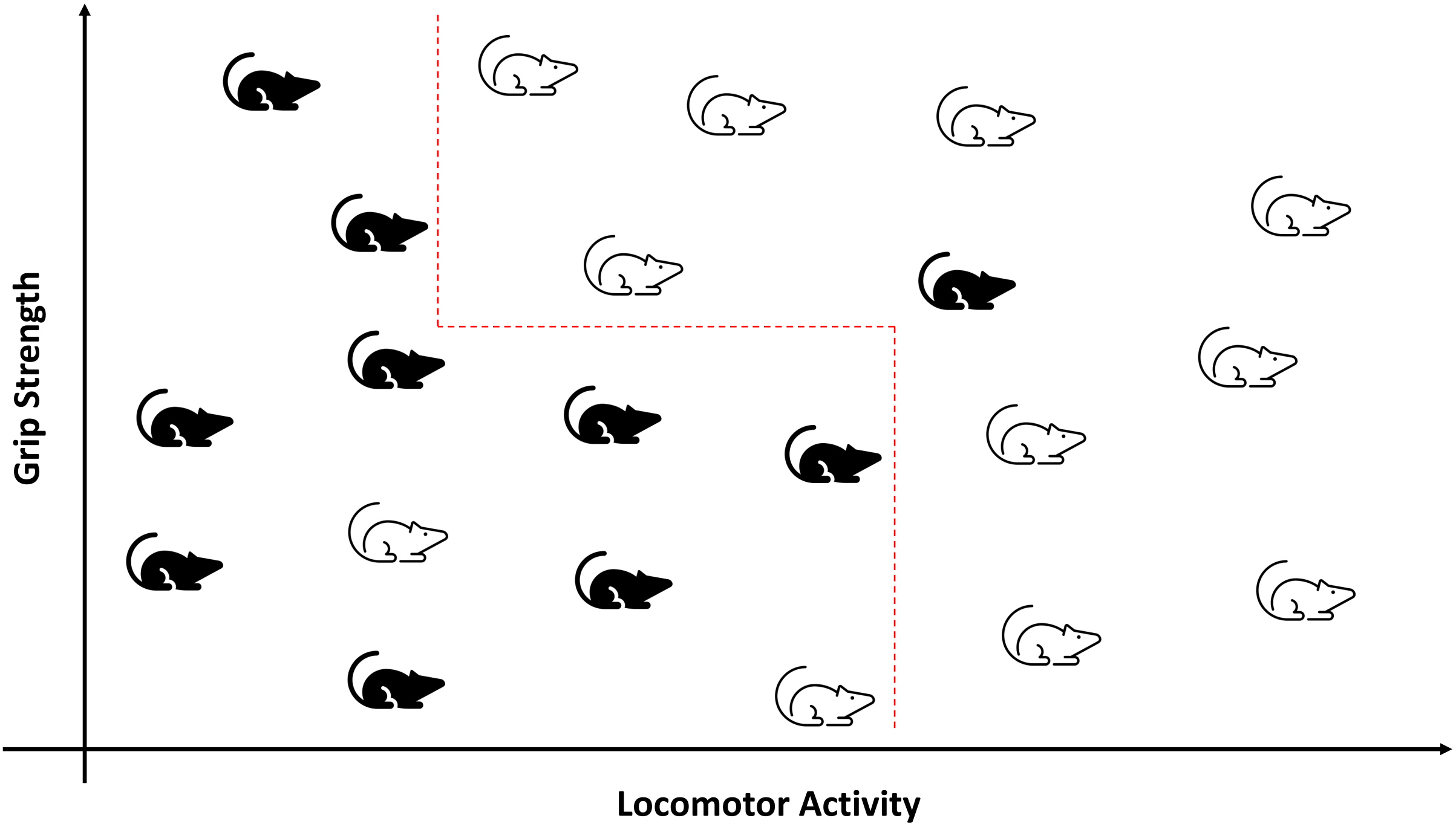
Mouse	Group	Locomotor Activity	Grip Strength
10	WT	24	0.1
7.2	Mutant	18	0.12
34	Mutant	17	0.22
20	WT	33	0.05
8.8	Mutant	15	0.3

...

...



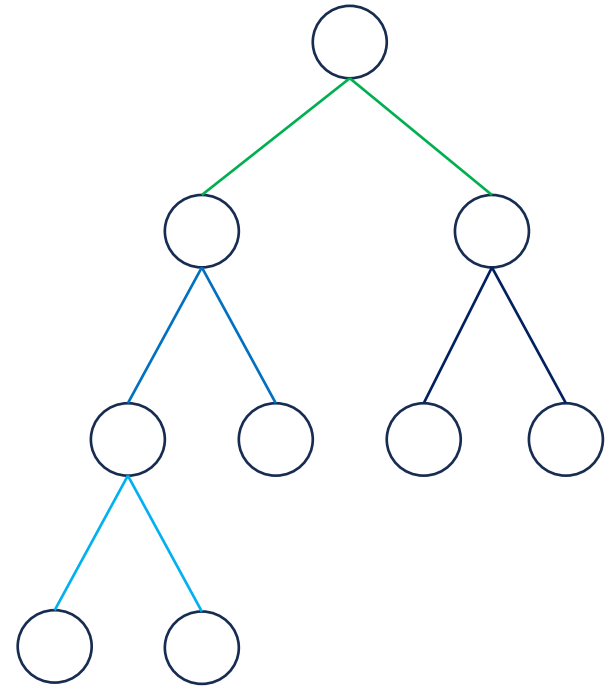
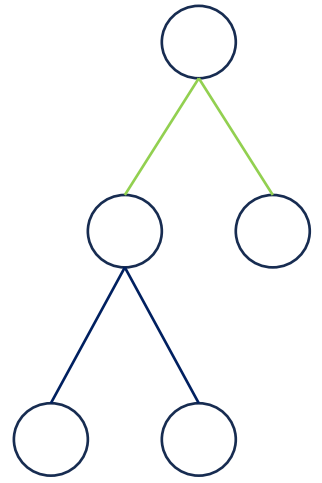
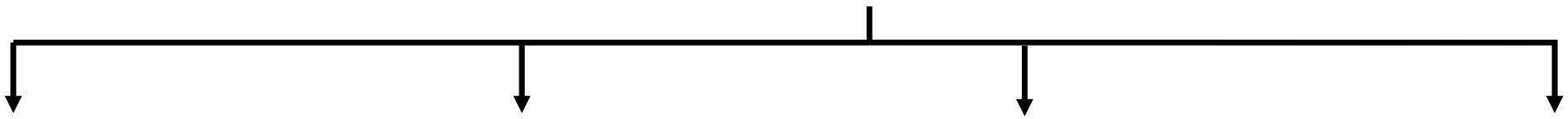
How Random Forest Works



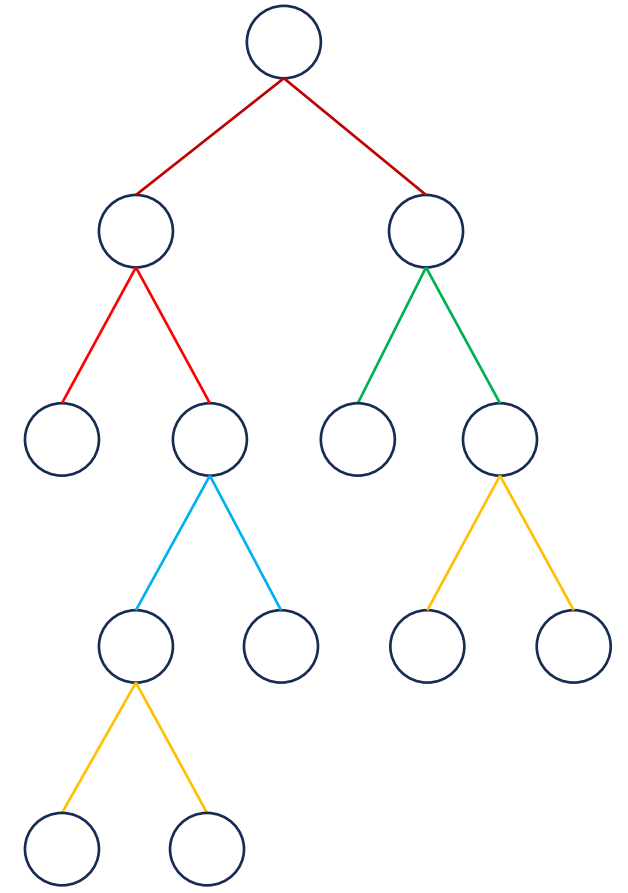
Mouse	Group	Locomotor Activity	Grip Strength
10	WT	24	0.1
7.2	Mutant	18	0.12
34	Mutant	17	0.22
20	WT	33	0.05
8.8	Mutant	15	0.3



Features



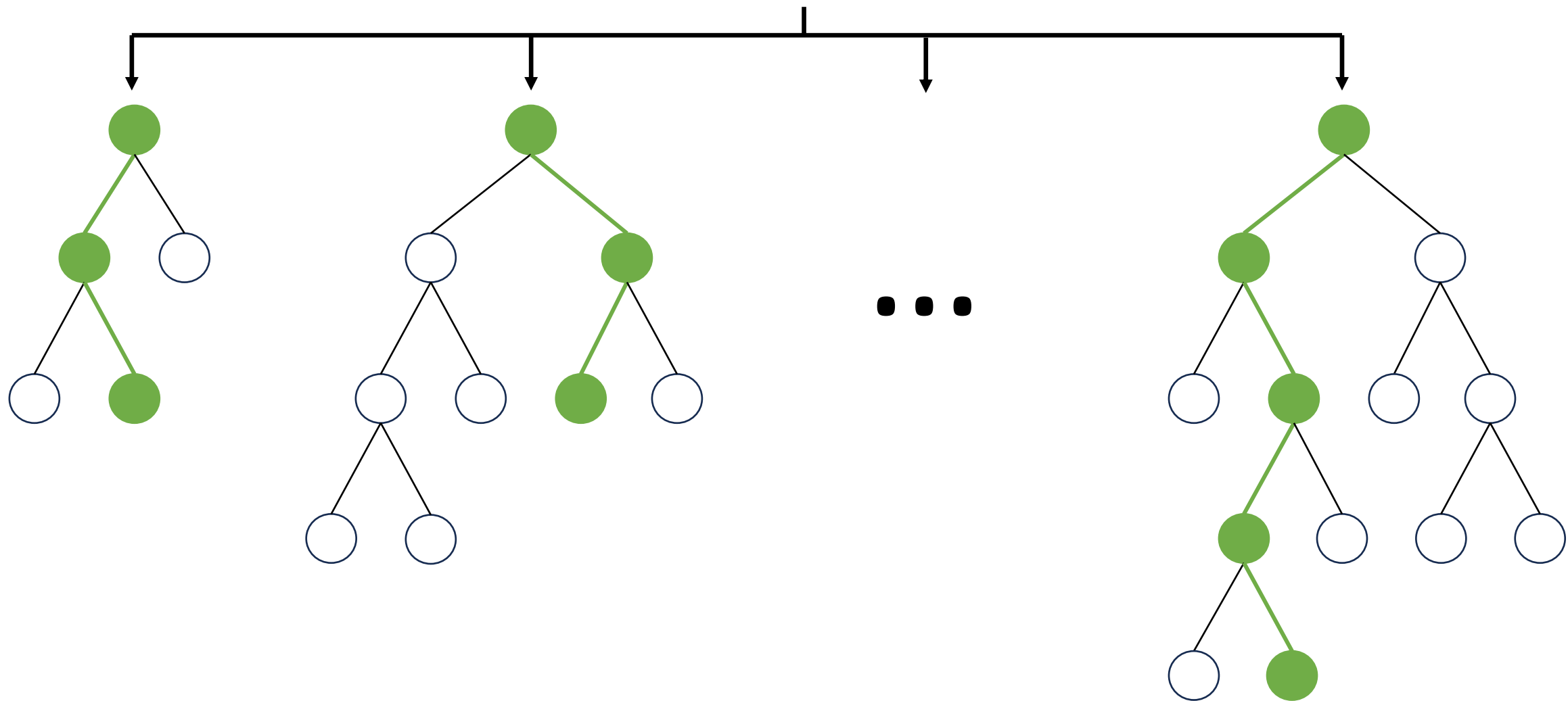
...



Tree 1

Tree 2

Tree T



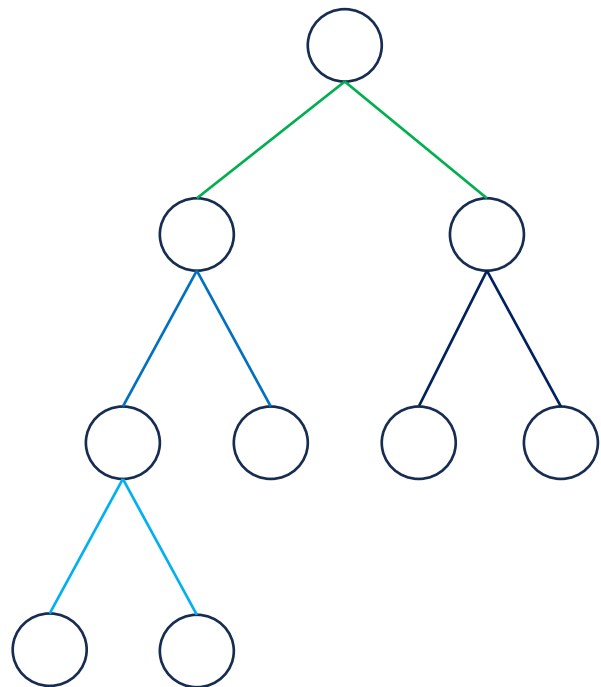
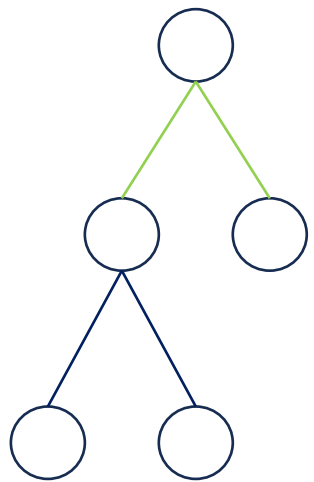
WT
 $p_1 = 0.55$

Mutant
 $p_2 = 0.76$

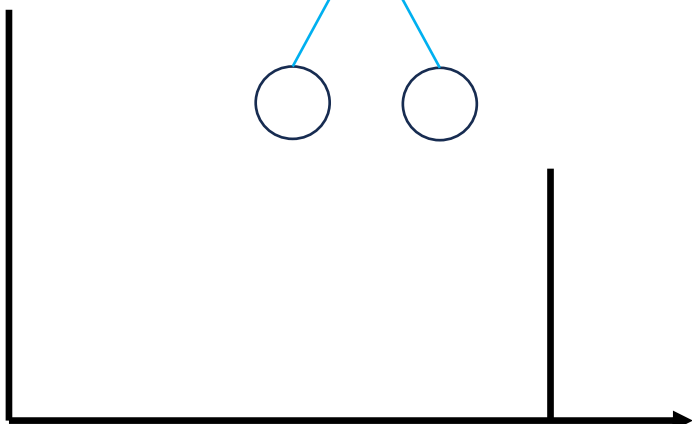
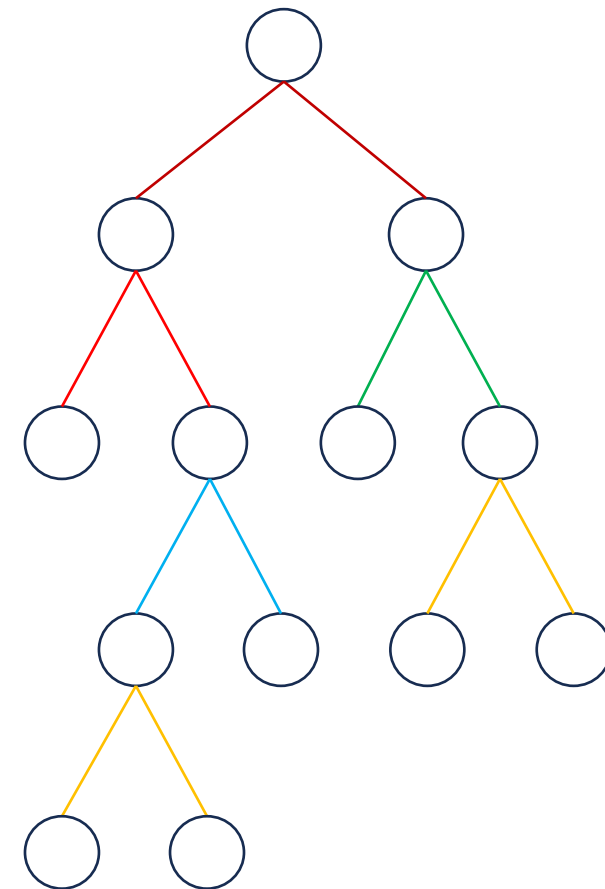
Mutant
 $p_T = 0.82$

Majority: Mutant
 $p = \text{Average } p_i$





...

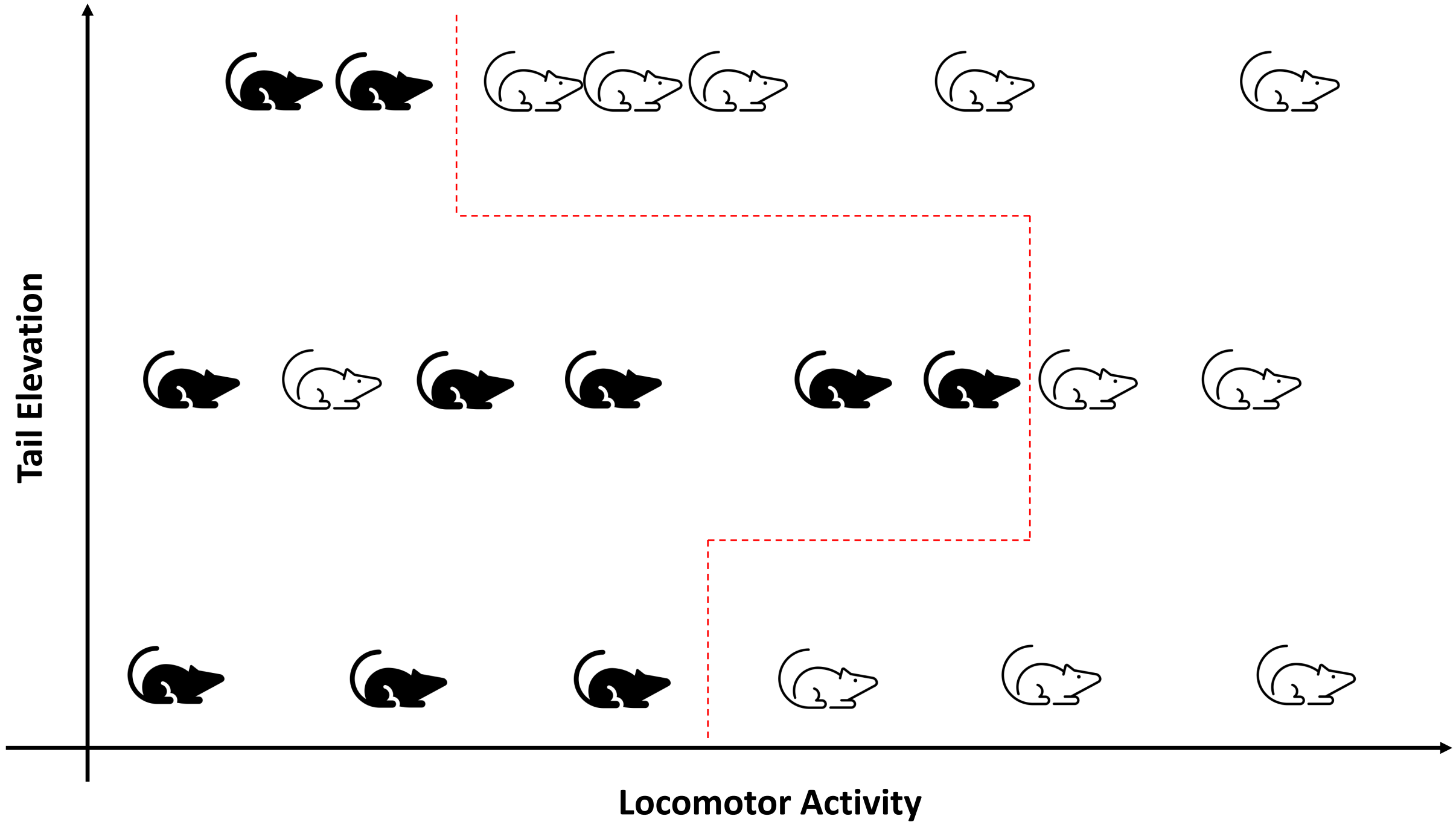


**Expected Fraction of the
Samples**

+

Decrease in Impurity





Summary

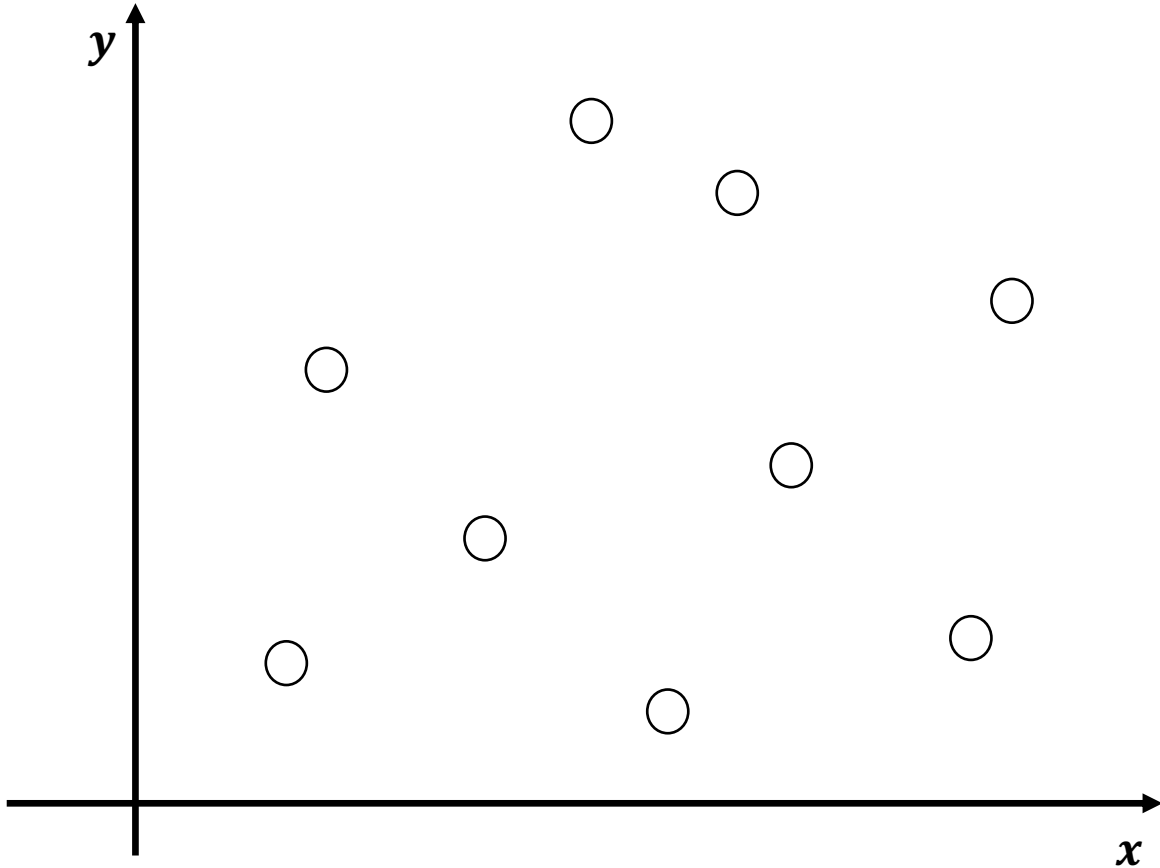
- Learned what tabular data is
- What splits are geometrically and within the algorithm
- How random forests trains on current data and tests on new points
- How random forest estimates feature importance

Independence Testing, K-Sample
Testing and Random Forest

Question: Is there a relationship?

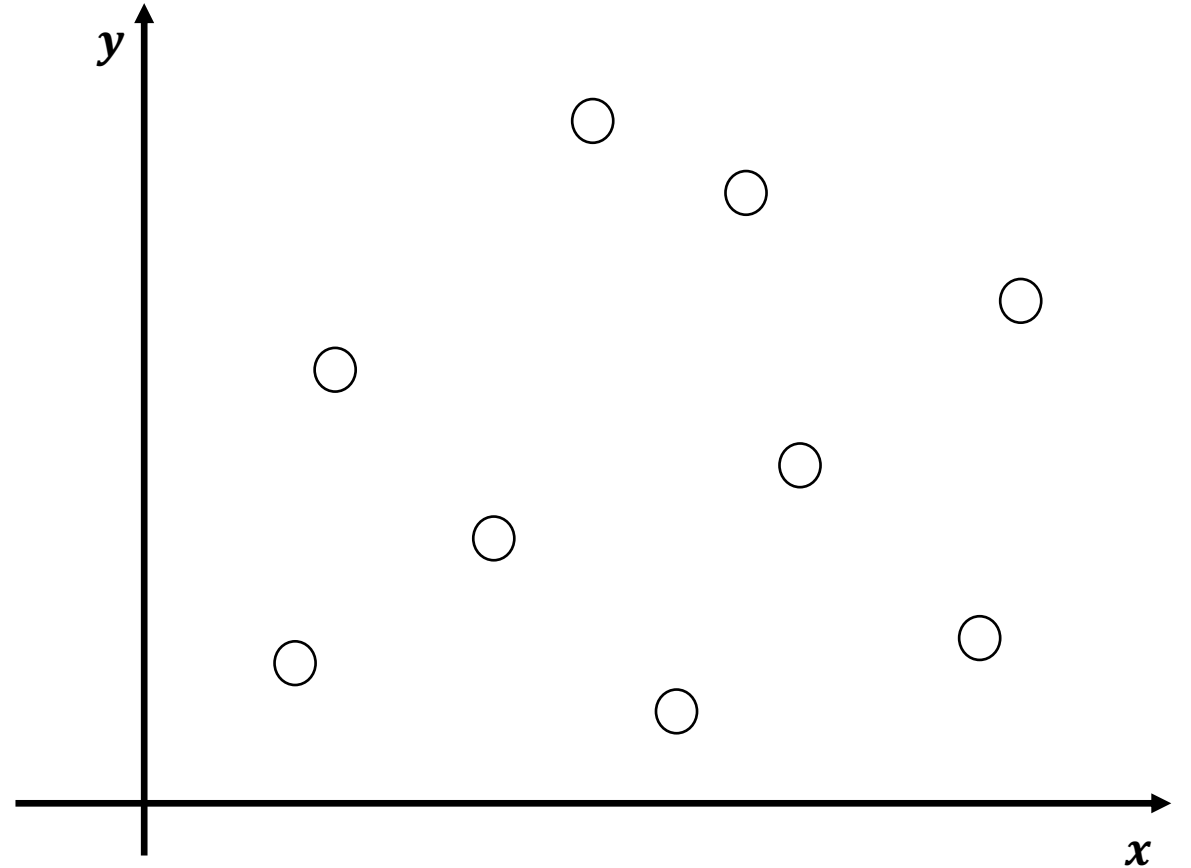
$$H_0 : \rho = 0$$

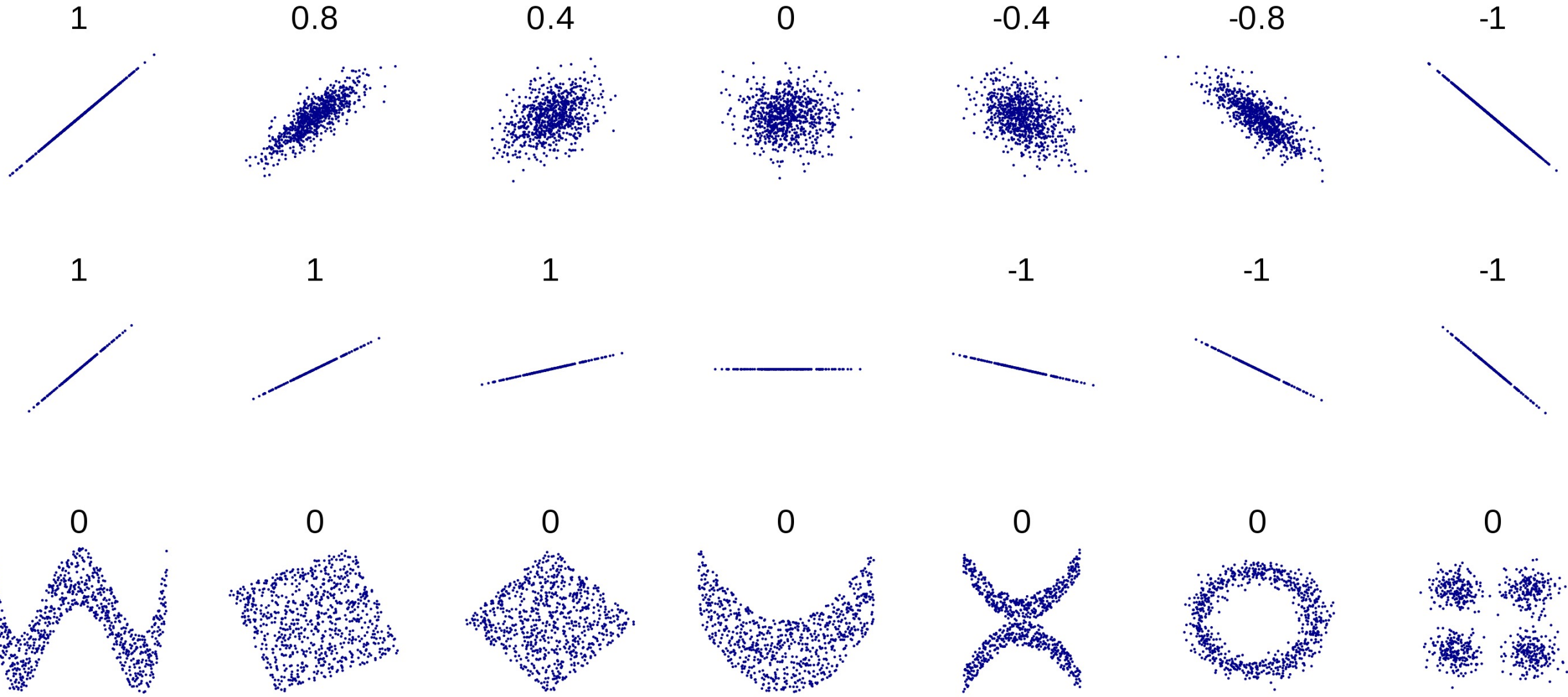
$$H_A : \rho \neq 0$$

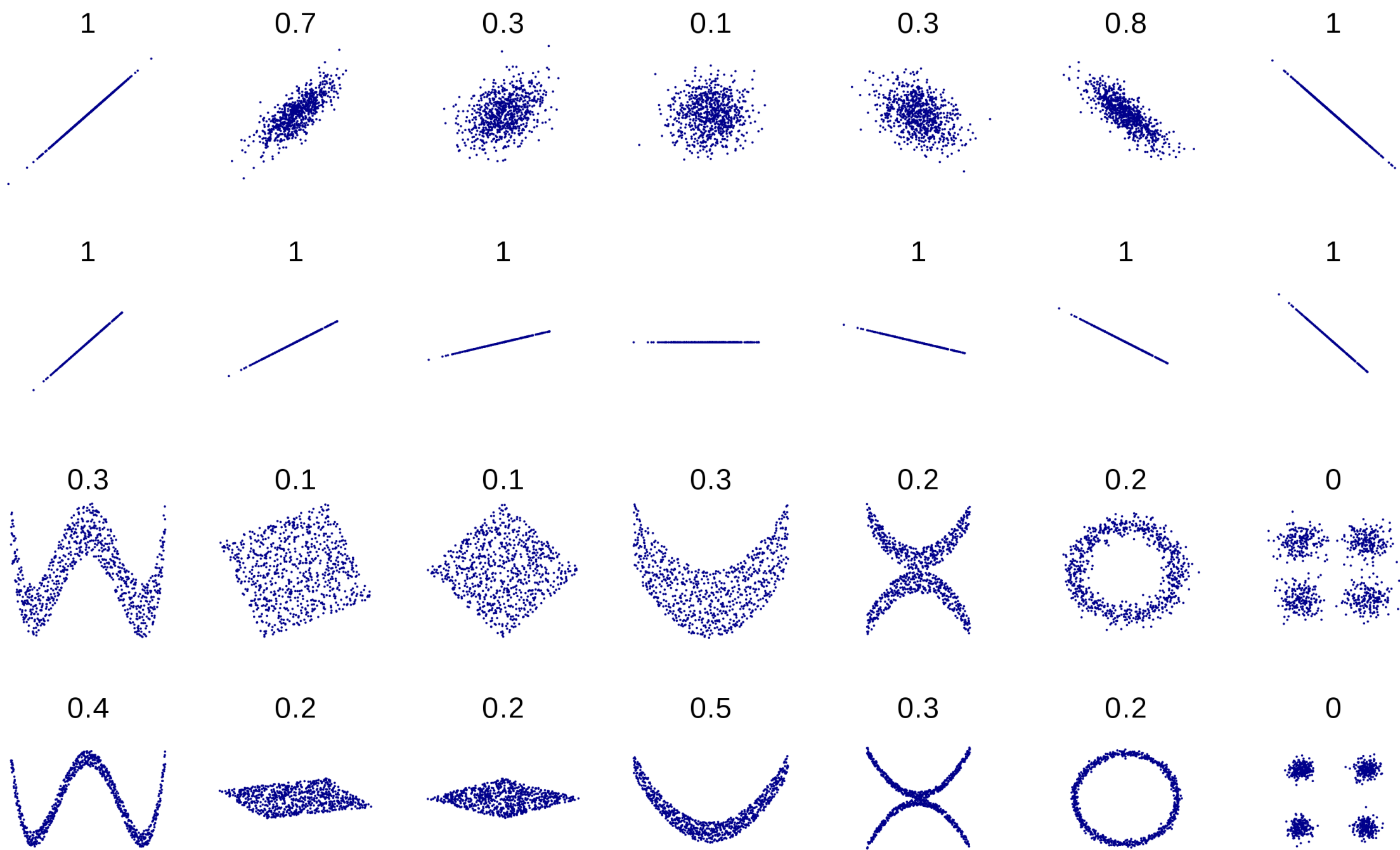


$$H_0 : F_{XY} = F_X F_Y$$

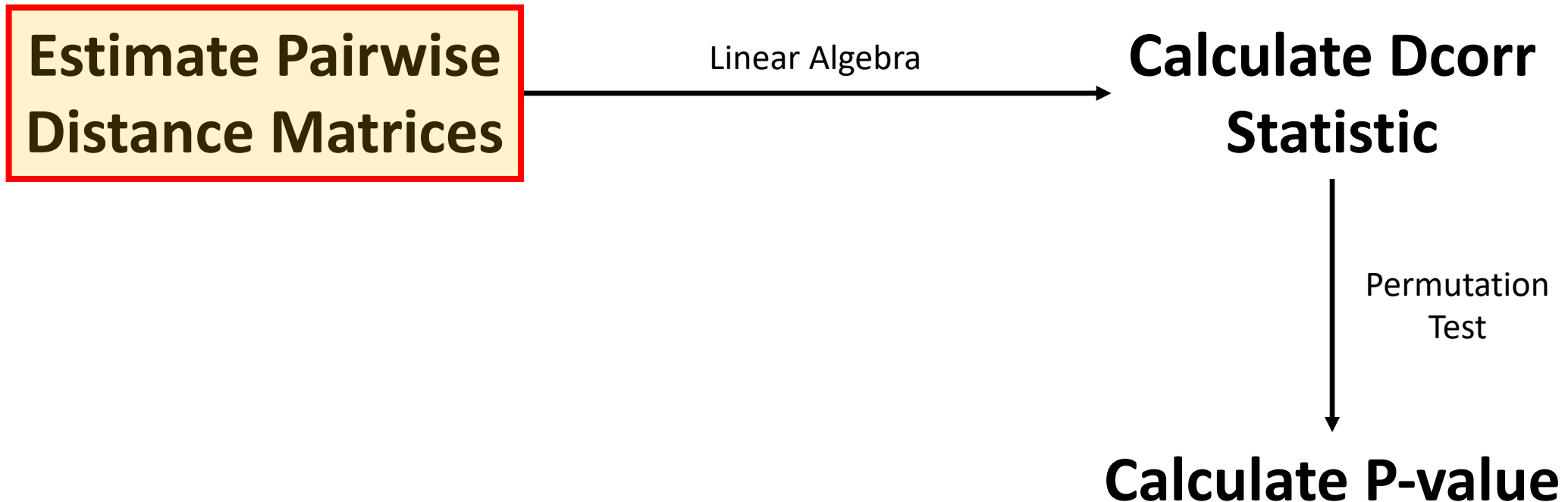
$$H_A : F_{XY} \neq F_X F_Y$$

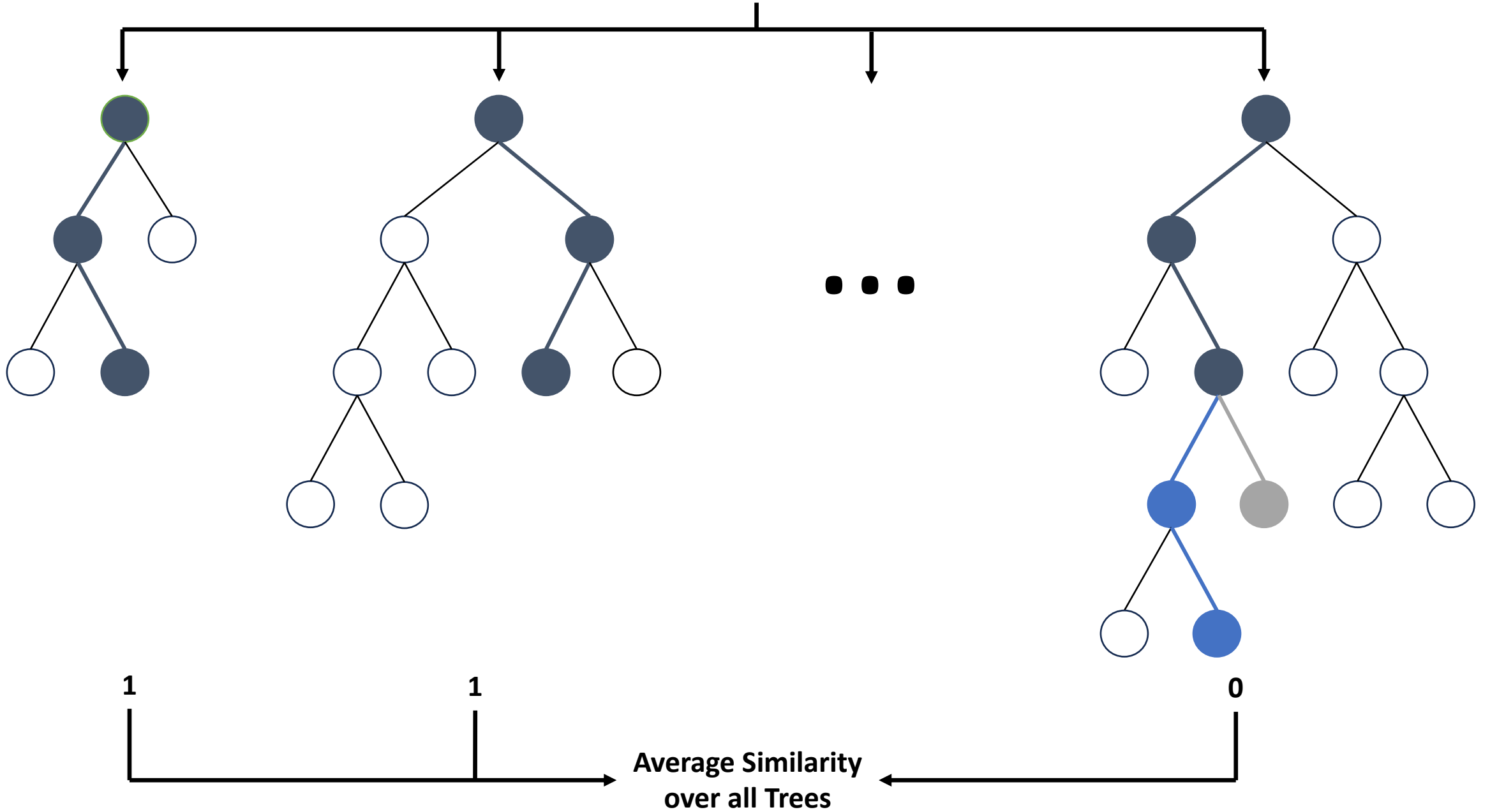






How to Compute Dcorr





How to Compute KMERF

**Compute
similarities via RF**

Exact
Equivalence

**Estimate Pairwise
Distance Matrices**

Linear Algebra

**Calculate Dcorr
Statistic**

RF Feature
Importance

**Calculate Feature
Importance**

Permutation
Test

Calculate P-value

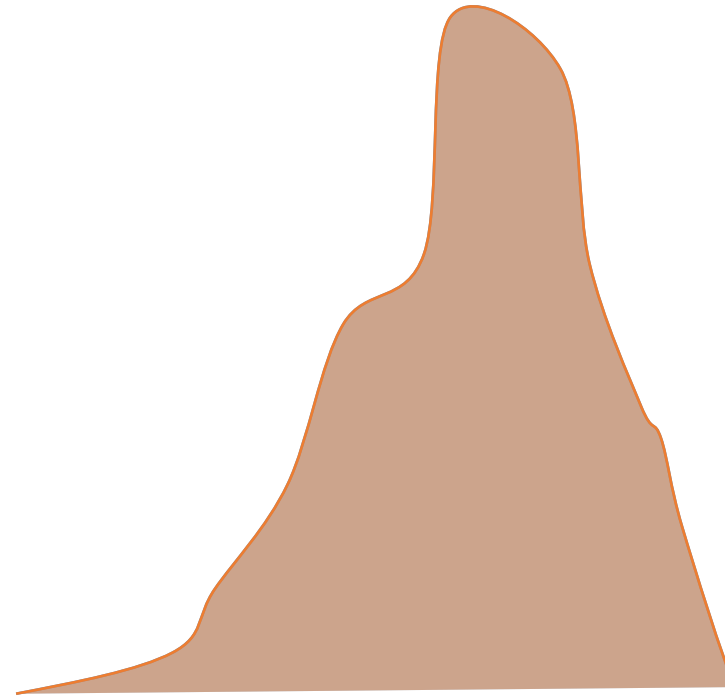
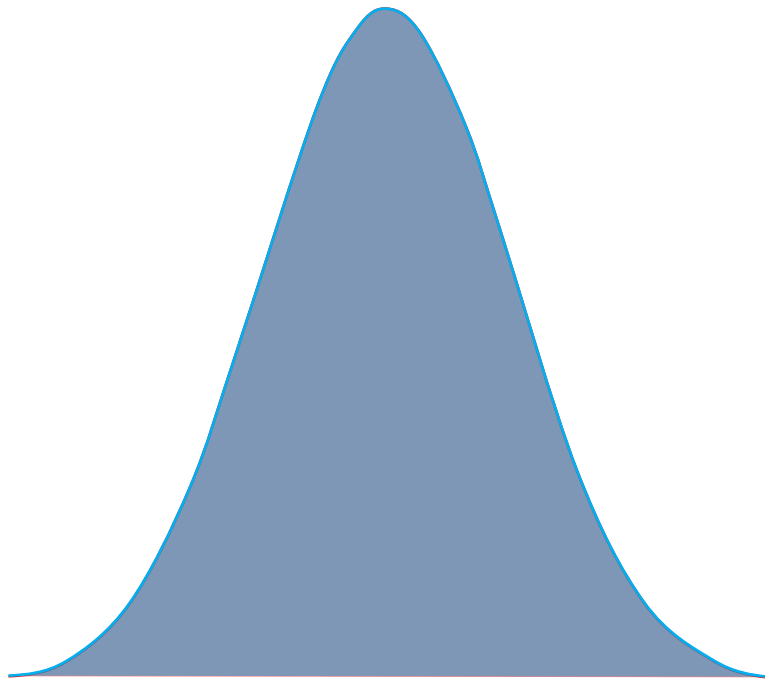
Question: Are they different?

$$H_0 : \mu_X = \mu_Y$$

$$H_A : \mu_X \neq \mu_Y$$

$$H_0 : F_X = F_Y$$

$$H_A : F_X \neq F_Y$$



Question: Is there a relationship?

=

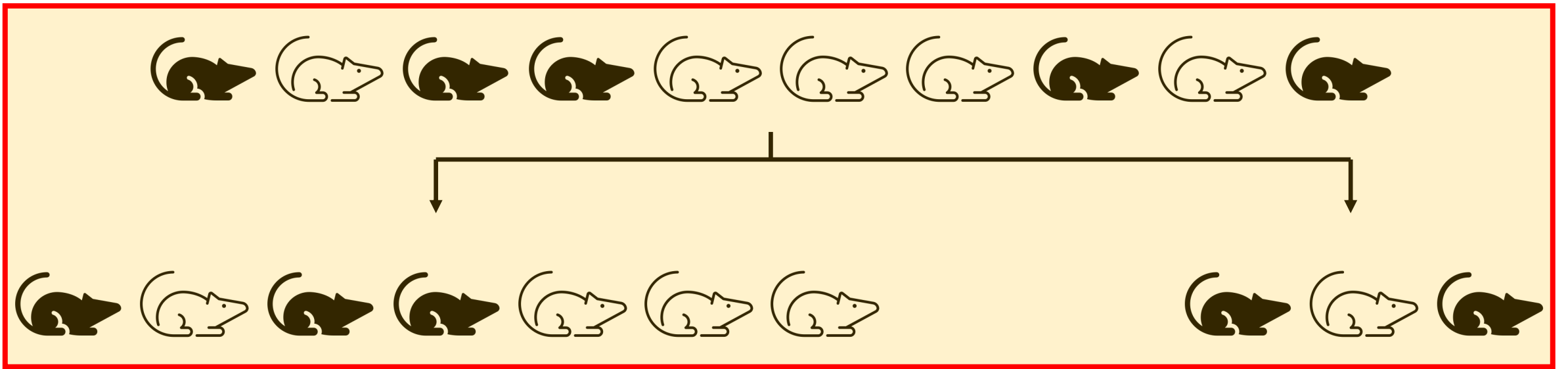
Question: Are they different?

(via a transformation of the data)

Summary

- Learned what independence and k-sample testing are
- Learned how Dcorr and KMERF works
- Found out what results are on our data

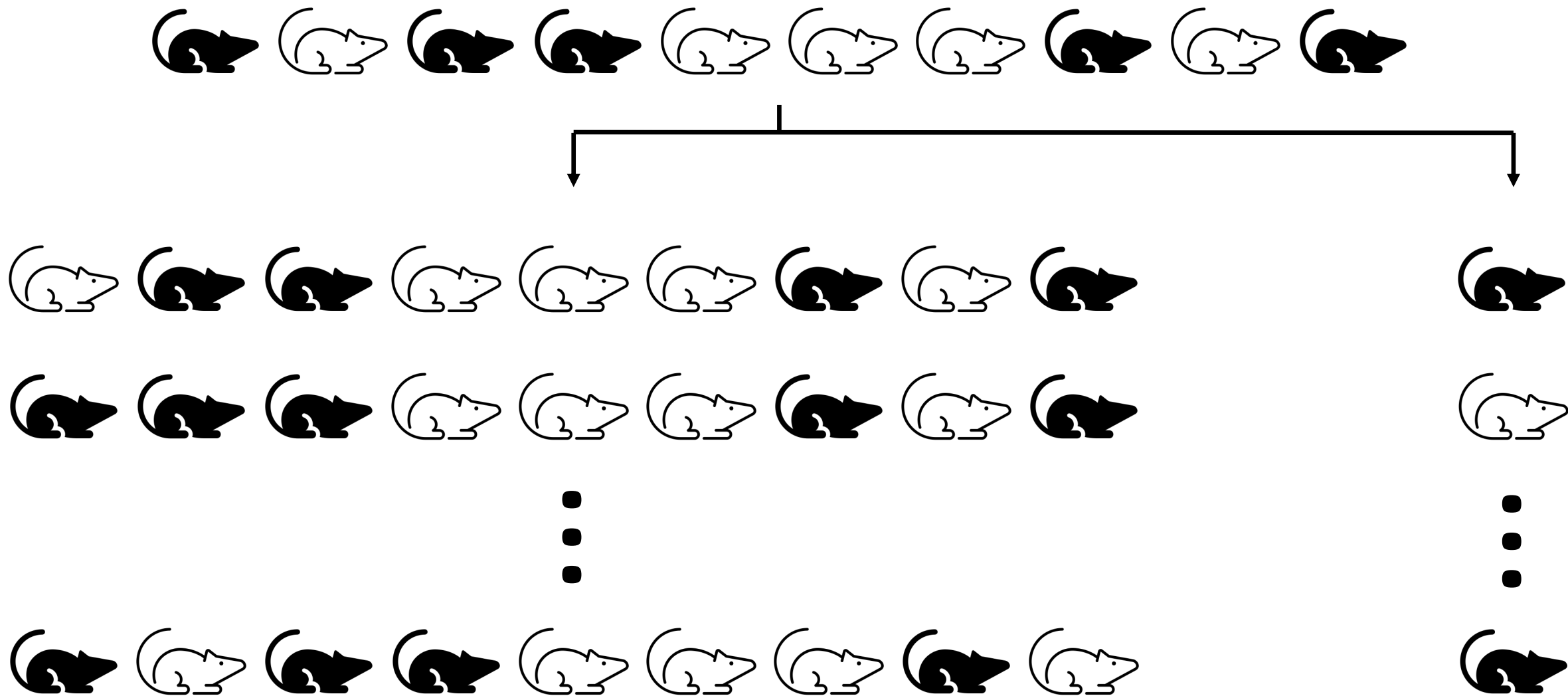
Methods to Evaluate Performance



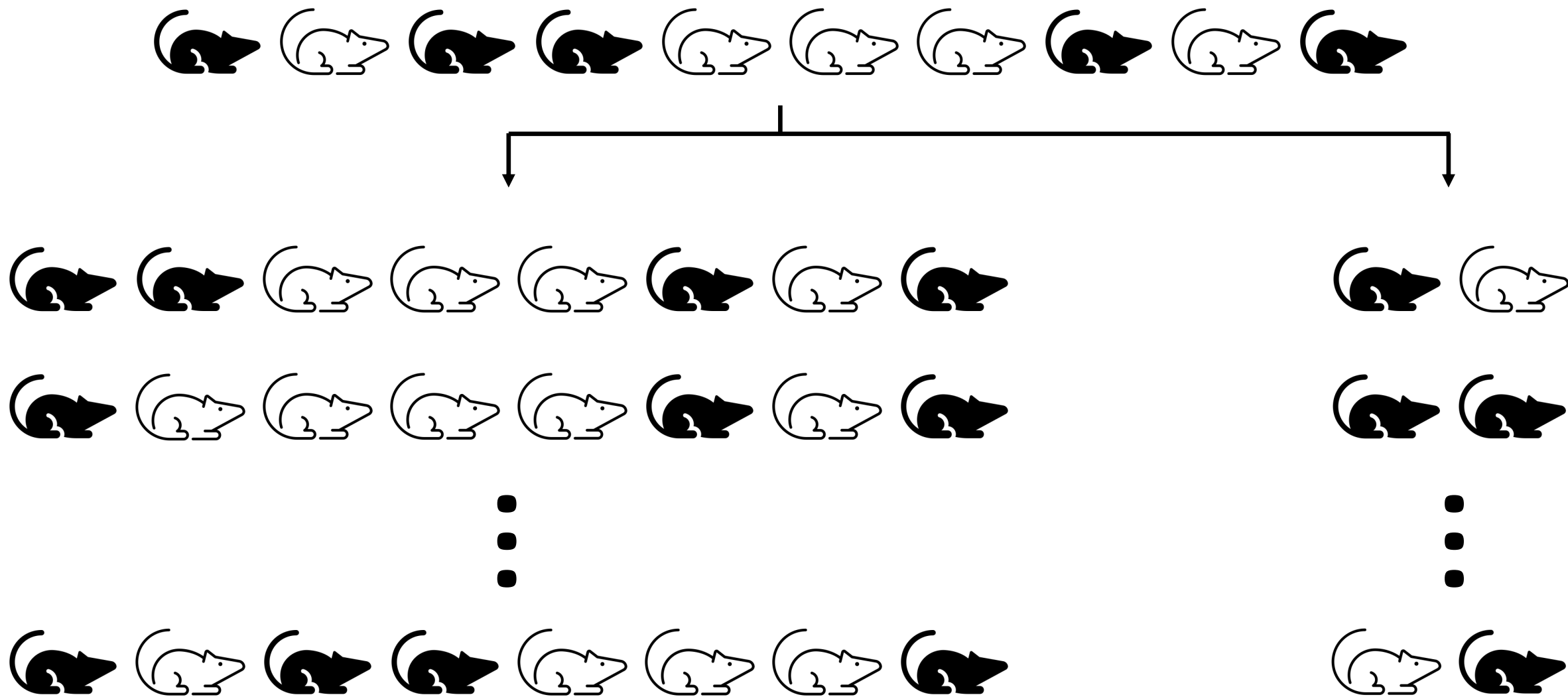
Train Classifier

Test Classifier

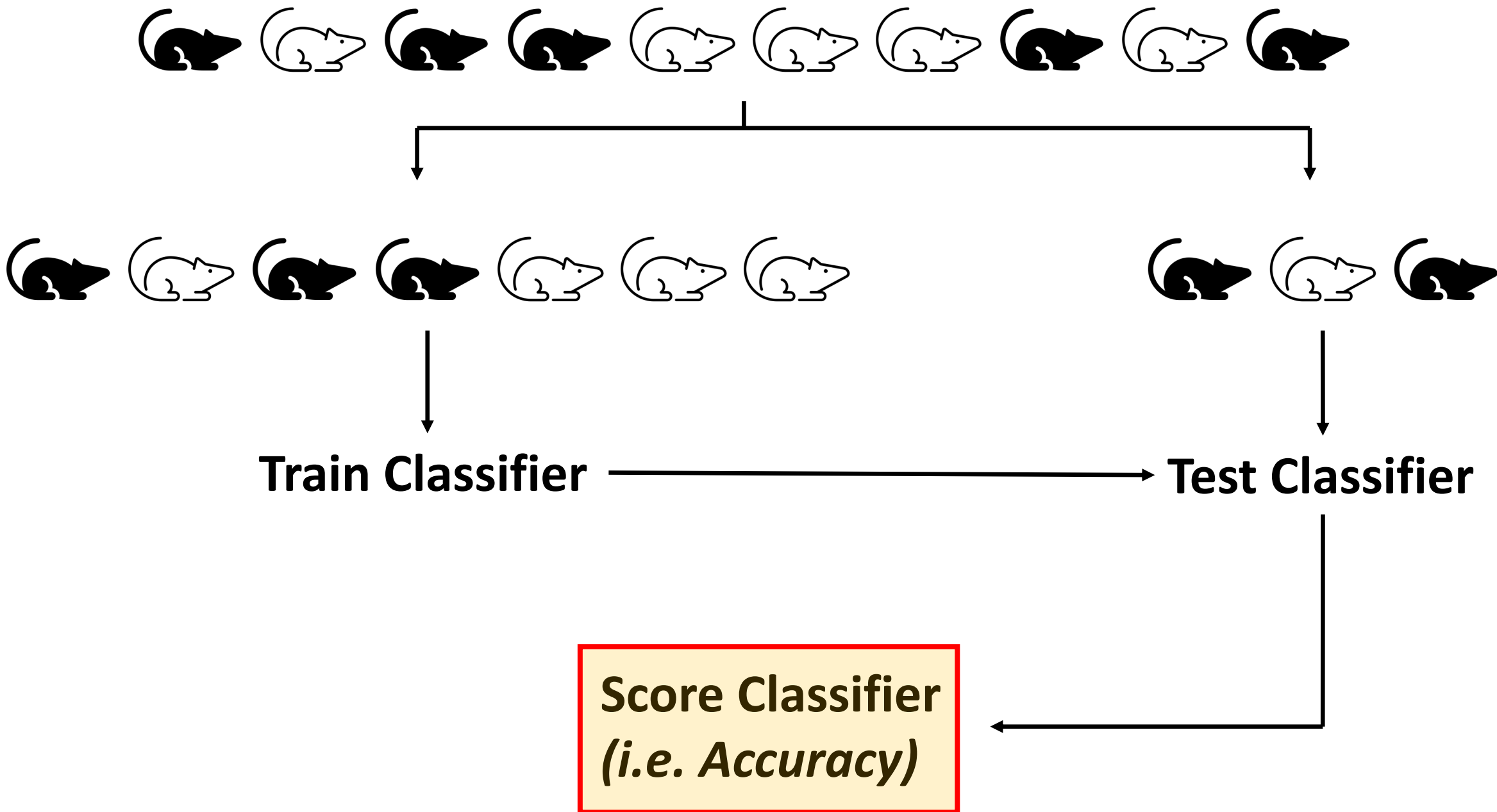
Score Classifier
(i.e. Accuracy)

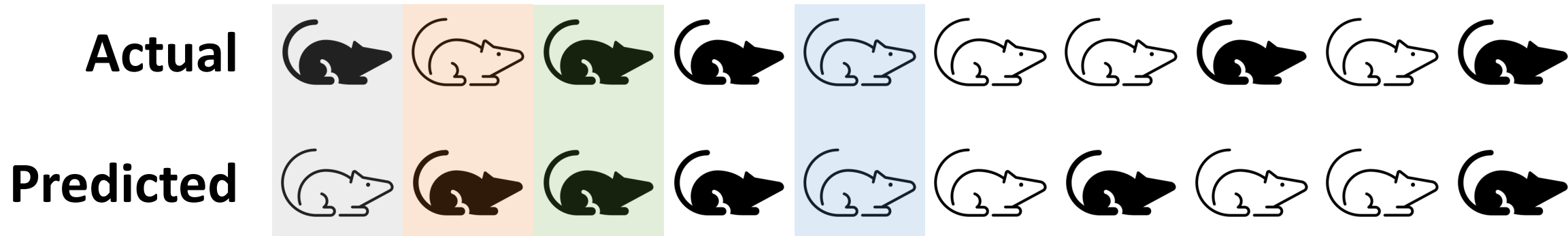


Leave One Out



K-Fold ($k = 5$)





Predicted

		Predicted	
		True (Pos)	False (Neg)
Actual	True (Pos)	True Positive (TP)	False Negative (FN)
	False (Neg)	False Positive (FP)	True Negative (TN)

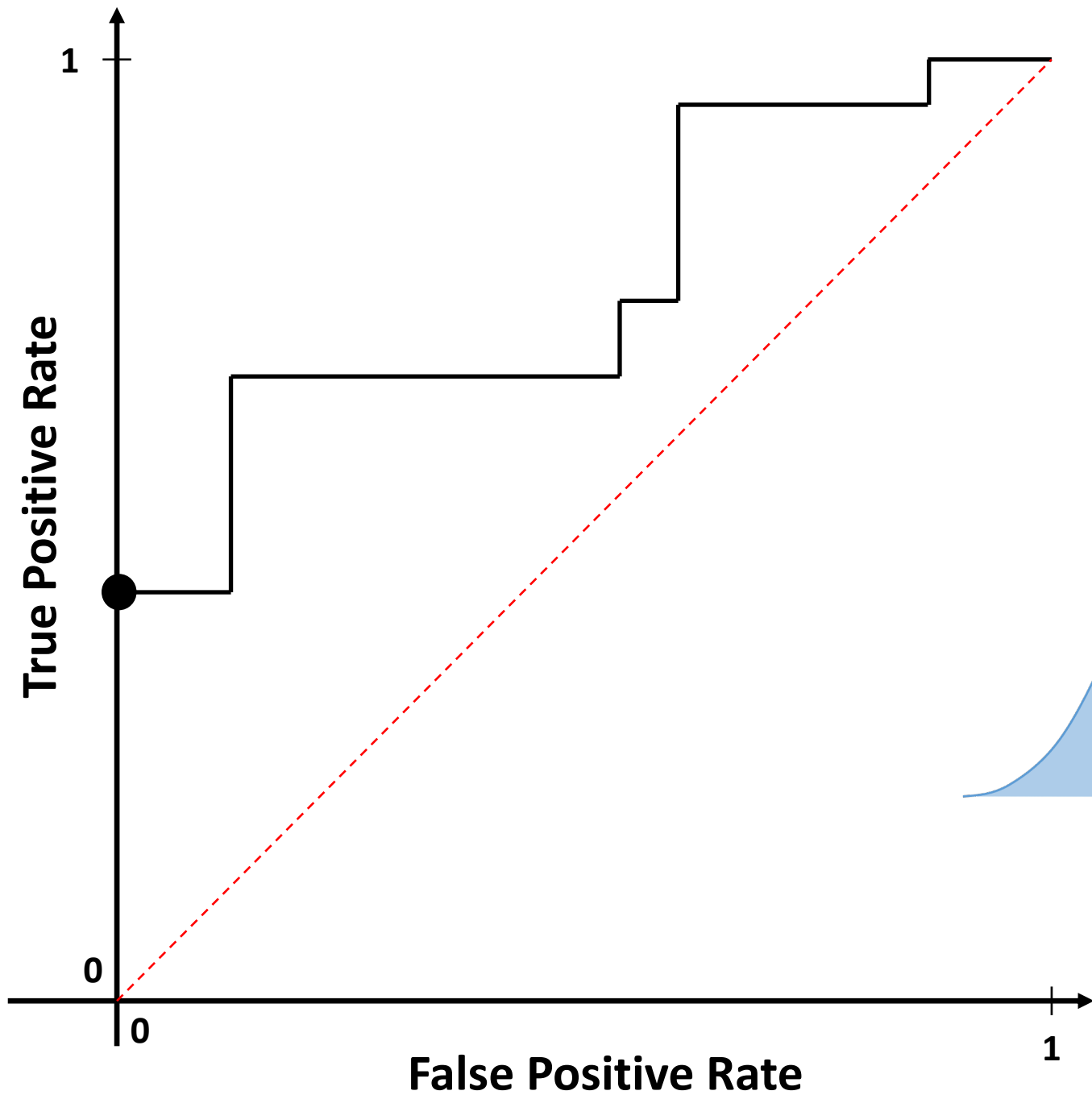


False Positive Rate
(1 – Specificity)

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

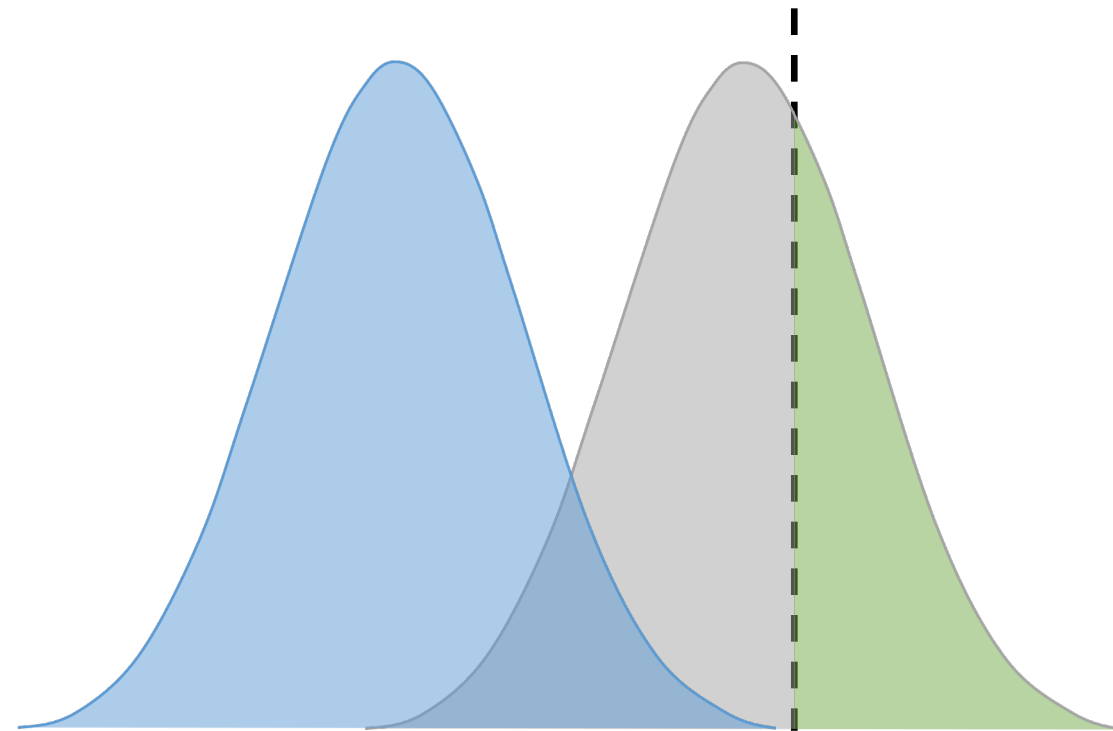
True Positive Rate
(Sensitivity)

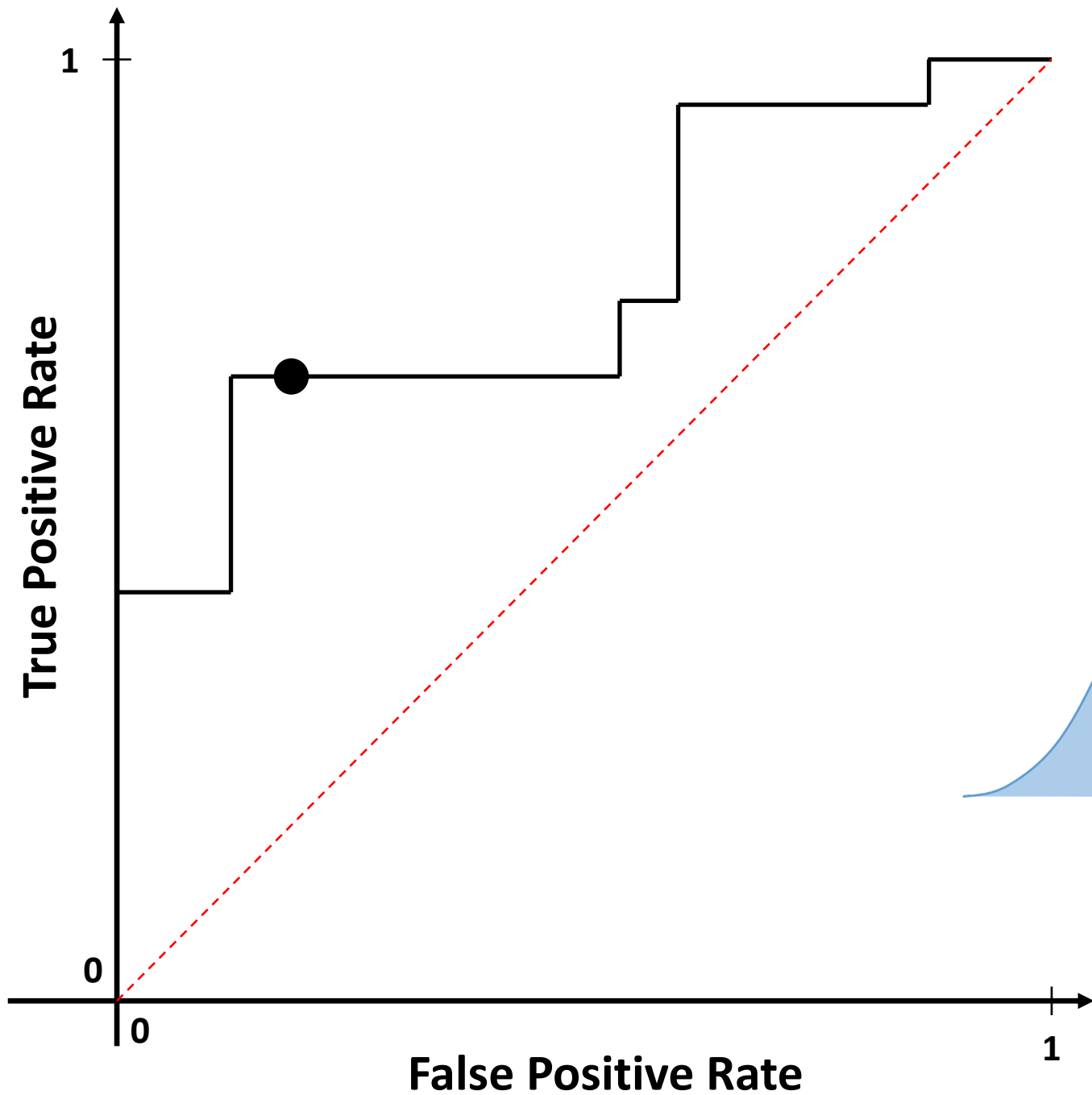
$$\frac{TP}{TP + FN}$$



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

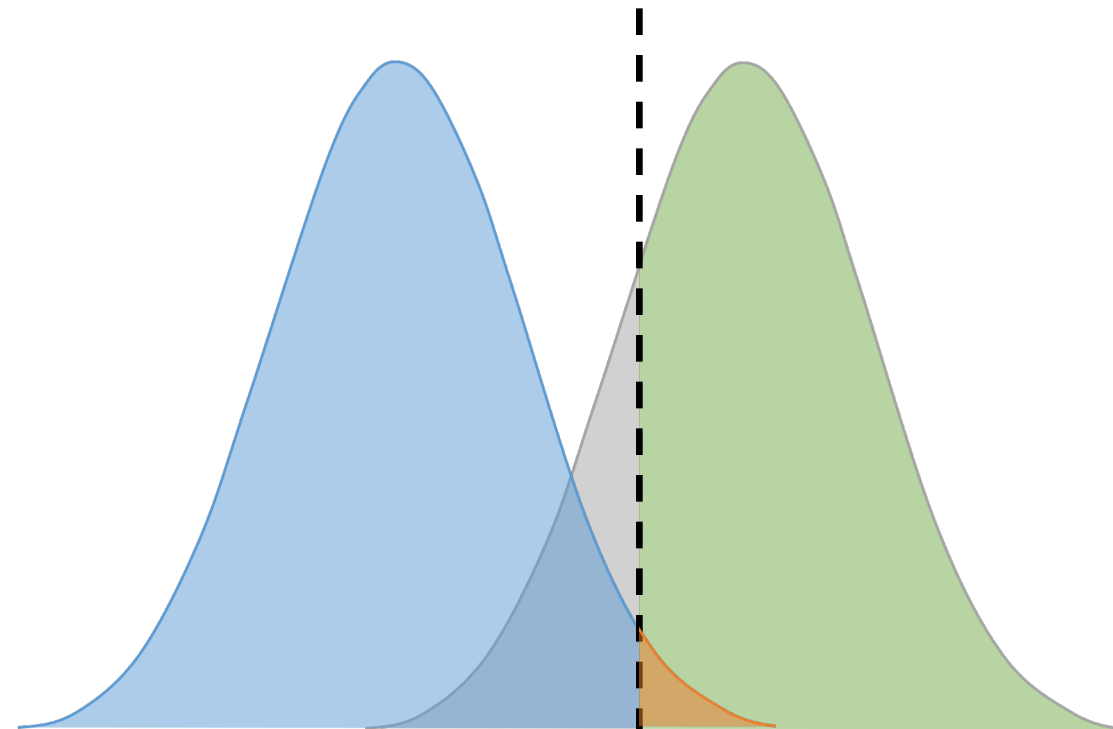
$$TPR = \frac{TP}{TP + FN}$$

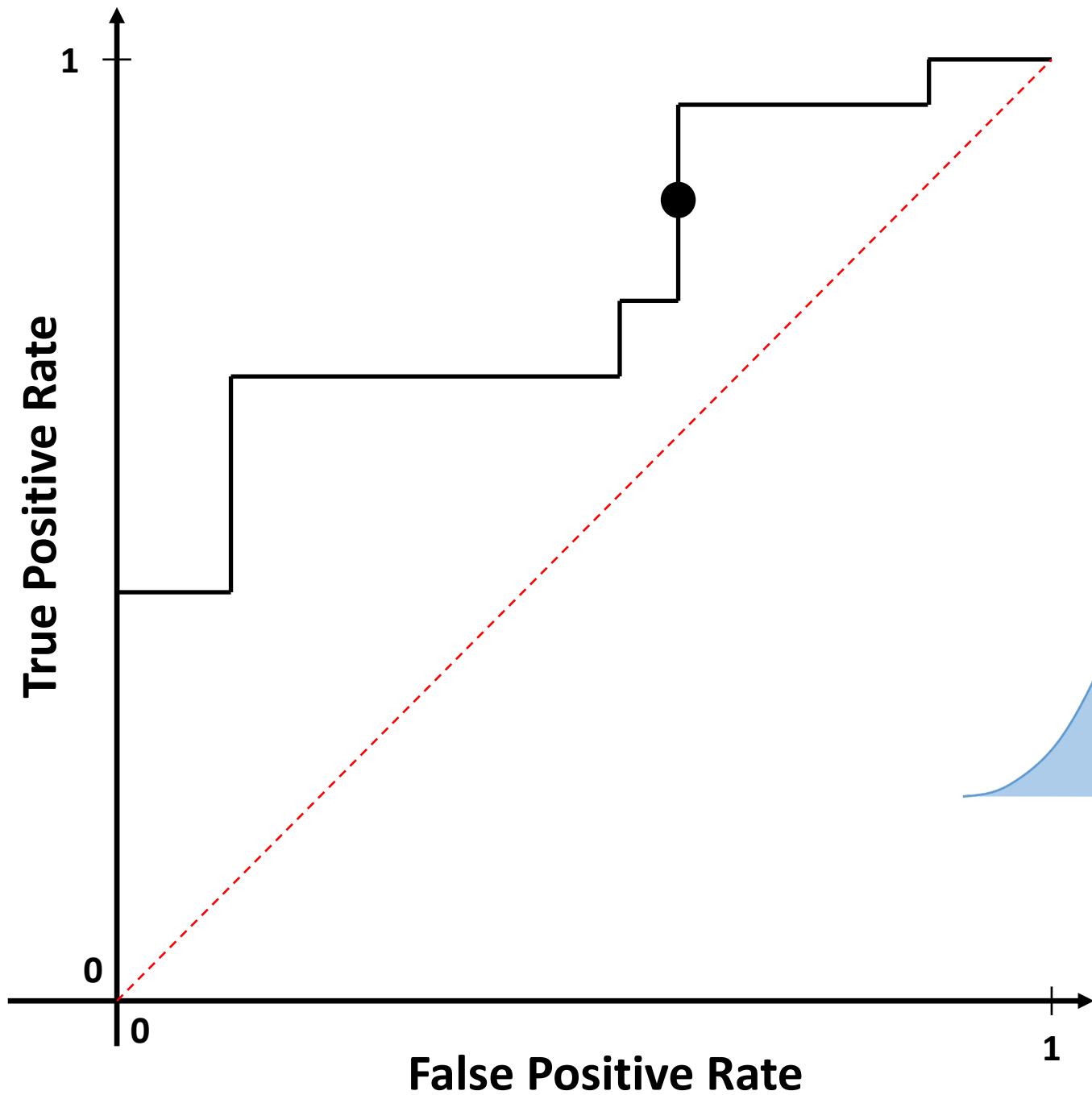




$$FPR = \frac{FP}{TN + FP}$$

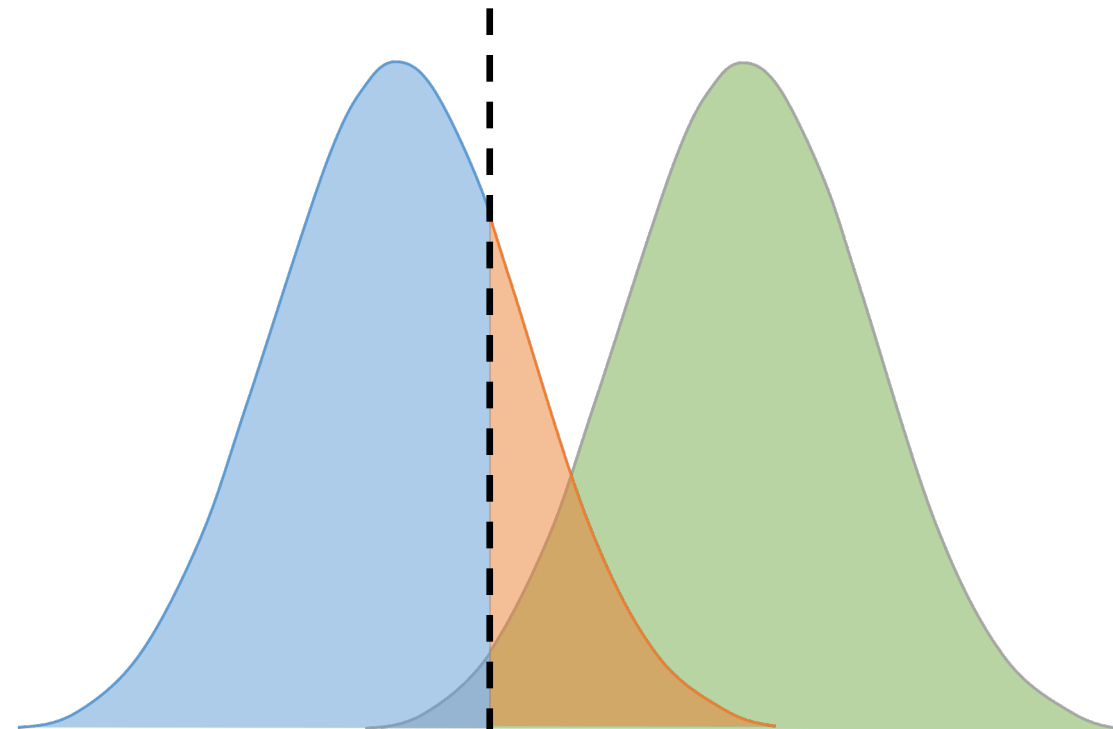
$$TPR = \frac{TP}{TP + FN}$$

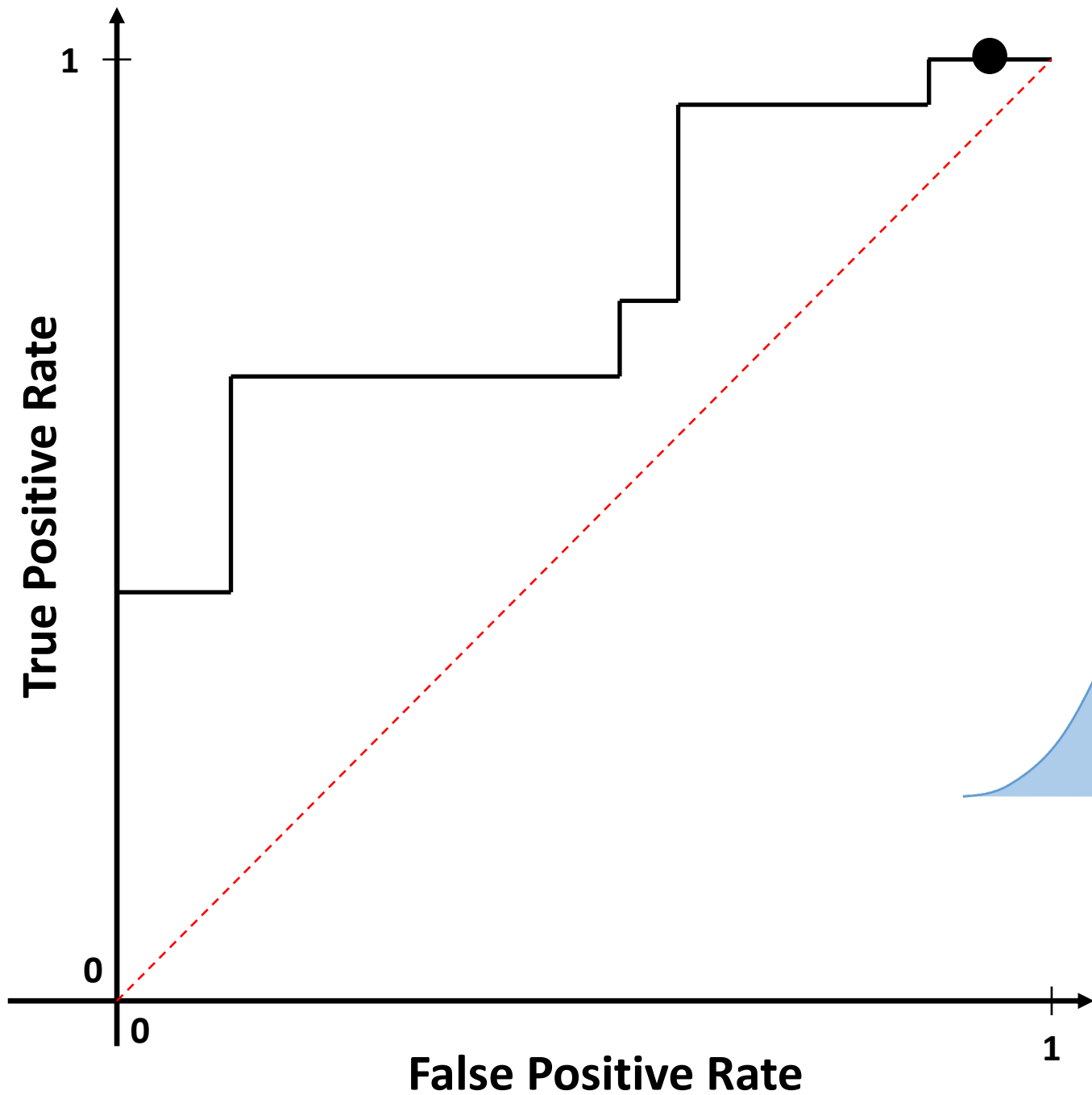




$$FPR = \frac{FP}{TN + FP}$$

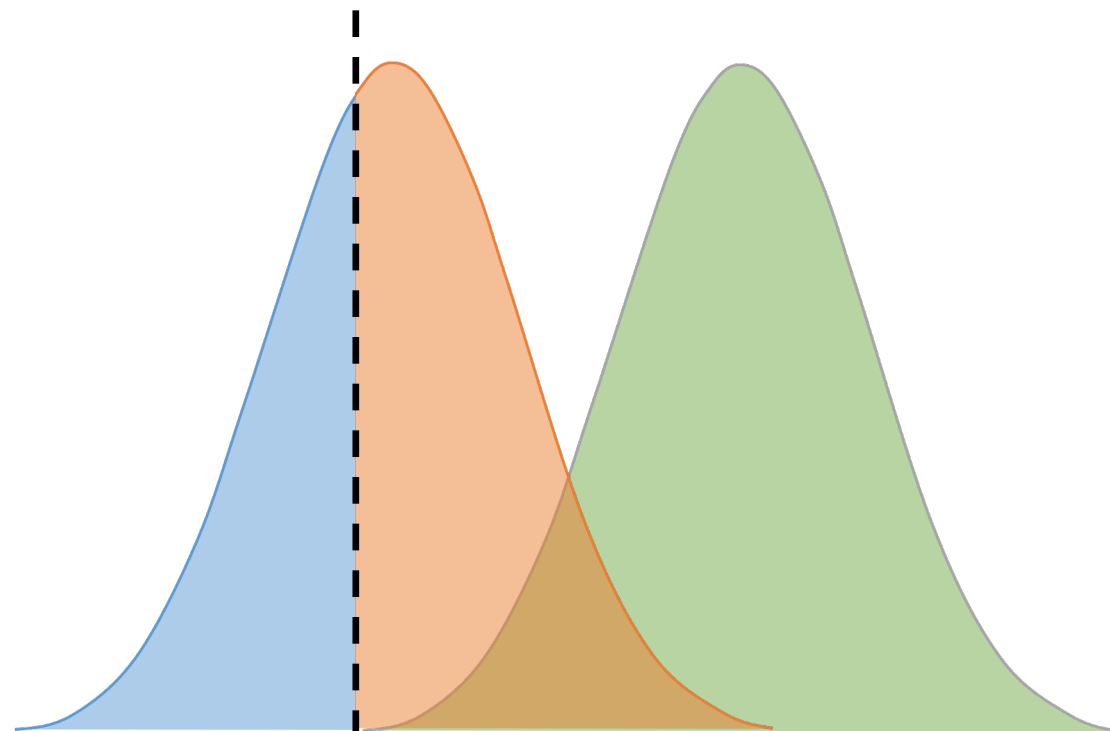
$$TPR = \frac{TP}{TP + FN}$$

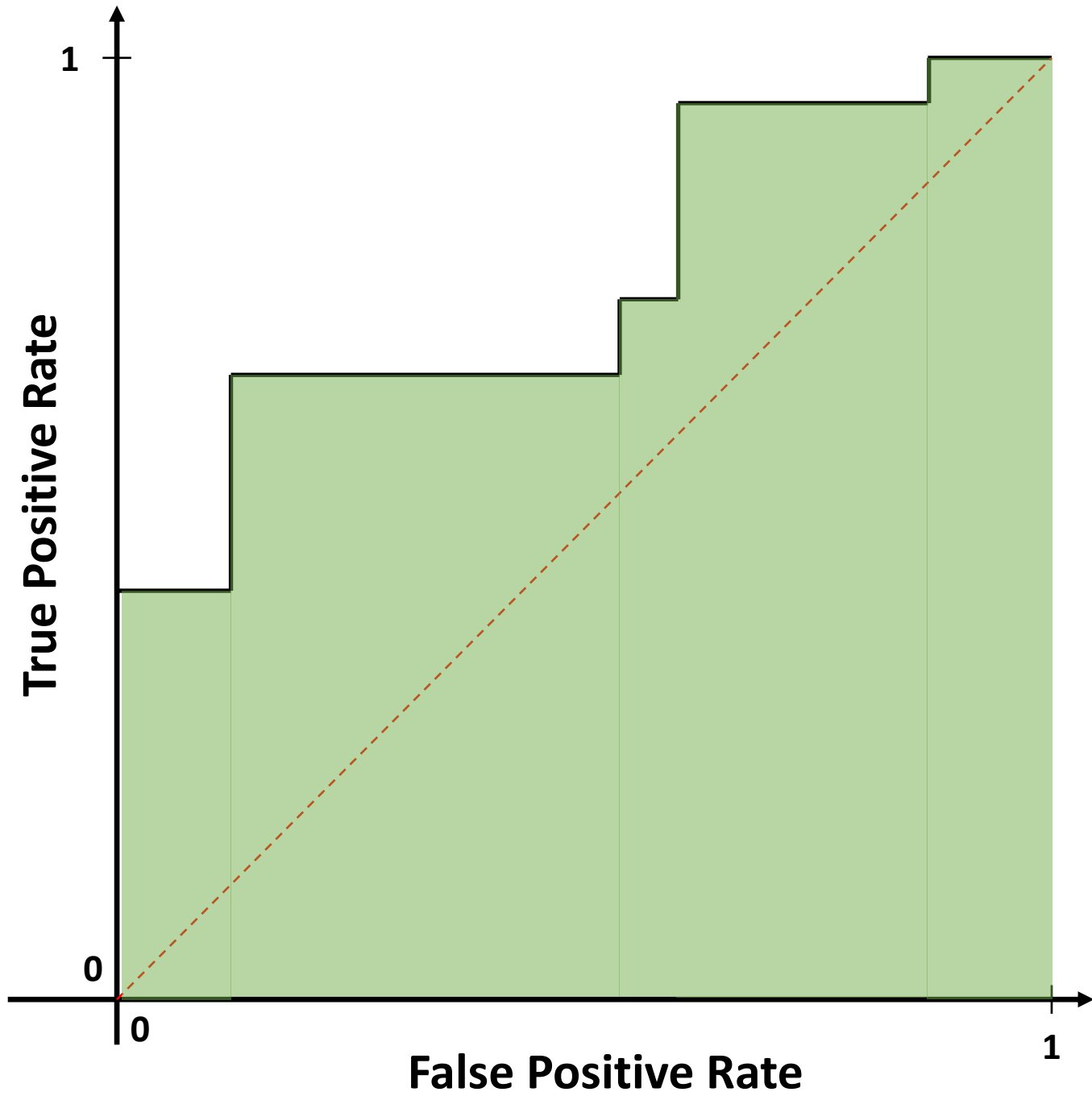


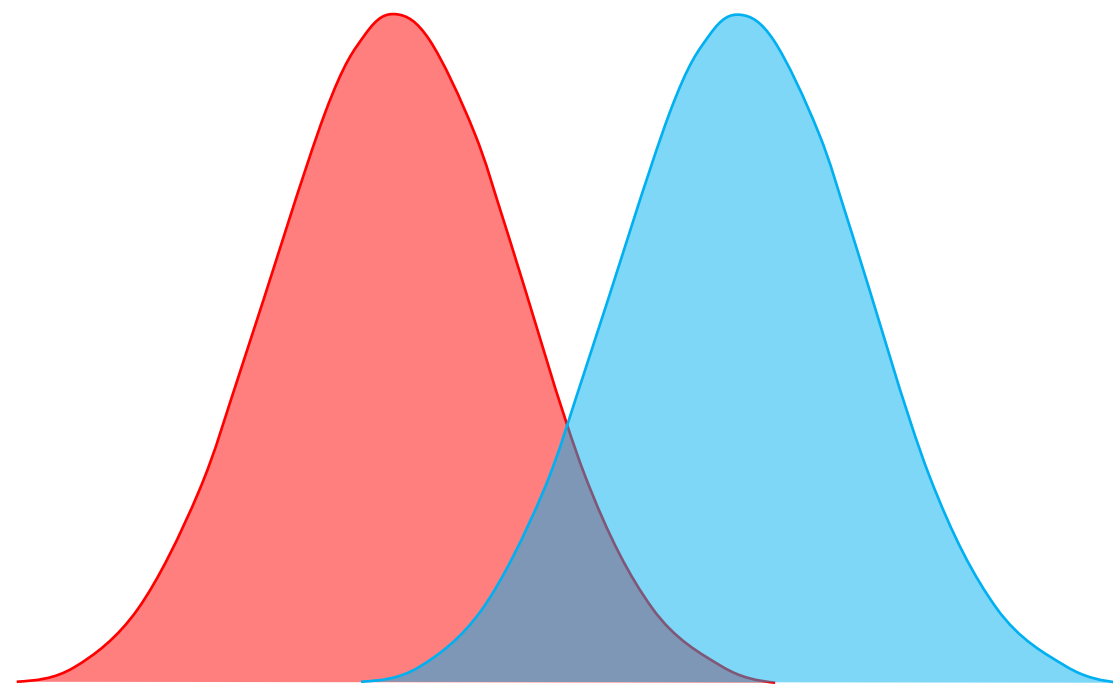
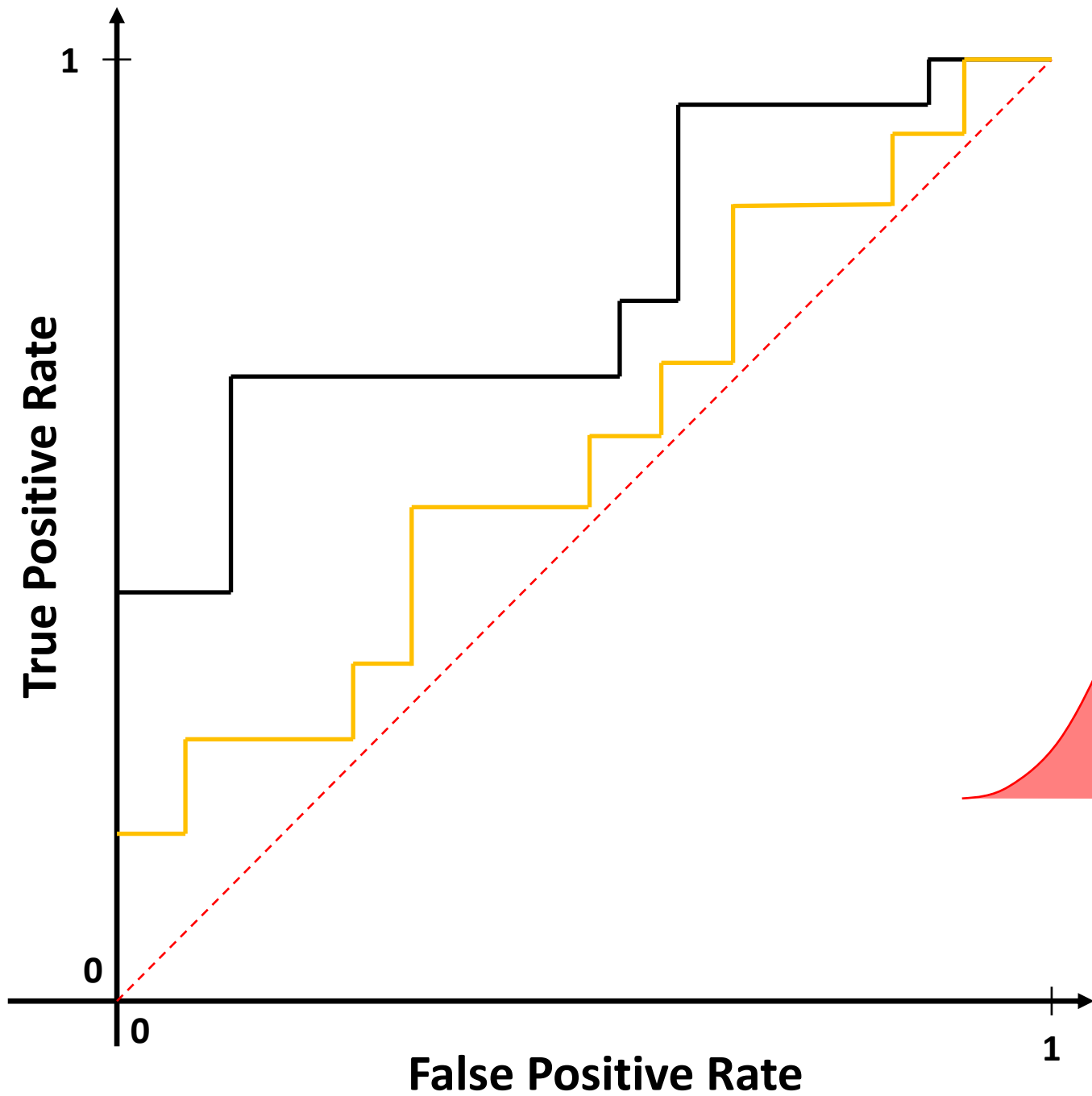


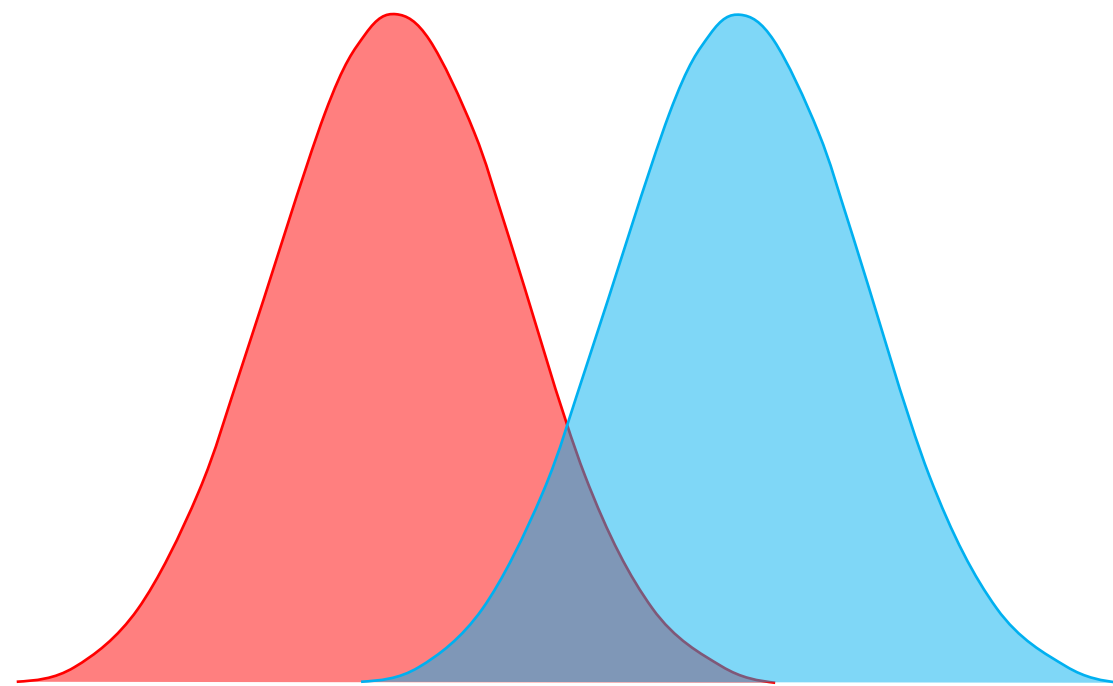
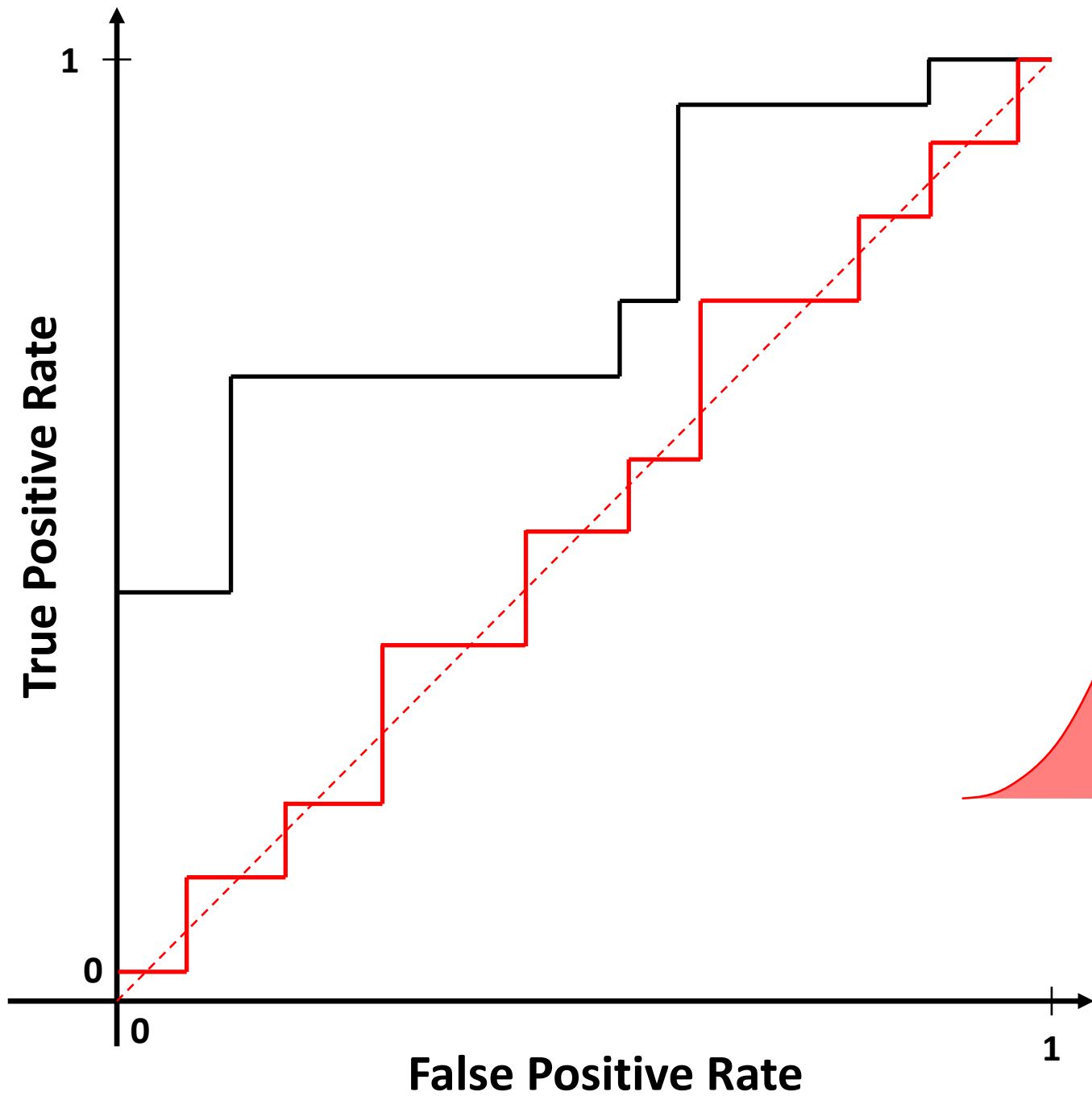
$$FPR = \frac{FP}{TN + FP}$$

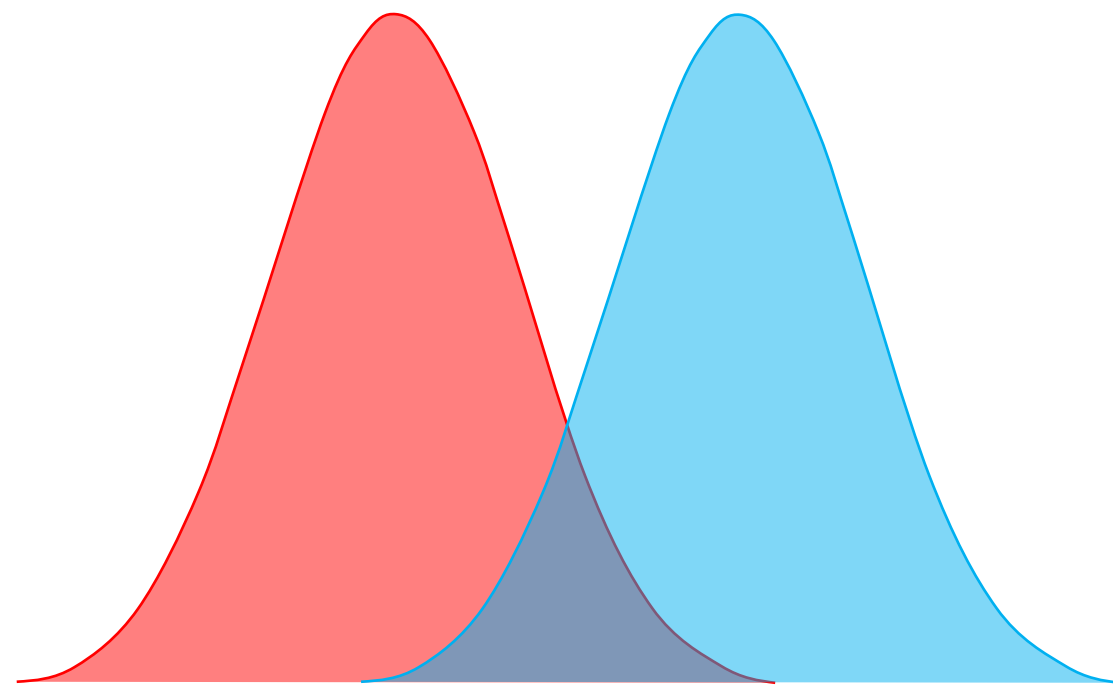
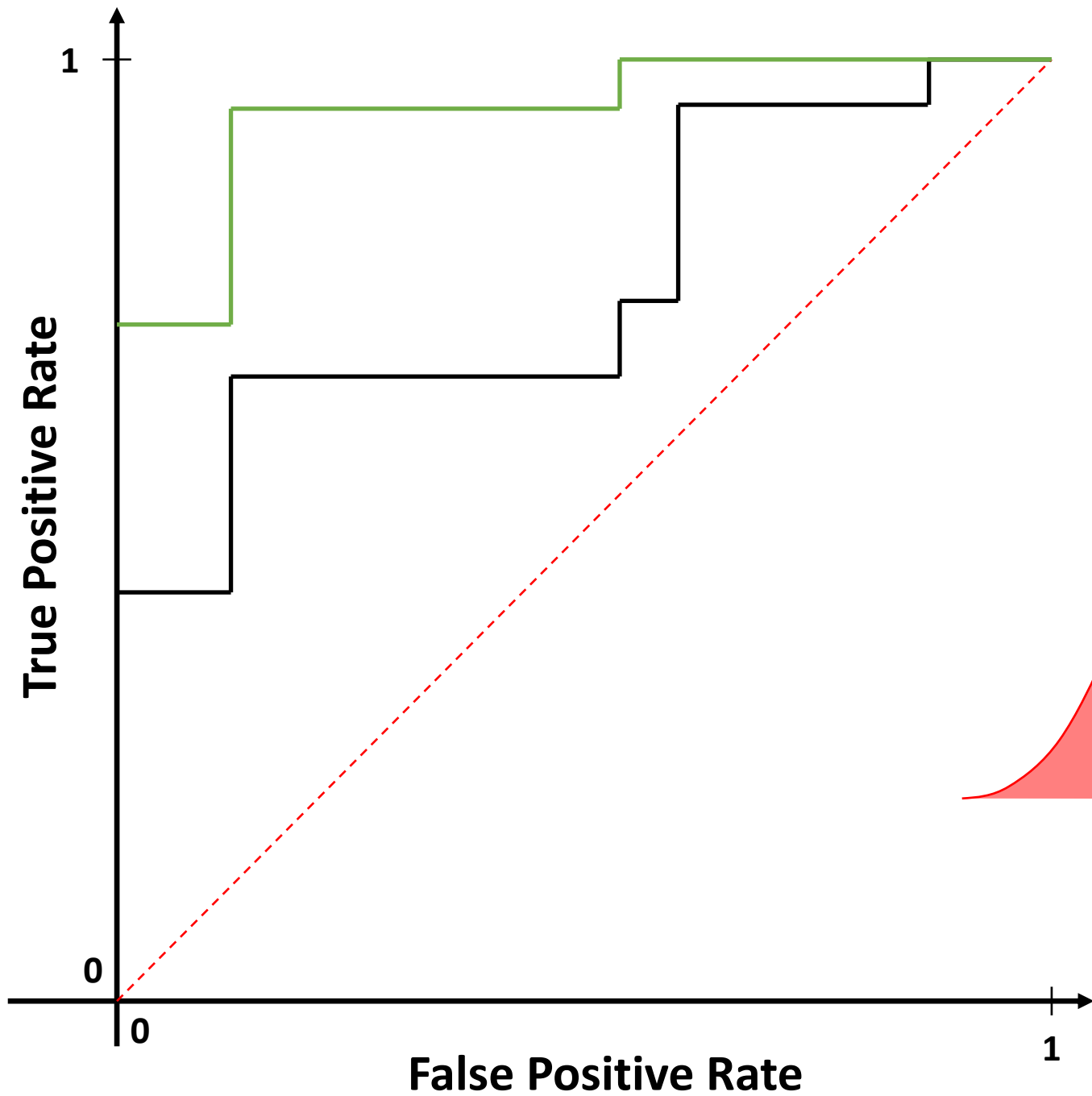
$$TPR = \frac{TP}{TP + FN}$$

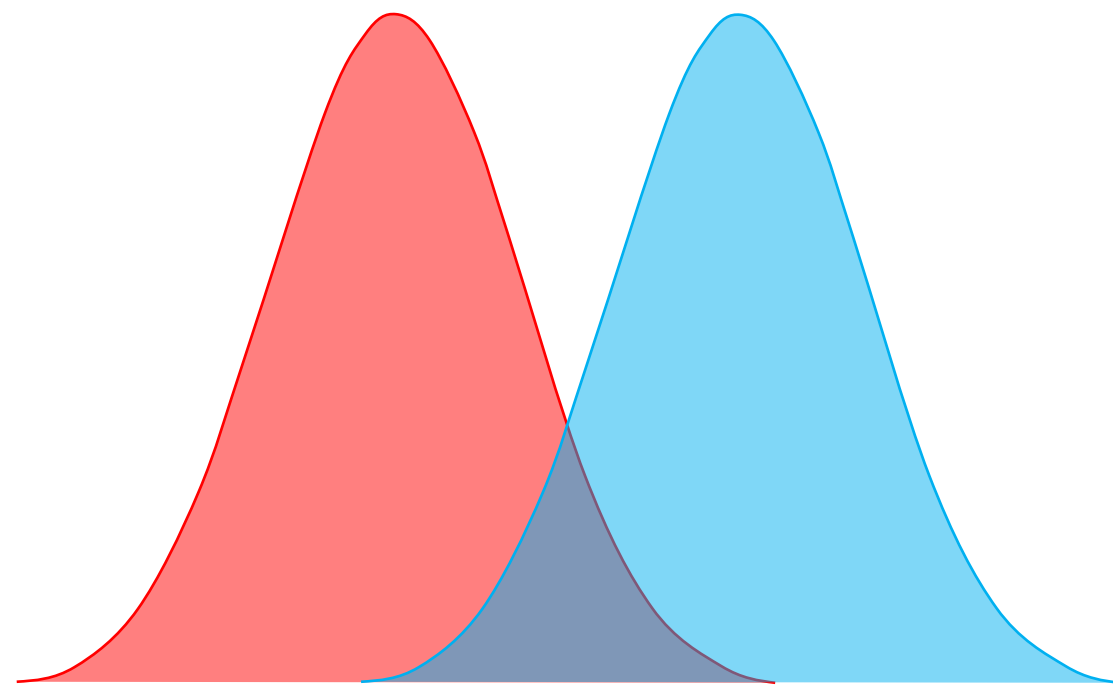
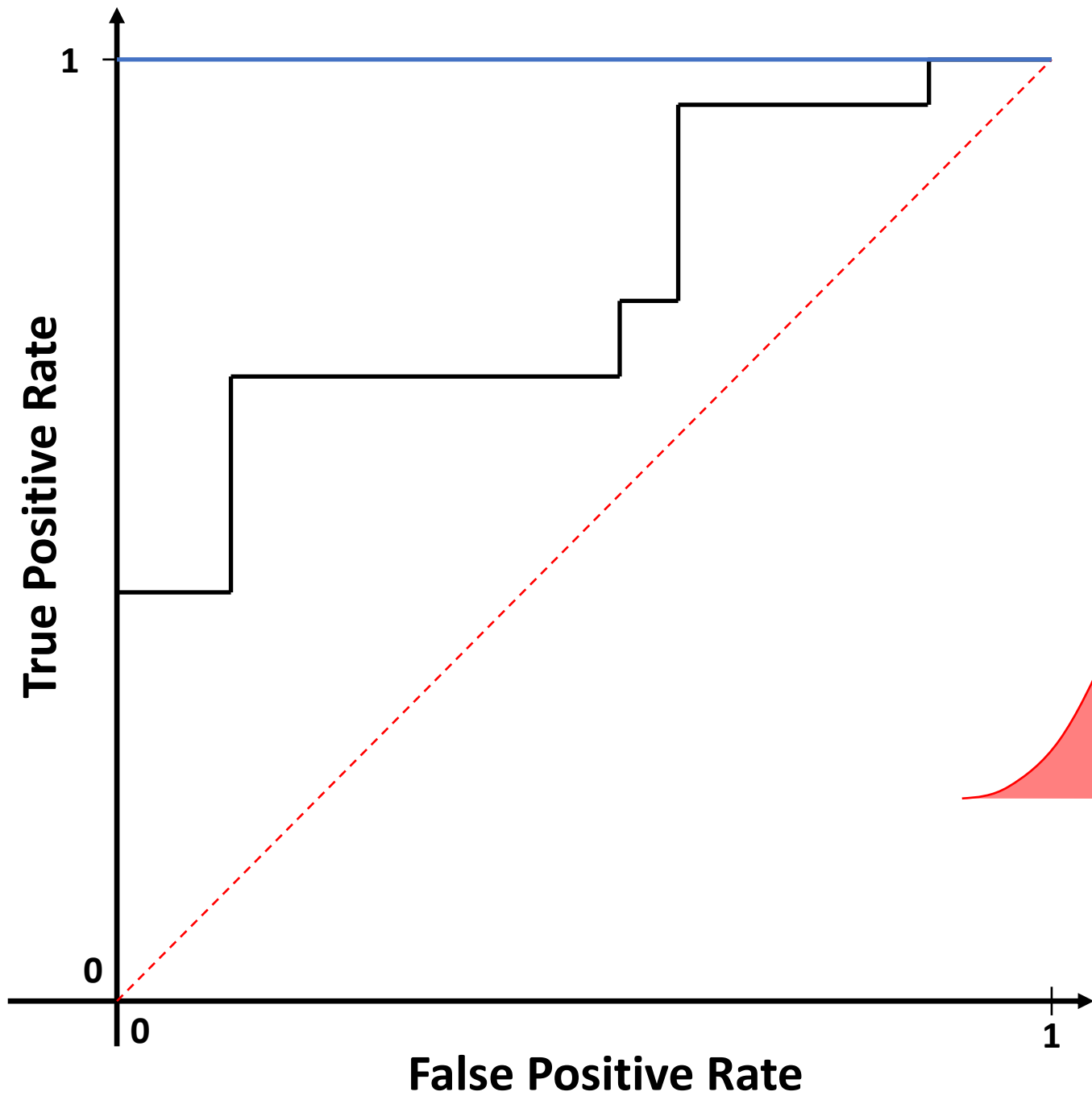








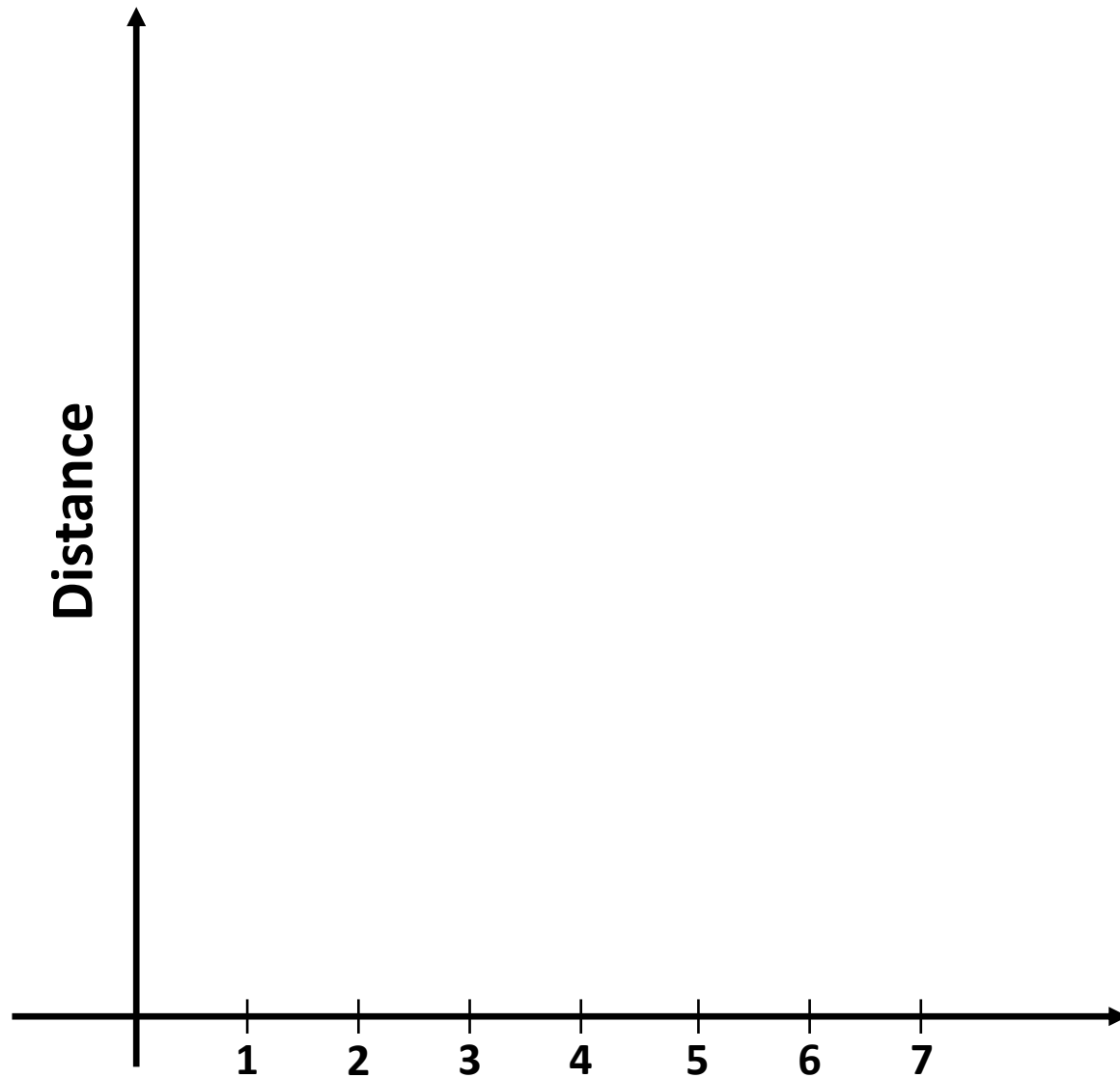
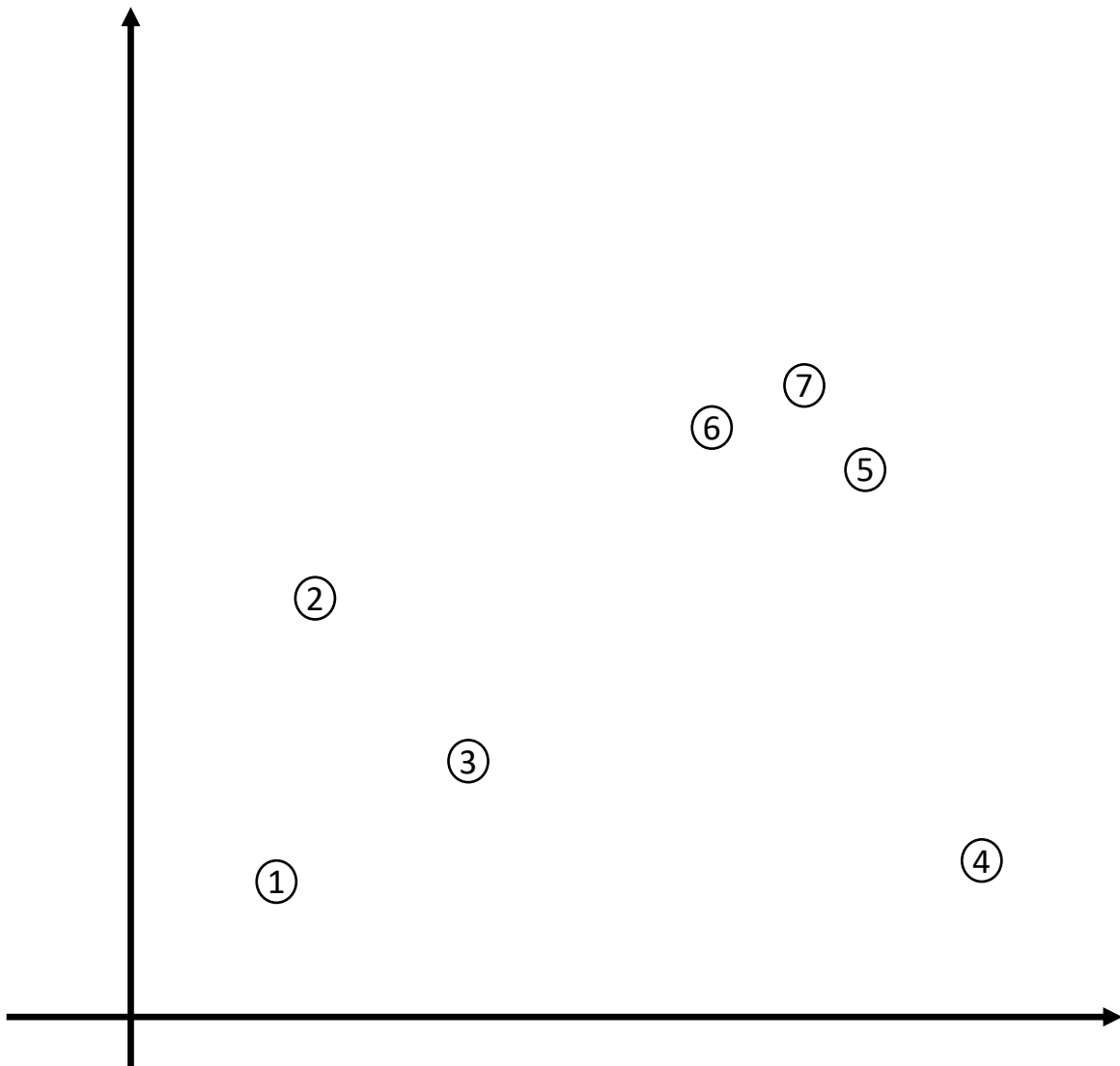


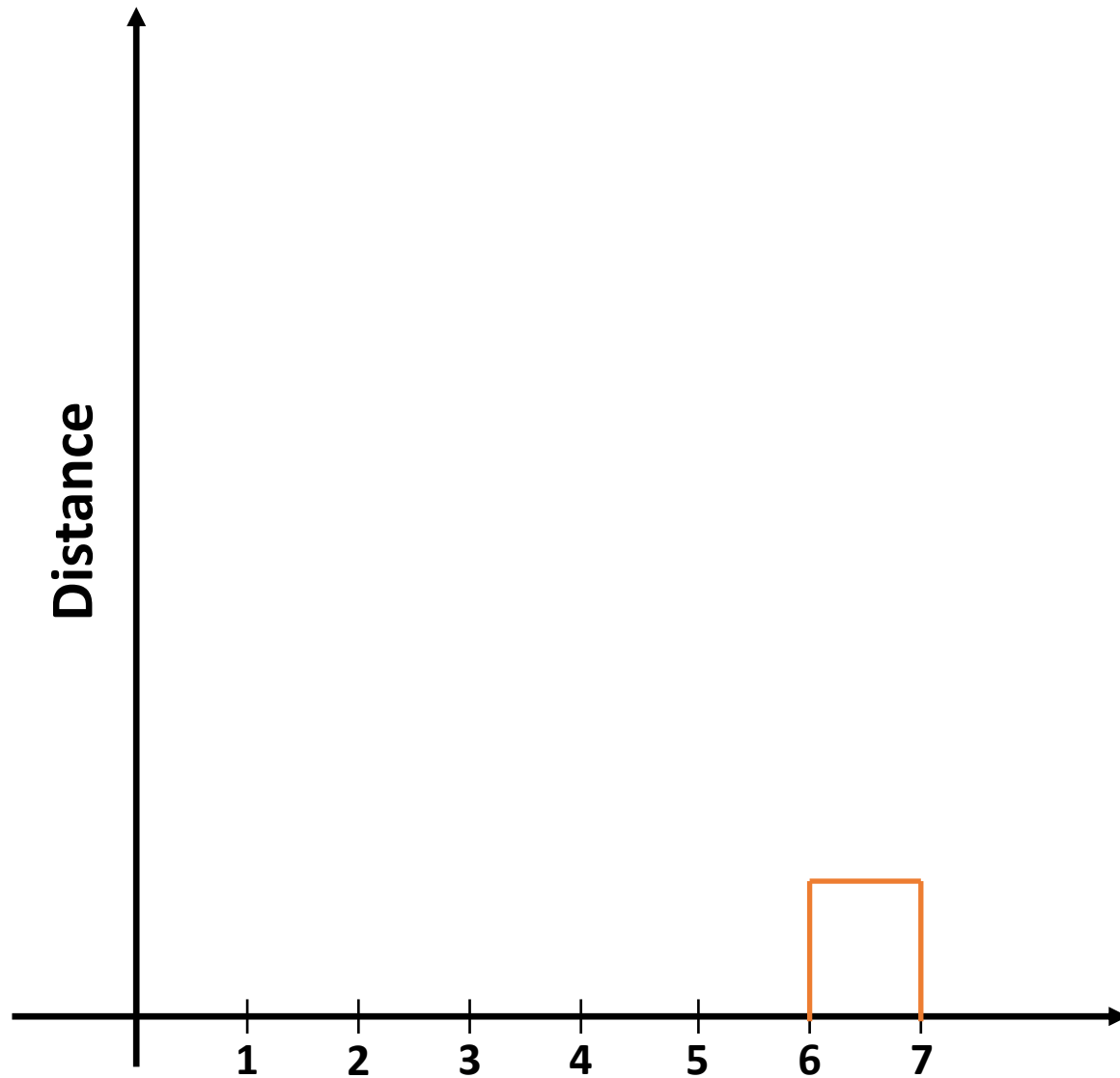
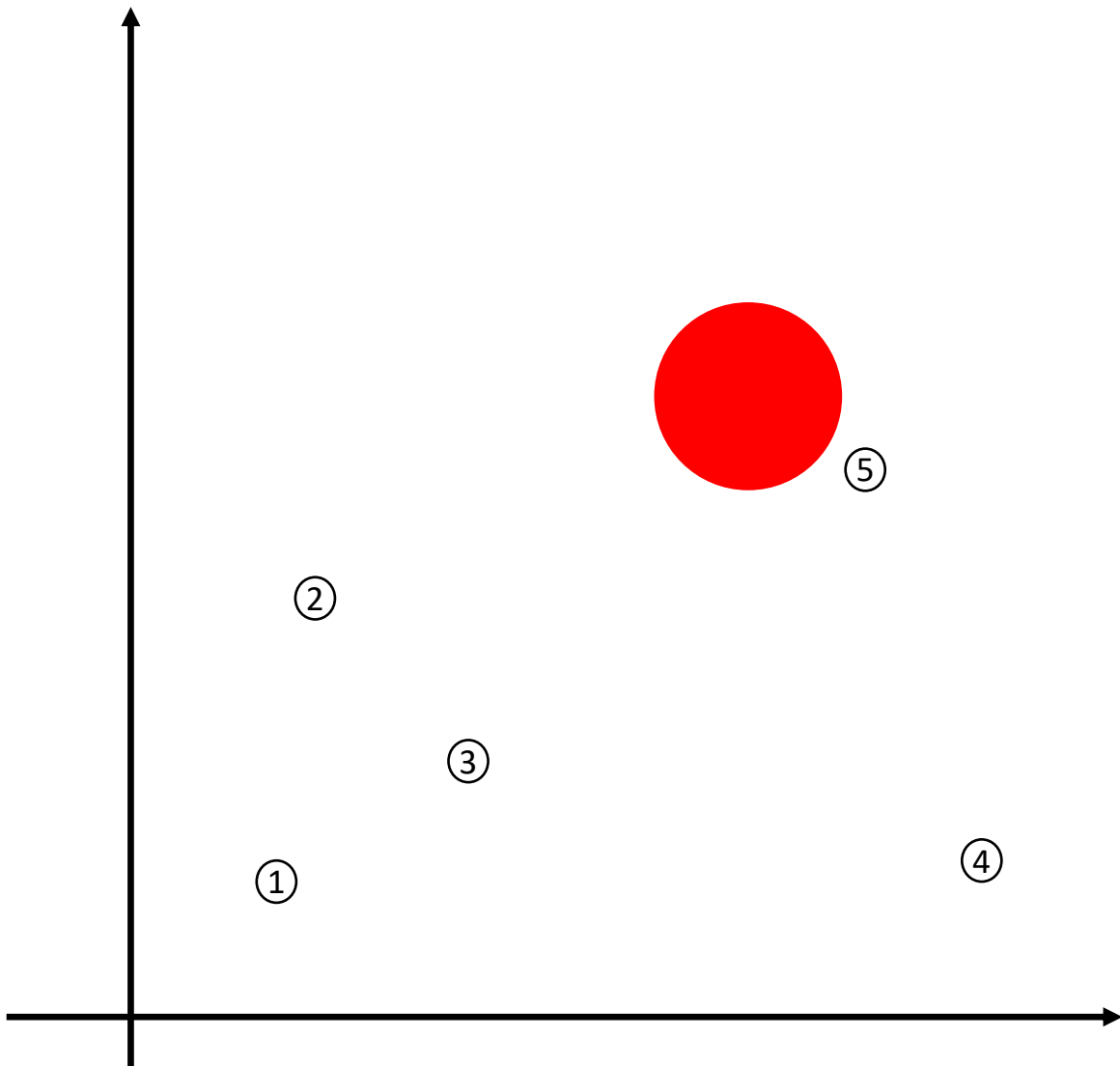


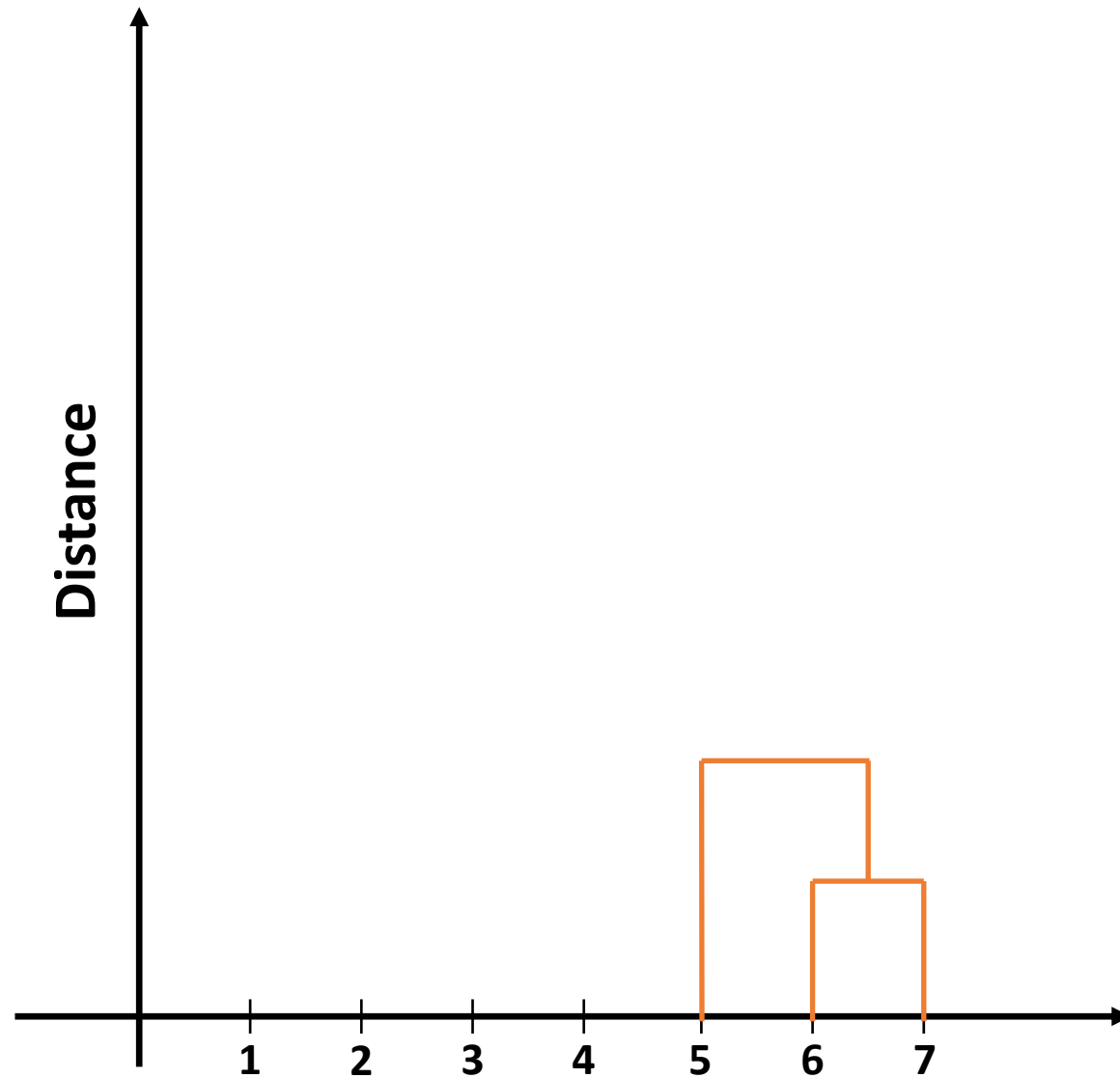
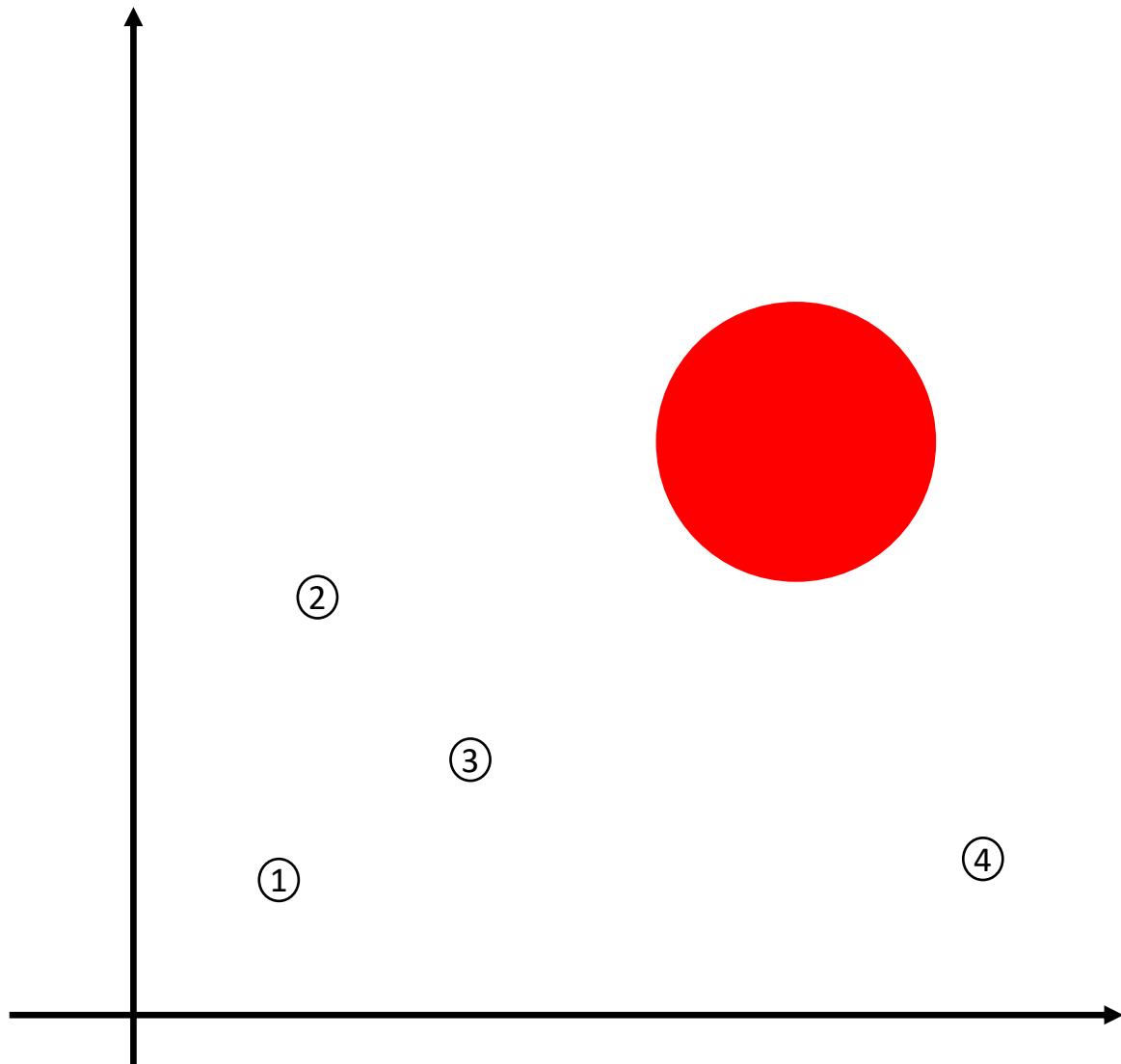
Summary

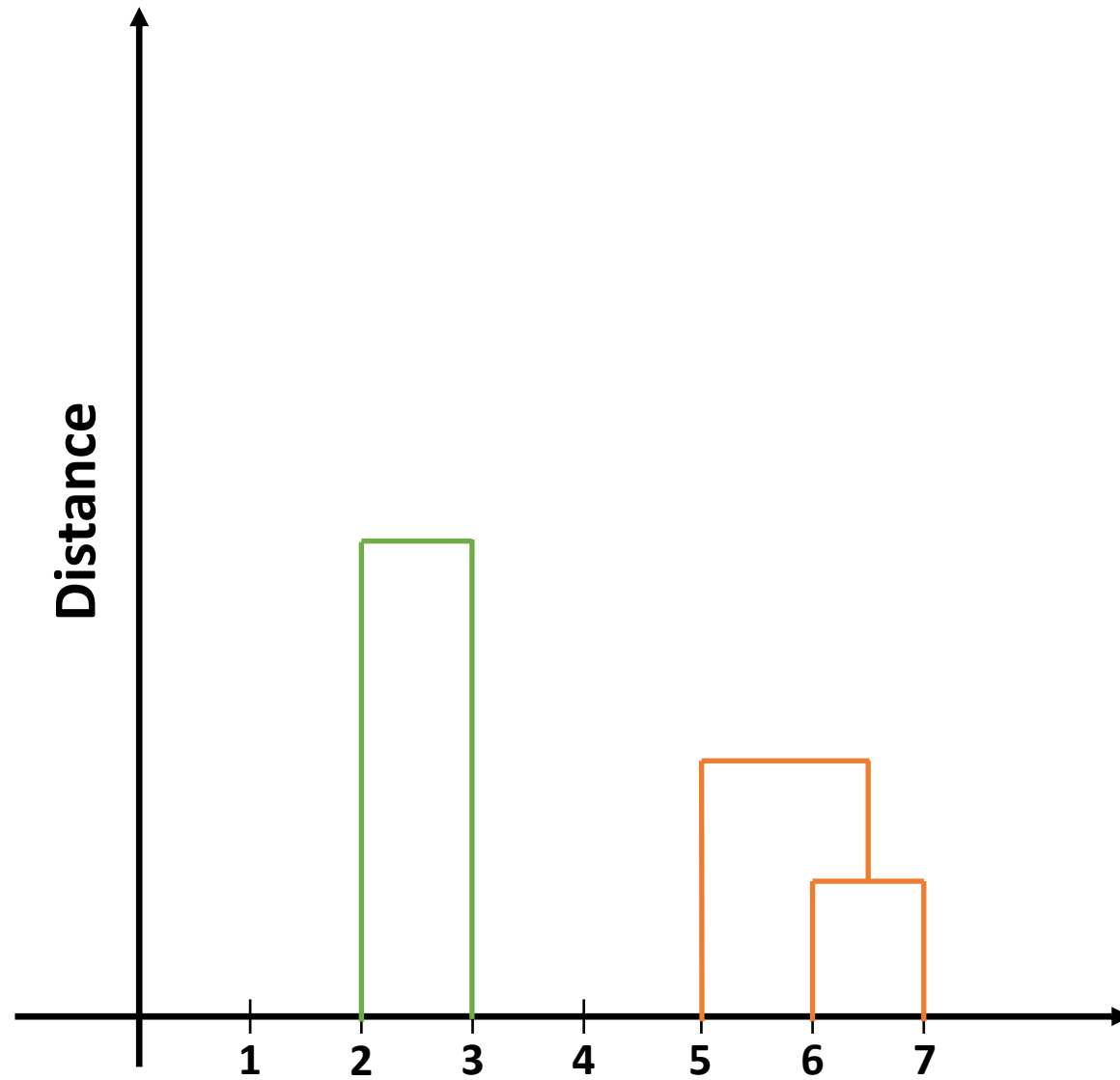
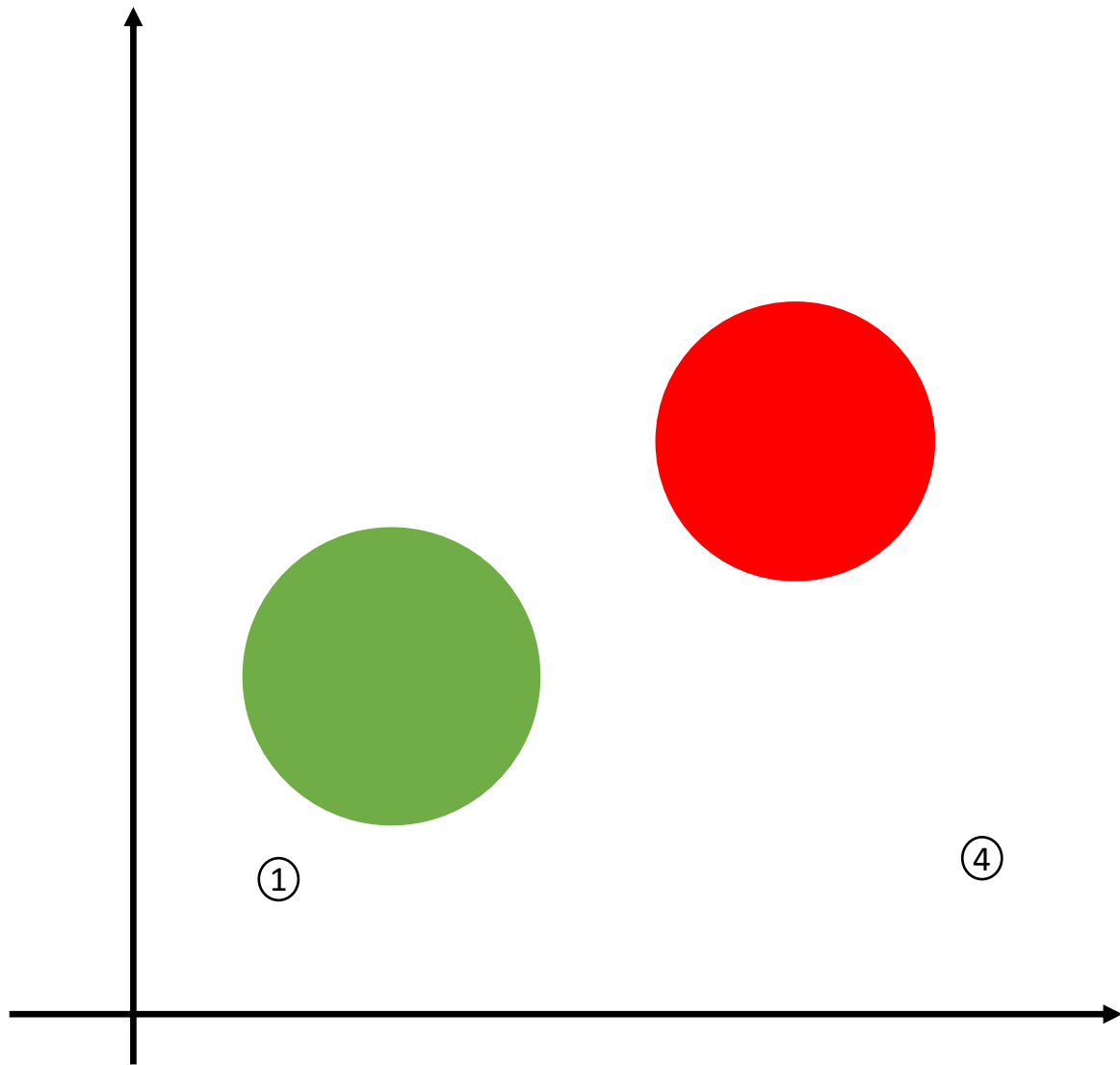
- Methods of performing cross-validation
- How to evaluate classifiers using classification accuracy and ROC
- What a good and bad classifier is with regards to AUC and ROC
- How this works with our data

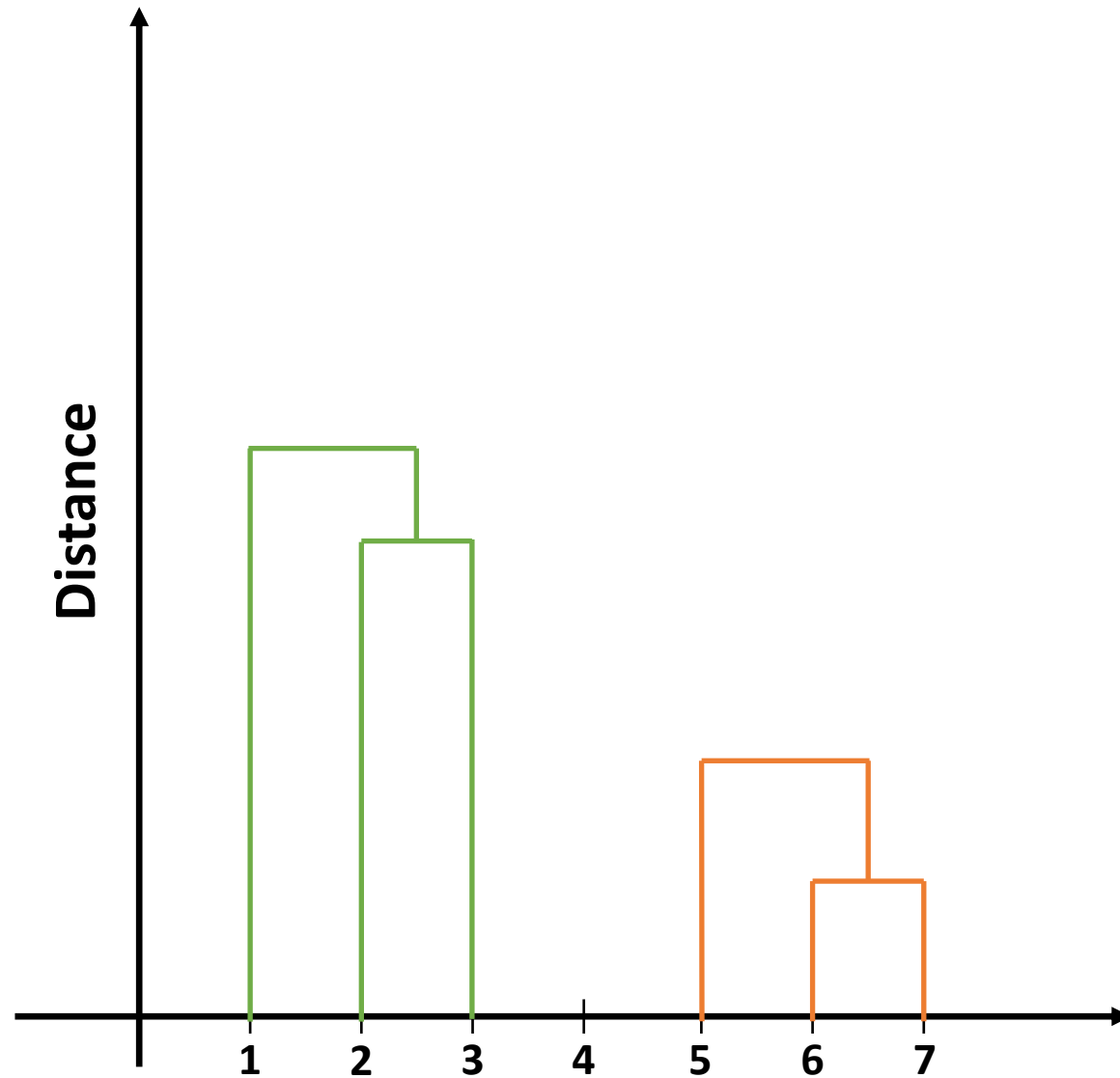
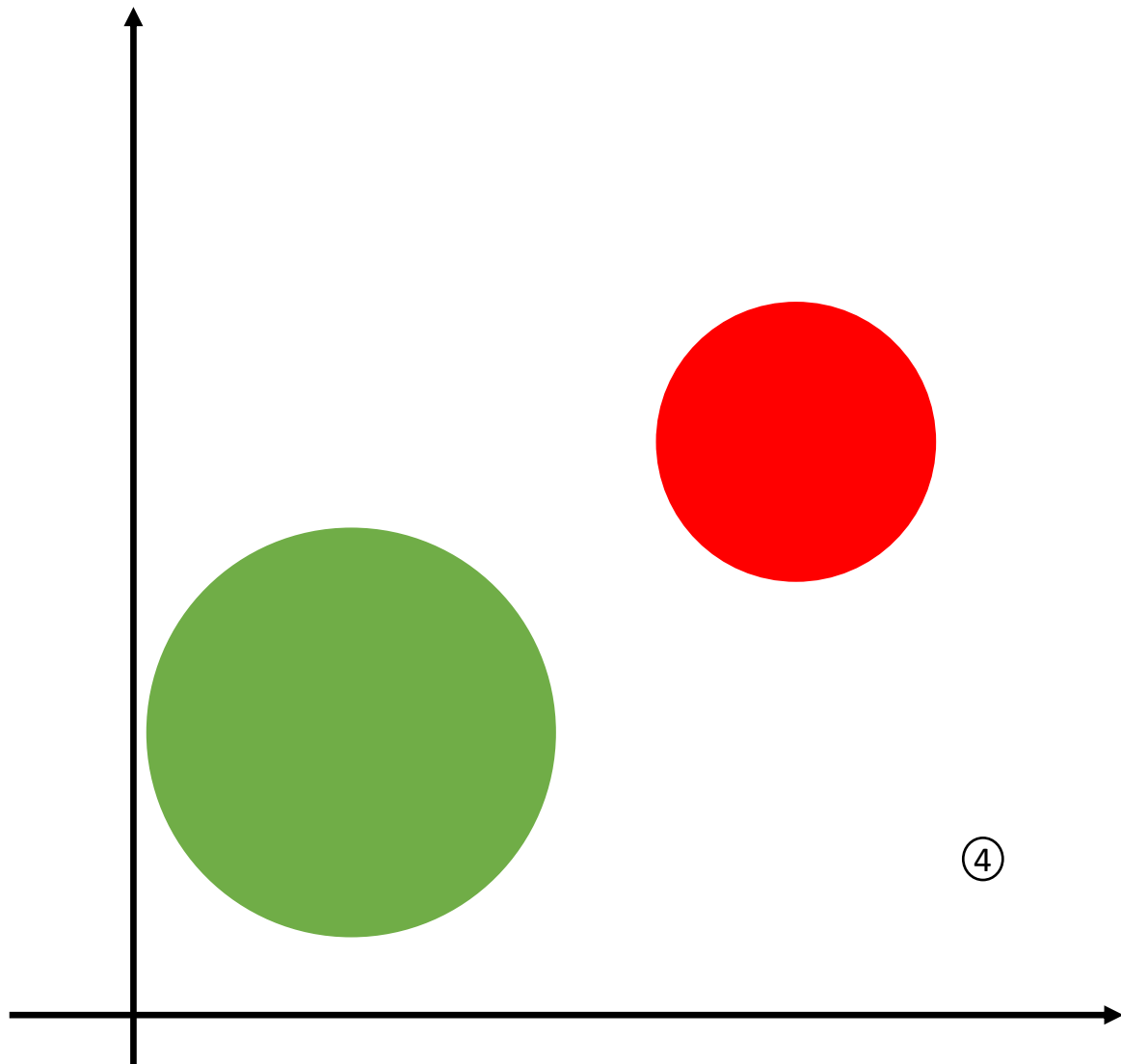
Hierarchical Clustering

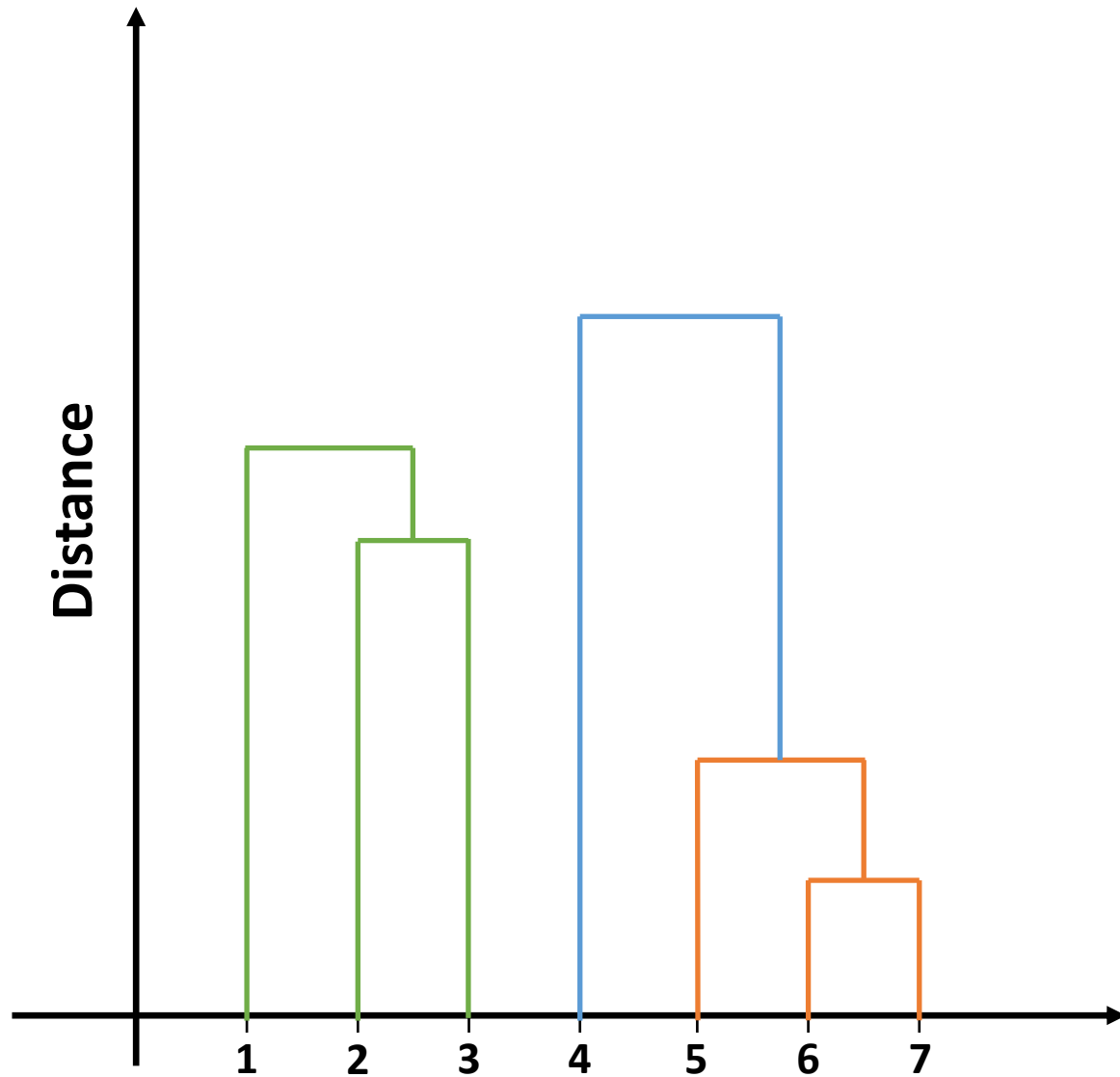
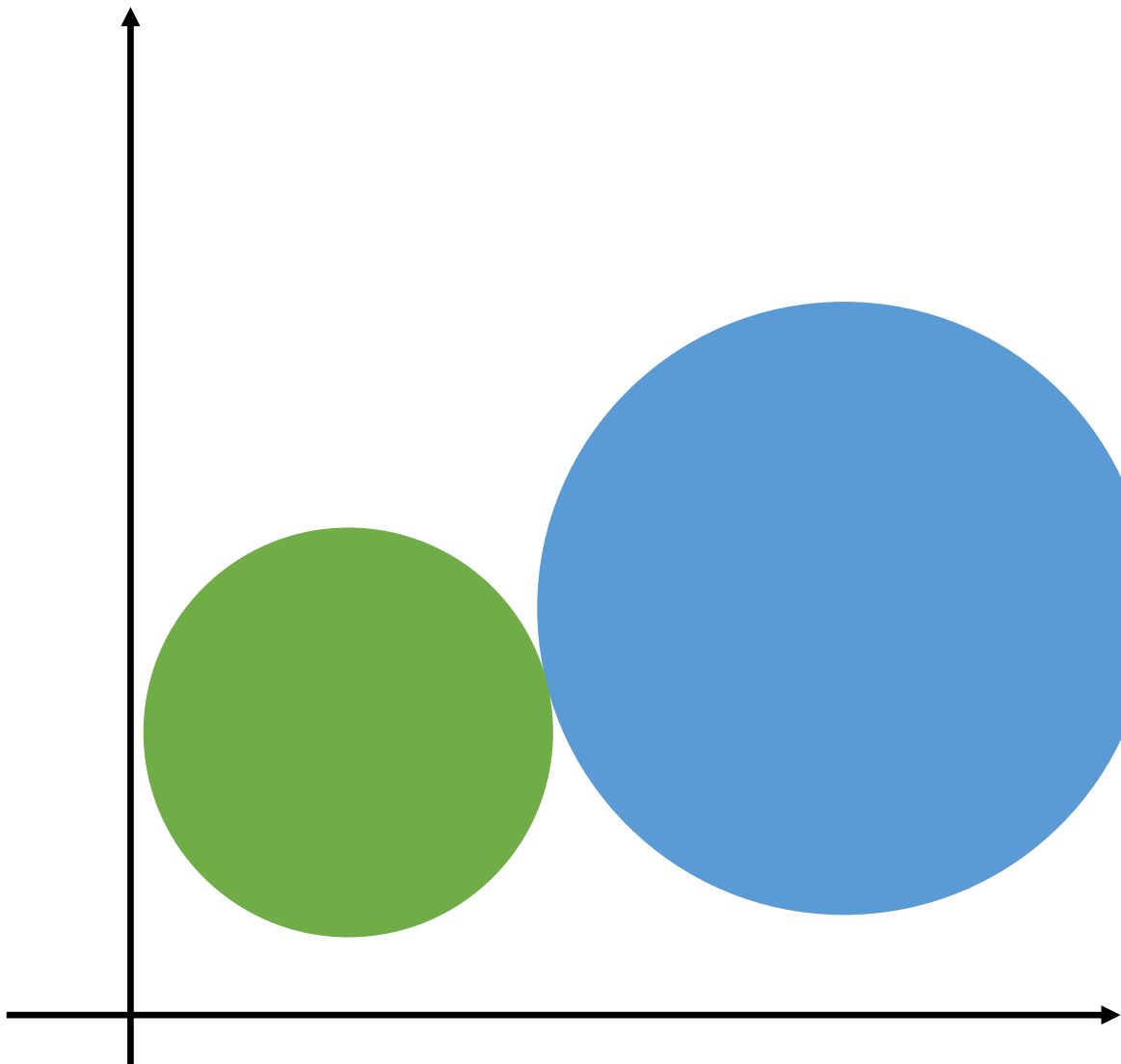


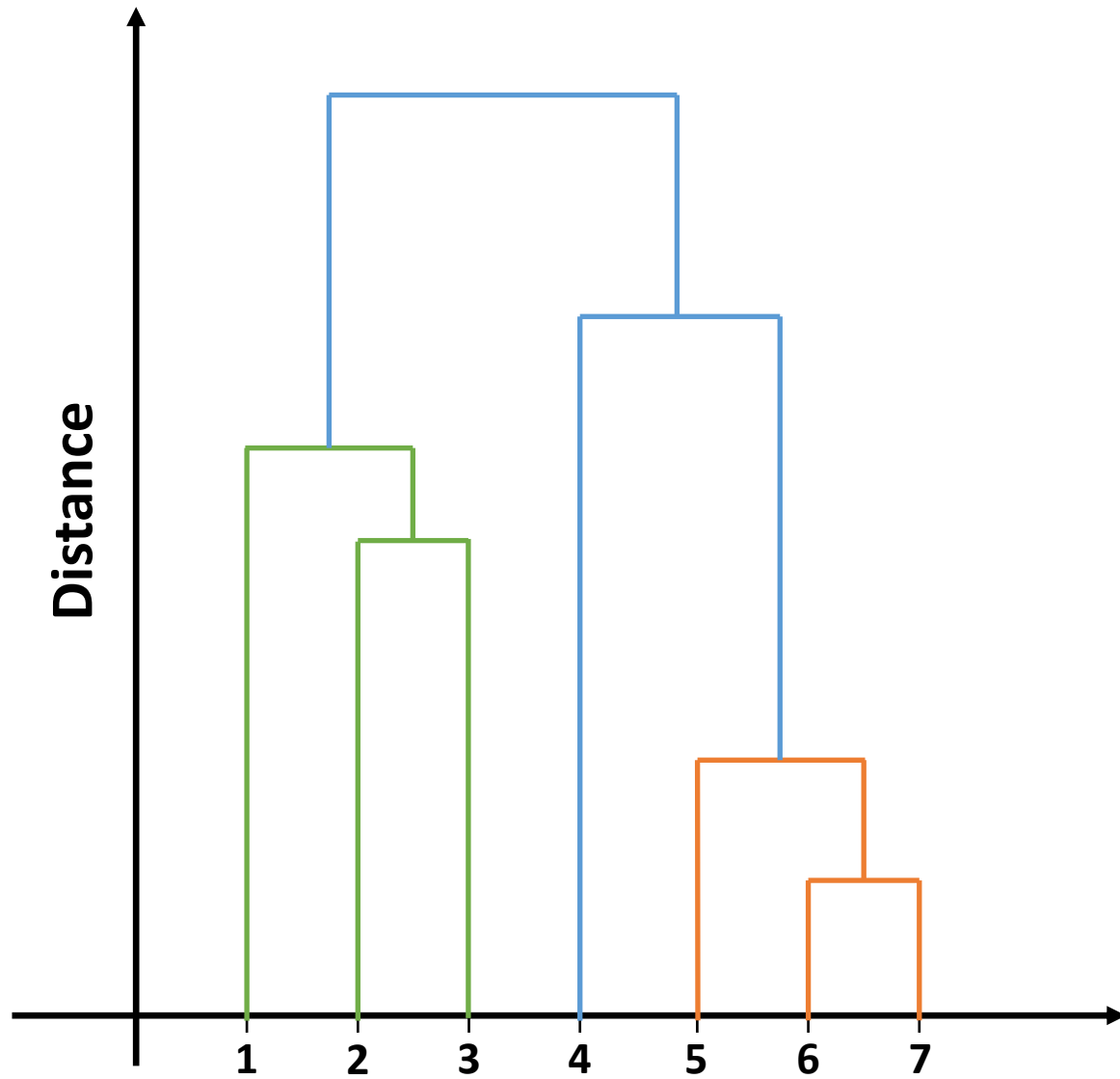
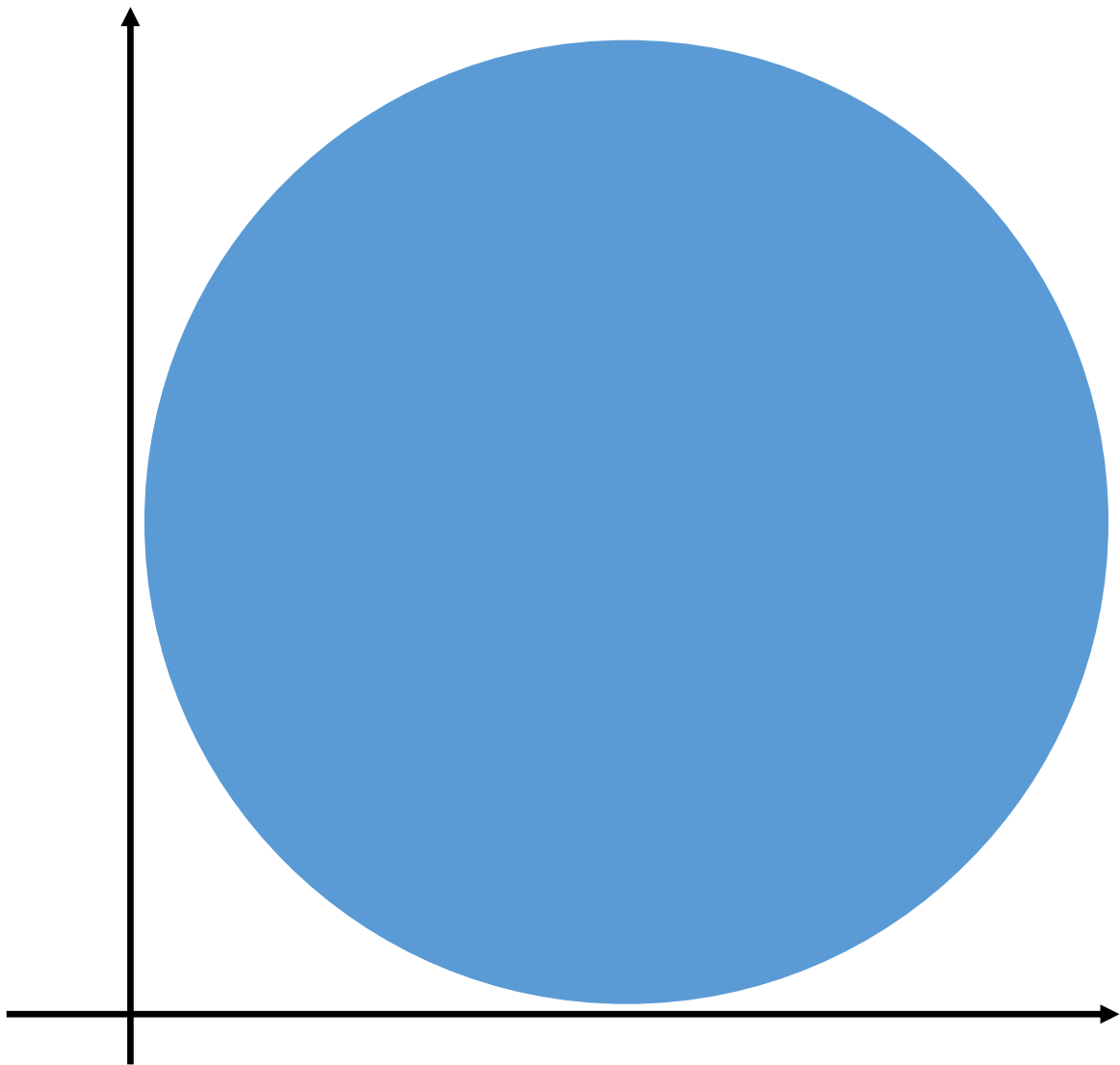


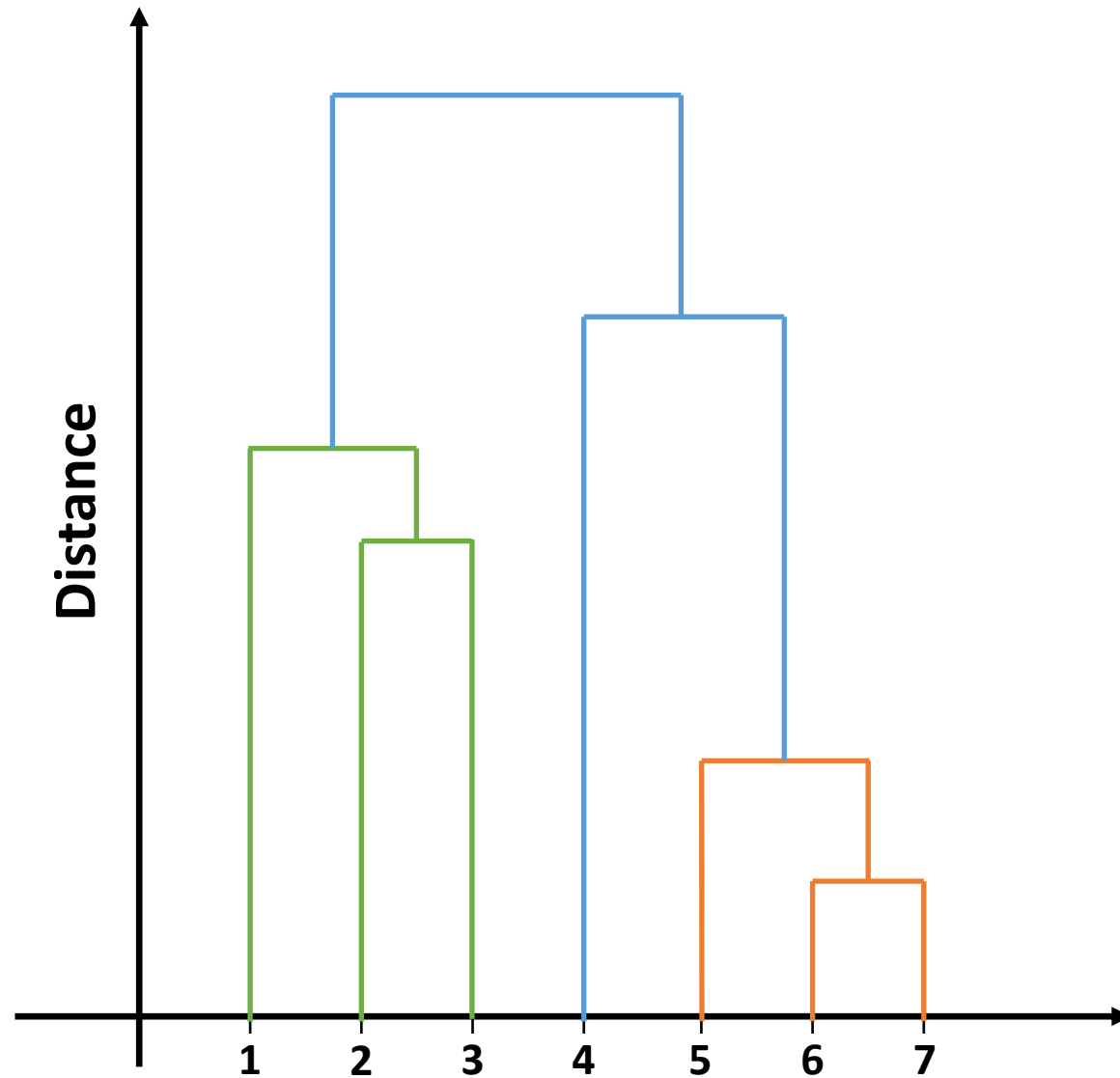
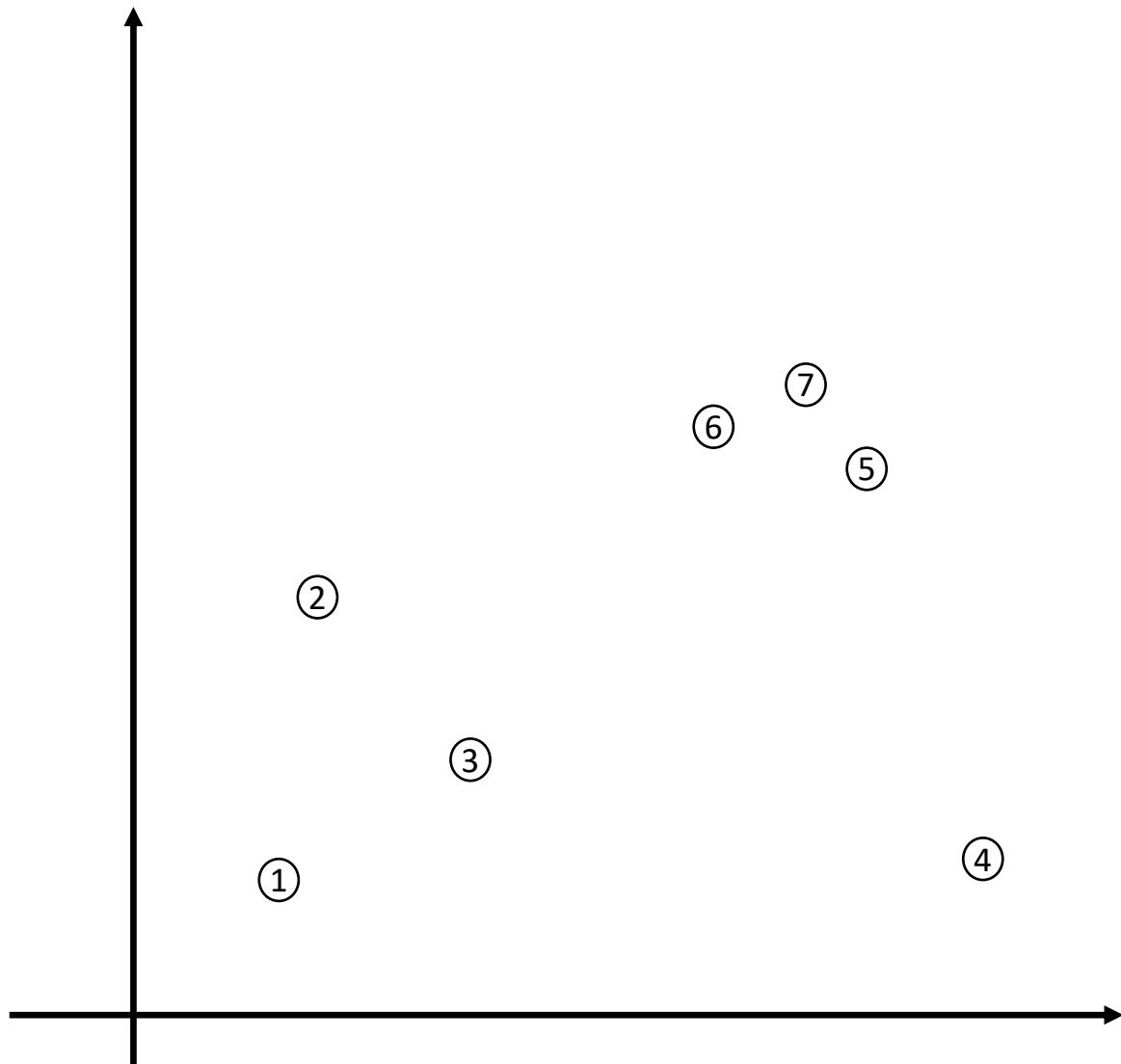












Summary

- Learned how to do hierarchical clustering

Acknowledgements

JHU

- Joshua Vogelstein
- Cencheng Shen

NIEHS

- Jesse Cushman
- Dalisa Kendricks
- Leslie Wilson
- Jariatu Stallone
- Sydney Fry
- DaNashia Thomas

Questions?

Input

Nearest Neighbors

Linear SVM

RBF SVM

Gaussian Process

Decision Tree

Random Forest

Neural Network

AdaBoost

Naïve Bayes

QDA

