

# Machine Learning Fundamentals

Zakarya Farou; Krisztian Buza

Department of Artificial Intelligence  
Eötvös Loránd University  
Budapest, Hungary

## Announcement

- Please send your presentation slides (about your selected topic) in PDF format to [buza@inf.elte.hu](mailto:buza@inf.elte.hu) ,  
Deadline: **18th November 2019, 8:00 am**
- Remember, you will have to present your topic on the  
**20th/27th November 2019**
- Check out the web page of the course ([www.biointelligence.hu/iml](http://www.biointelligence.hu/iml))  
**for the schedule of the presentations**

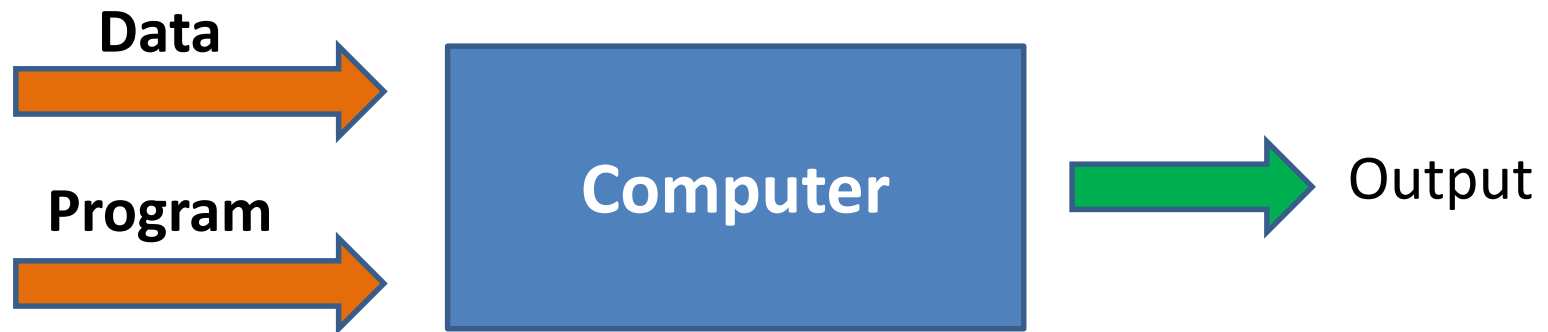
# What is Machine Learning?

- Machine learning is a form of AI that enables a system to learn from data rather than through explicit programming

# What is Machine Learning?

- Machine learning is a form of AI that enables a system to learn from data rather than through explicit programming

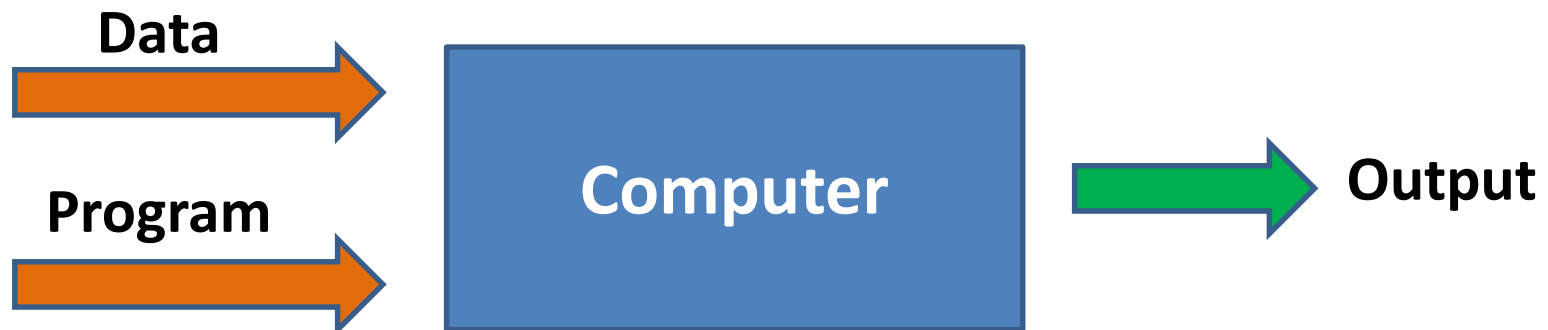
## Traditional Programming



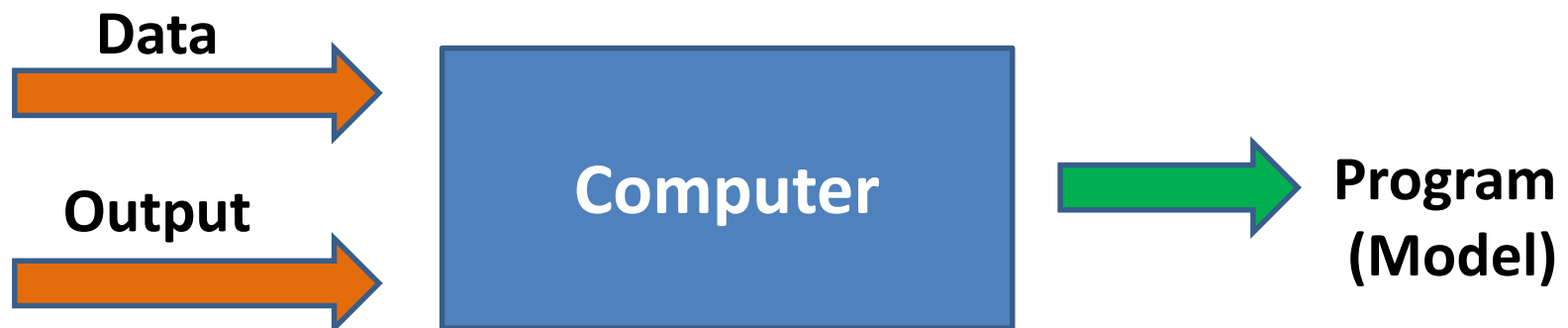
# What is Machine Learning?

- Machine learning is a form of AI that enables a system to learn from data rather than through explicit programming

## Traditional programming



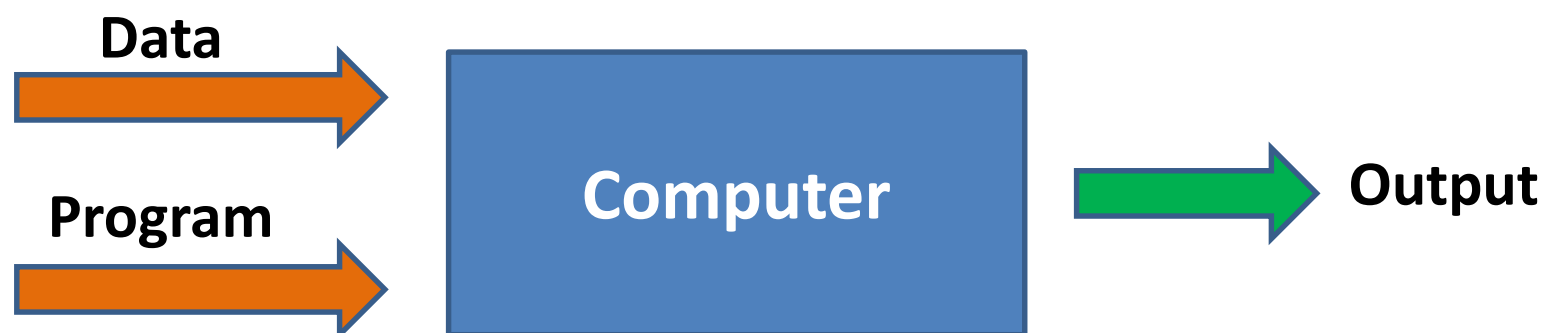
## Machine Learning (Supervised)



# What is Machine Learning?

- Machine learning is a form of AI that enables a system to learn from data rather than through explicit programming

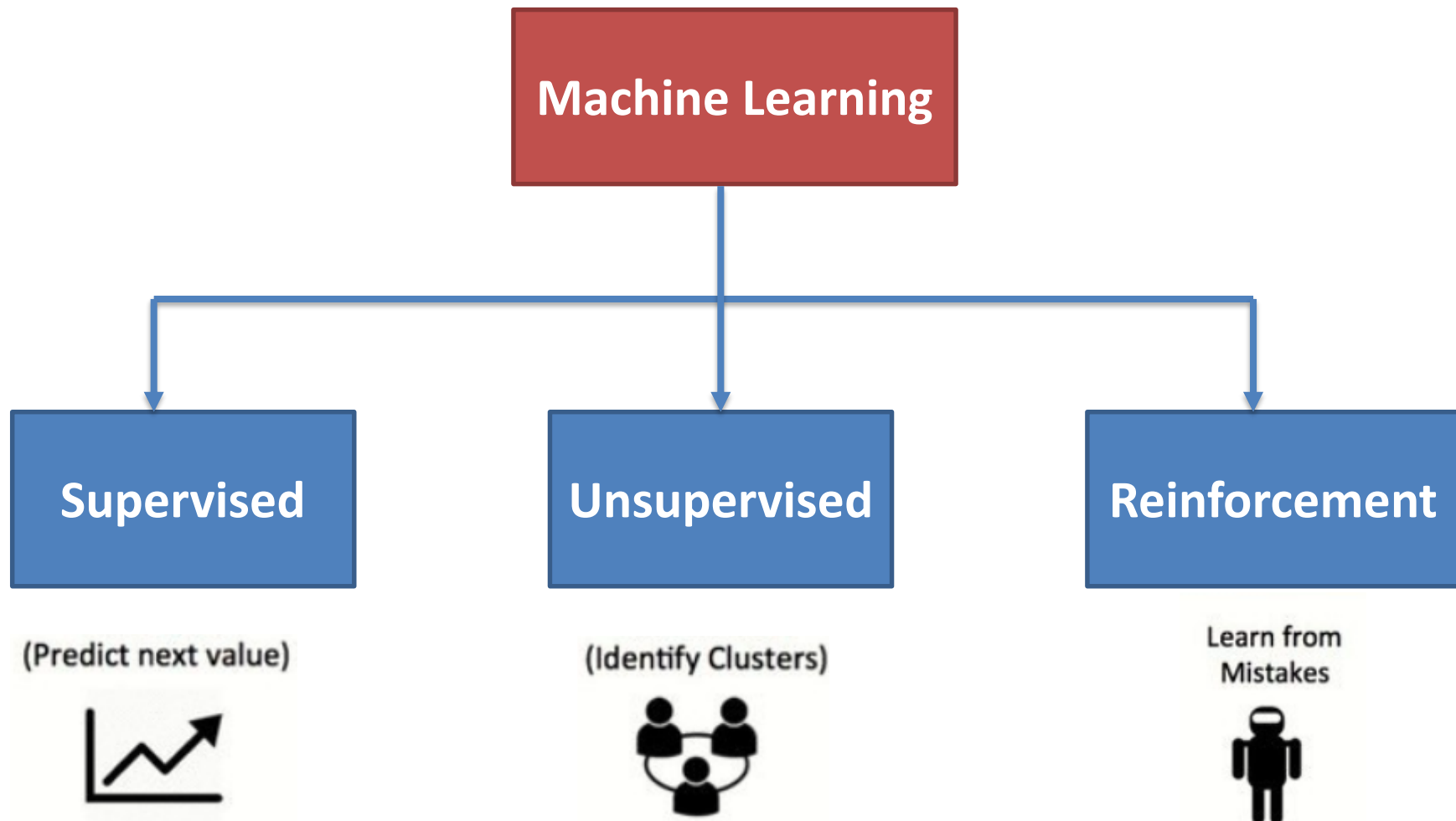
## Traditional programming



## Machine Learning (Unsupervised)



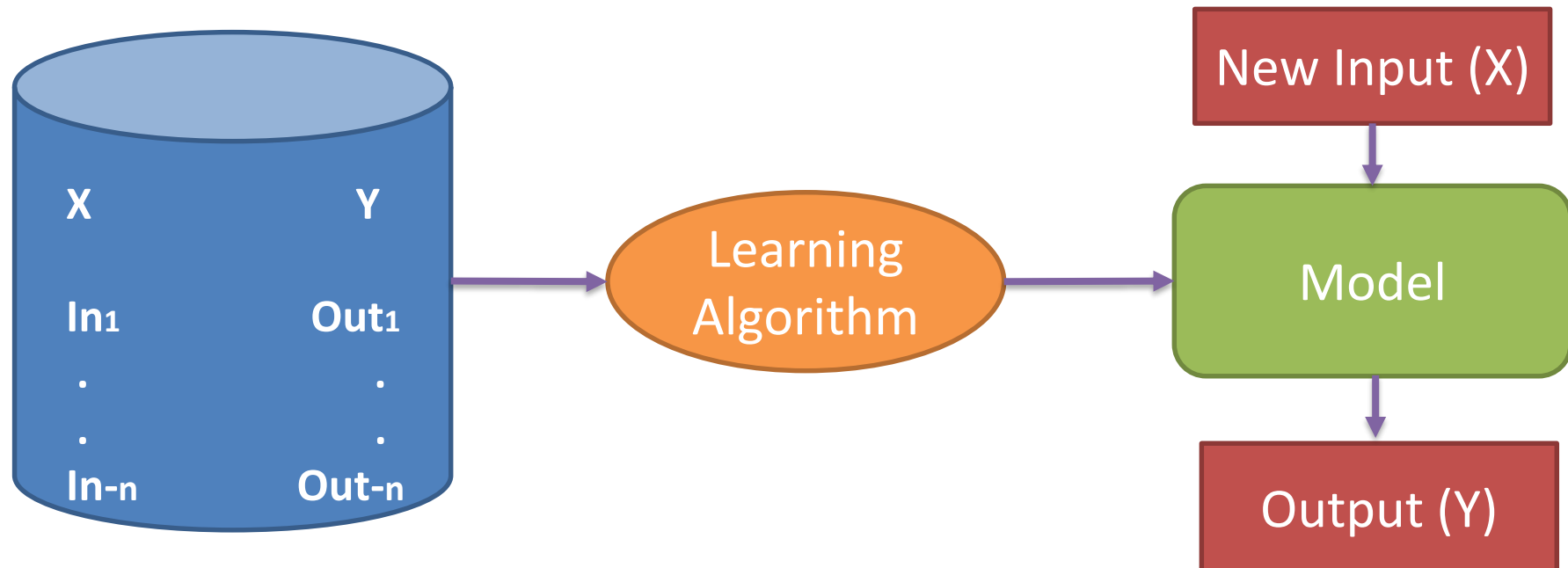
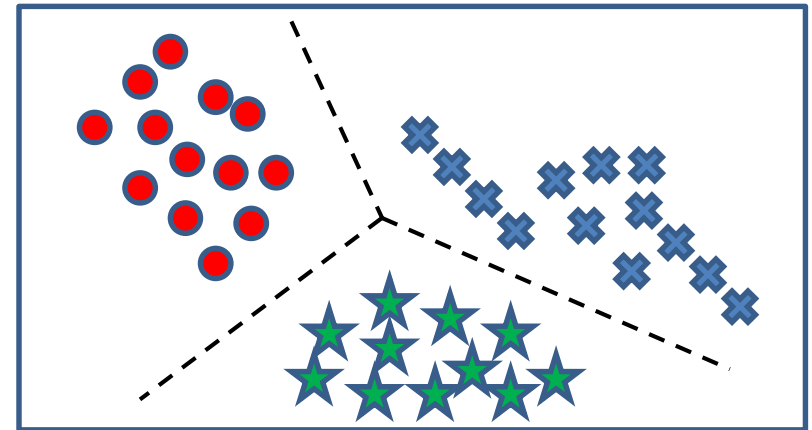
# Types of Machine Learning



# Types of Machine Learning

## Supervised Learning

In Supervised Learning, an algorithm learns from a dataset

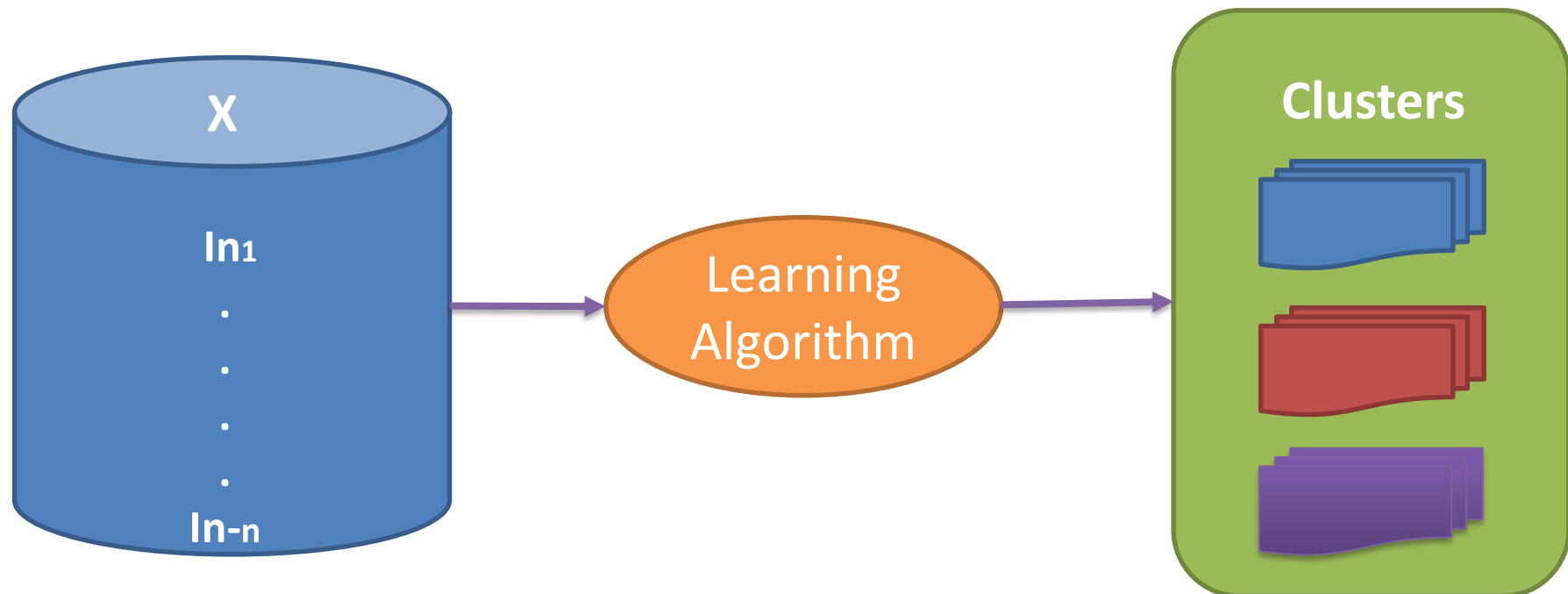
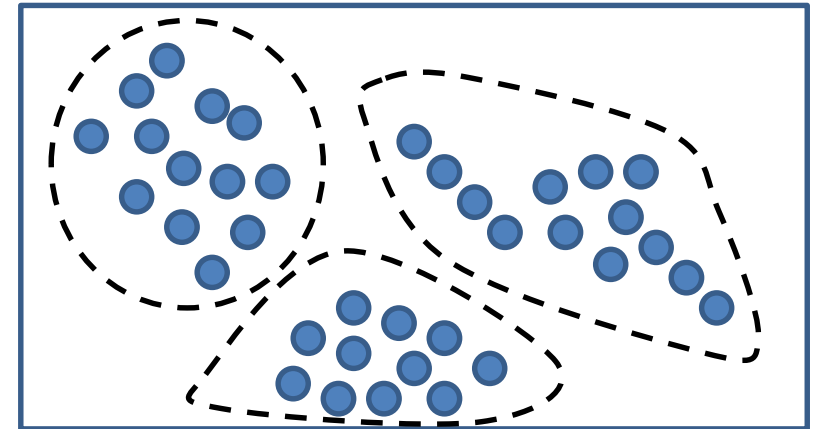




# Types of Machine Learning

## Unsupervised Learning

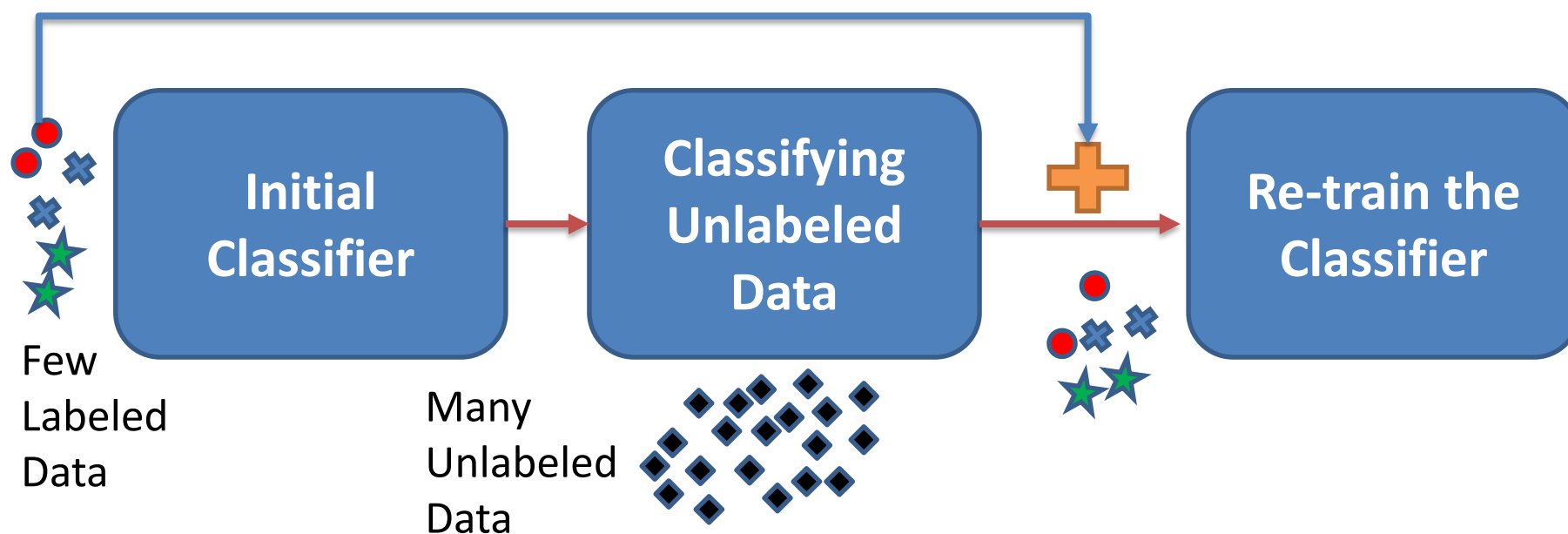
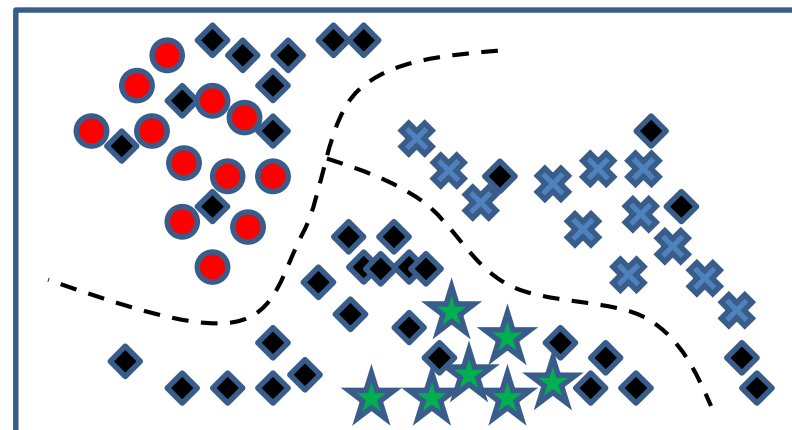
In Unsupervised Learning, an algorithm is provided with only the data



# Types of Machine Learning

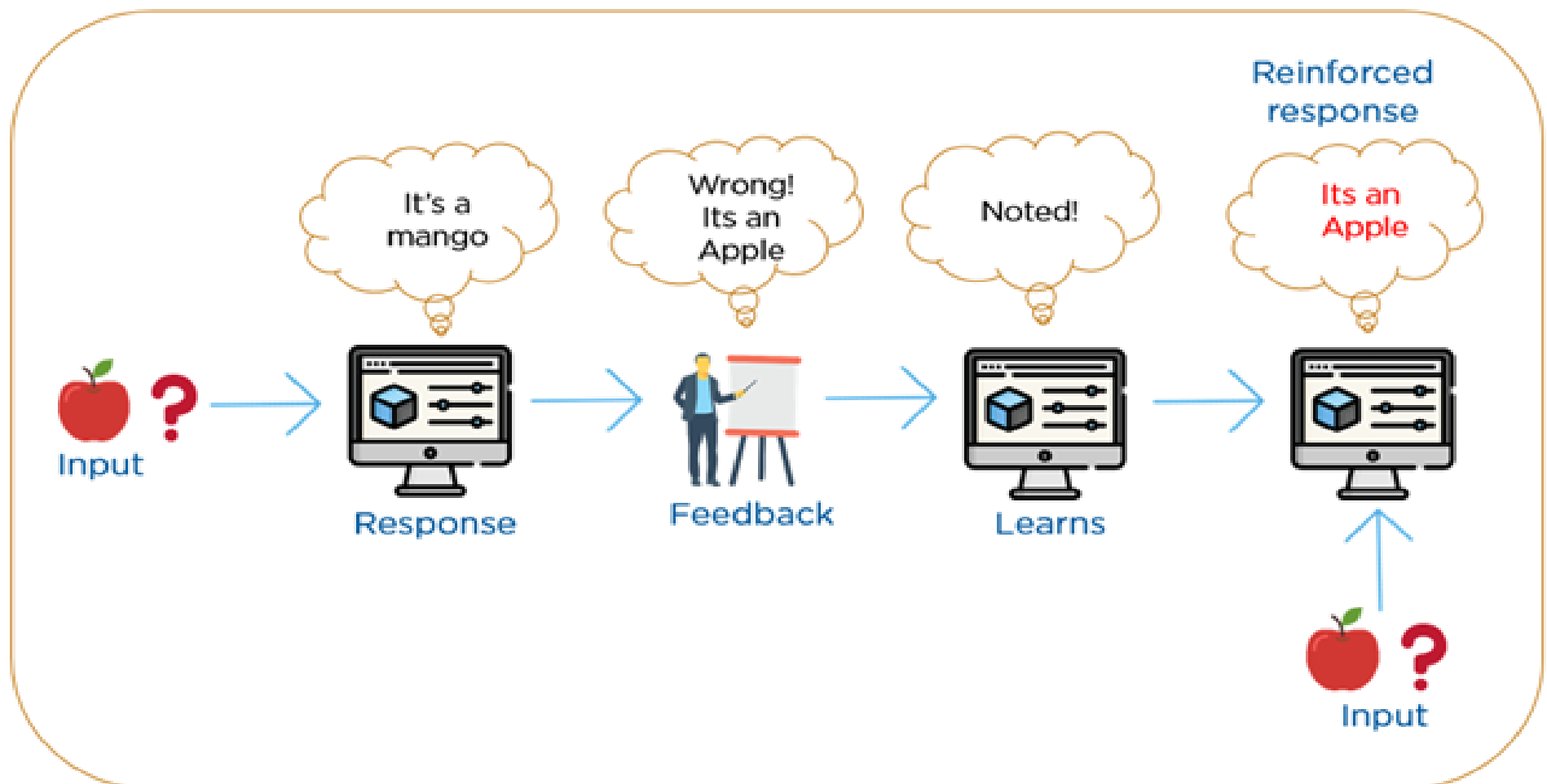
## Semi-supervised Learning

Semi-Supervised Learning is, basically, the combination of Supervised and Unsupervised learning



# Types of Machine Learning

## Reinforcement Learning



# Key Elements of Machine Learning

- Based on state-of-the-art , all machine learning algorithms today are made up of three components. They are as follows:

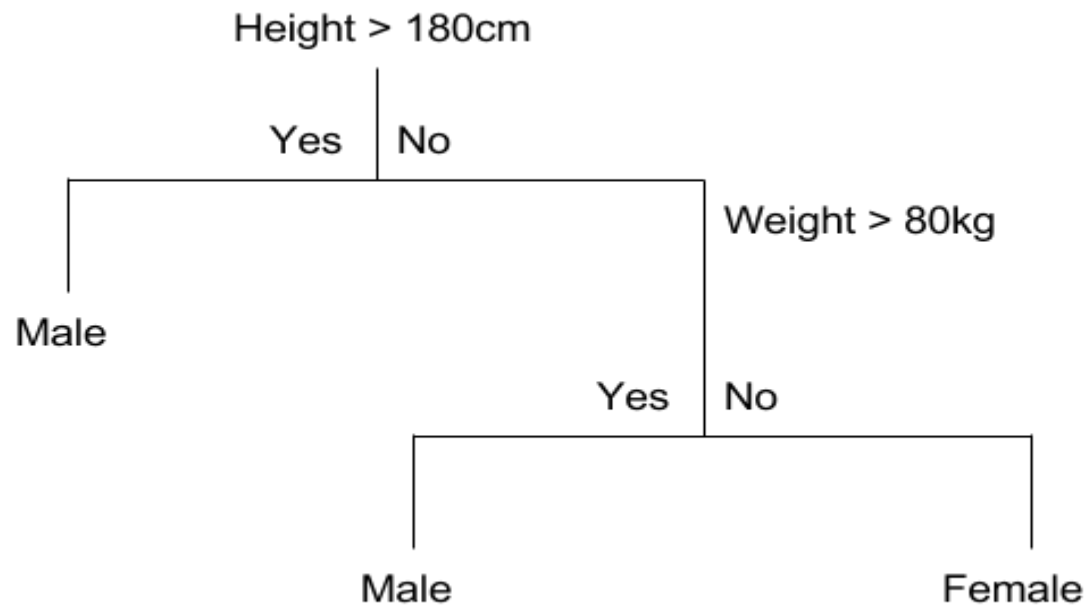
**Representation**

**Evaluation**

**Optimization**

# Key Elements of Machine Learning

## Representation



Information represented by decision trees

# Key Elements of Machine Learning

## Representation

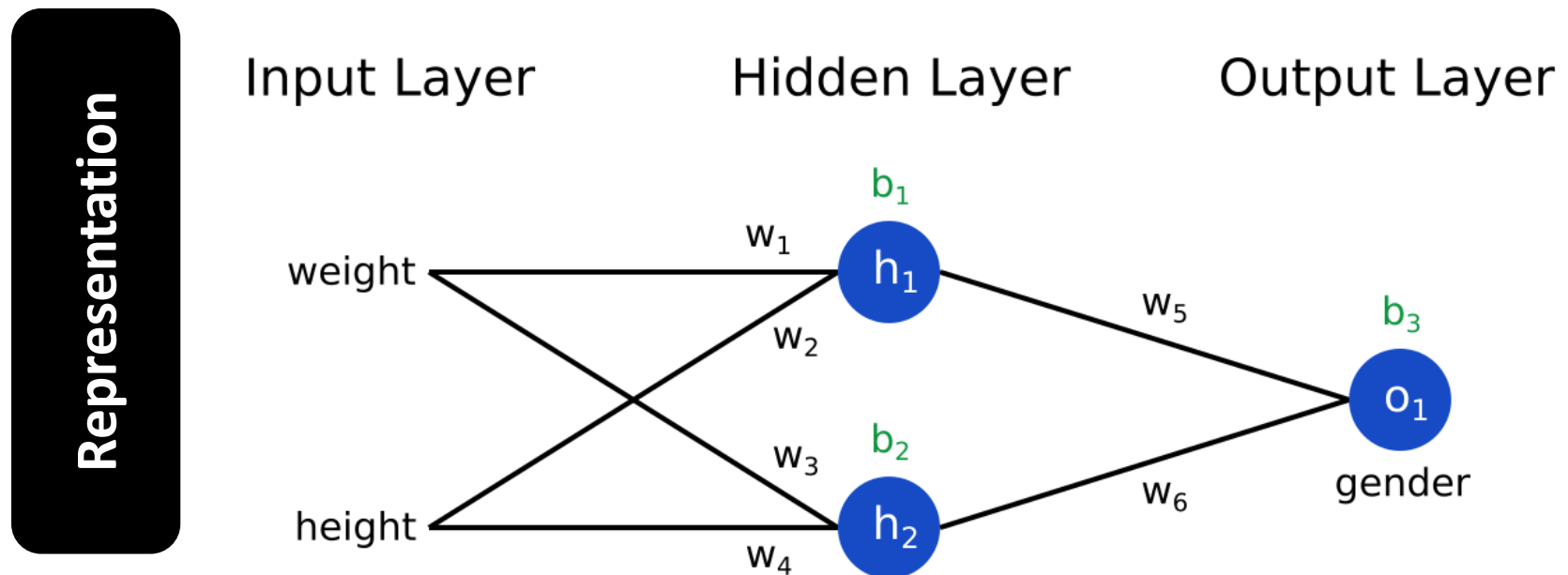
**Height** > 180 cm, then it's probably a male

**Height** > 180 cm and **Weight** > 80 kg, then it's probably a male

**Height** ≤ 180 cm and **Weight** ≤ 80 kg , then it's probably a female

Information represented by set of rules

# Key Elements of Machine Learning



Information represented by neural networks

# Key Elements of Machine Learning

## Evaluation

### Performance Metrics for Classification problems

1. Confusion Matrix
2. Accuracy
3. Precision
4. Recall (sensitivity)
5. F1 Score



# Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...



**Imagine that we have this medical data**

# Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...



**We have got some clinical measurements**

# Compare Machine Learning Methods

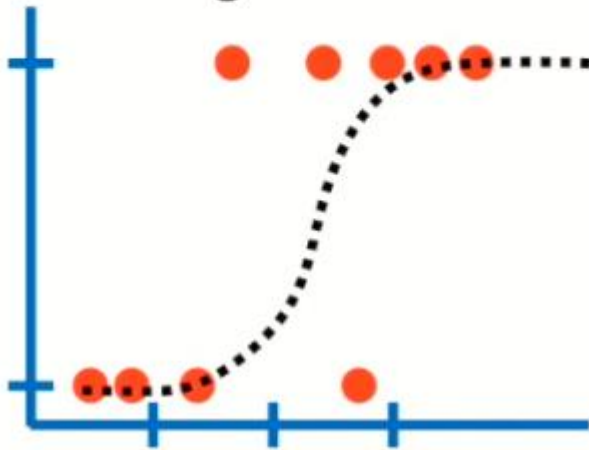
Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...



And we want to apply a ML method to predict whether someone will develop heart disease or not

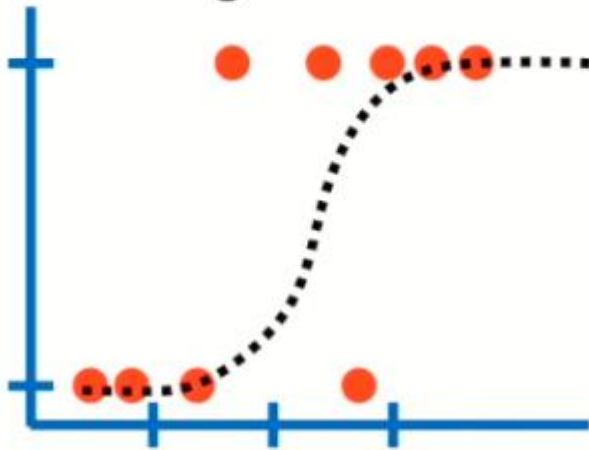
# Compare Machine Learning Methods

We could use Logistic Regression...



# Compare Machine Learning Methods

We could use Logistic Regression...

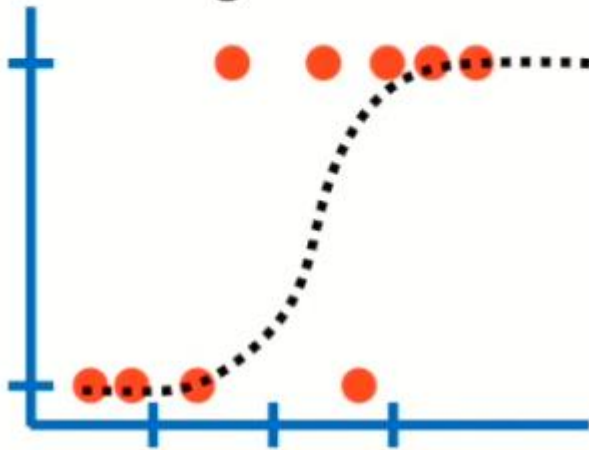


...or K-nearest neighbors...



# Compare Machine Learning Methods

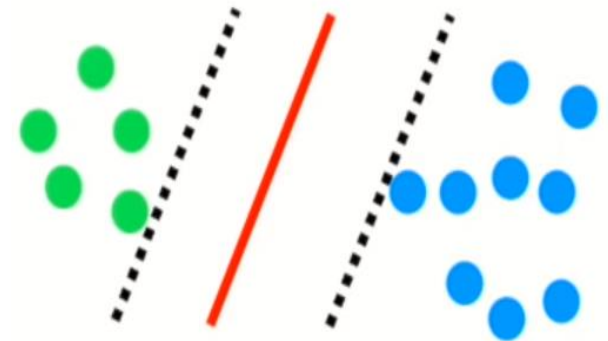
We could use Logistic Regression...



...or K-nearest neighbors...

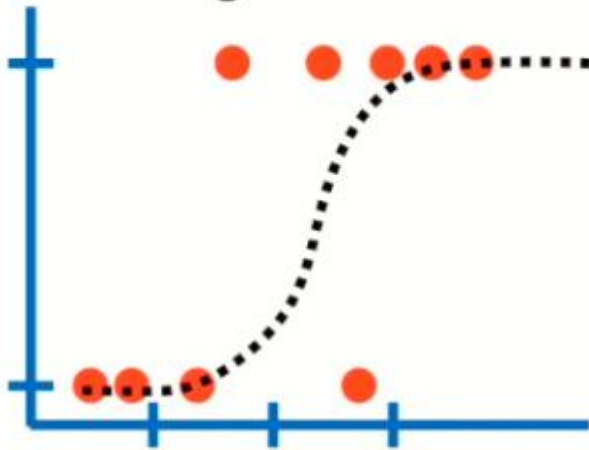


...or support vector machines (SVM)...



# Compare Machine Learning Methods

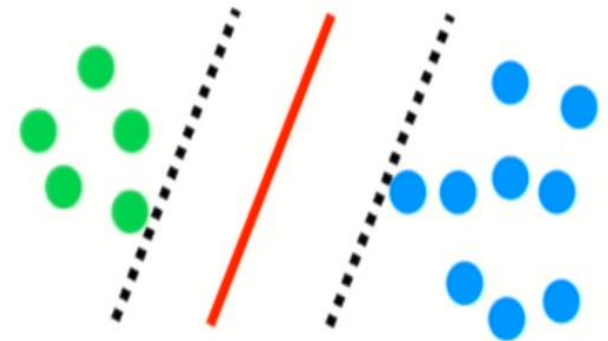
We could use Logistic Regression...



...or K-nearest neighbors...



...or support vector machines (SVM)...



**How to decide which one works better with our data ?**

# Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
...	...	...	...	...



**Training Data**

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	140	No
...	...	...	...	...



**Testing Data**

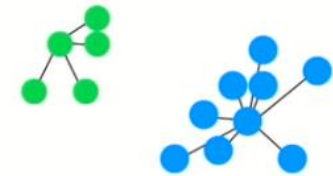


# Compare Machine Learning Methods

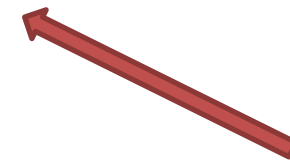
## Confusion Matrix

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
...	...	...	...	...

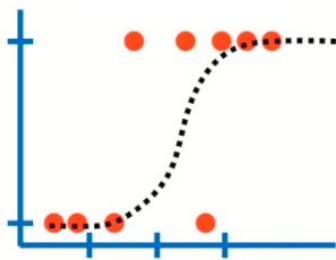
K-nearest neighbors



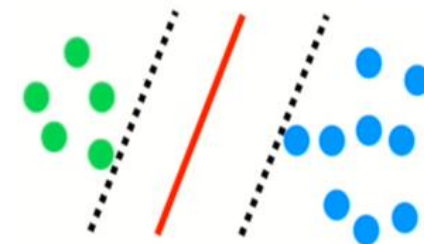
Training Data



Logistic Regression



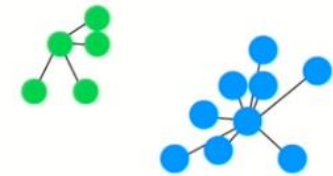
support vector  
machines (SVM)



# Compare Machine Learning Methods

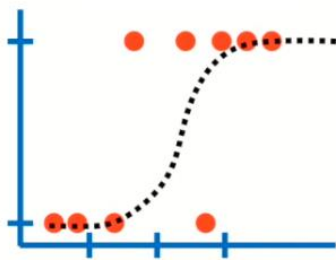
Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	140	No
...	...	...	...	...

K-nearest neighbors

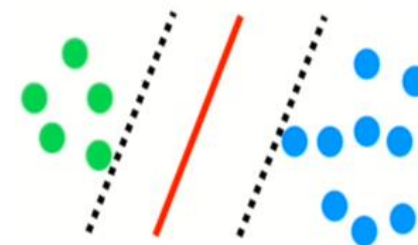


Testing Data

Logistic Regression



support vector machines (SVM)

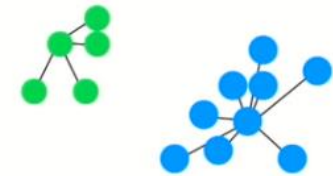


# Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	Yes			
...	...			

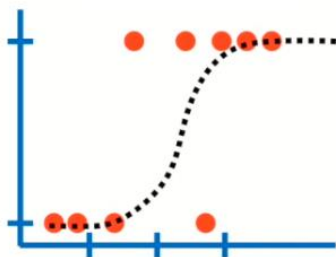
Now we need to summarize how each method performed on the Testing data

K-nearest neighbors

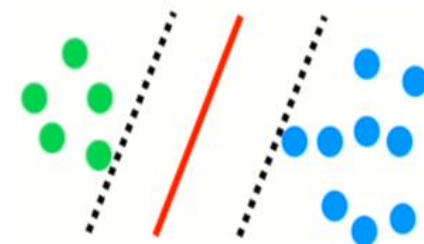


Testing Data

Logistic Regr



support vector machines (SVM)



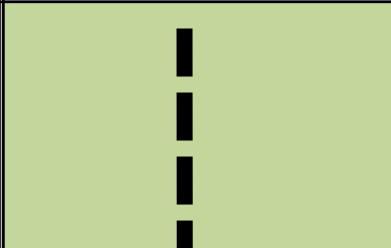
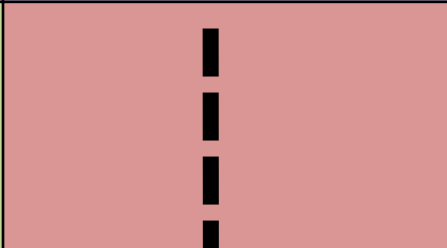
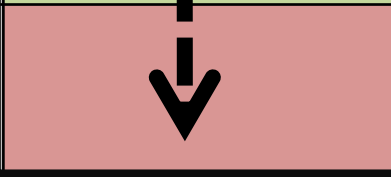
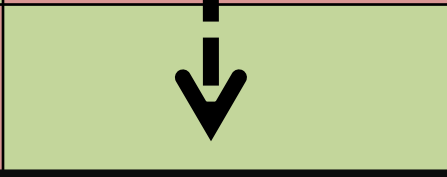


# Confusion Matrix

	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease		
Does Not Have Heart Disease		

The rows correspond to what the ML algorithm predicted

# Confusion Matrix

		Actual	
		Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease			
Does Not Have Heart Disease			

The columns correspond to the known truth

# Confusion Matrix

Actual			
	Has Heart Disease	Does Not Have Heart Disease	
Predicted	Has Heart Disease		
	Does Not Have Heart Disease		

The **Green Boxes** tell us how many times the samples were correctly classified by the algorithm

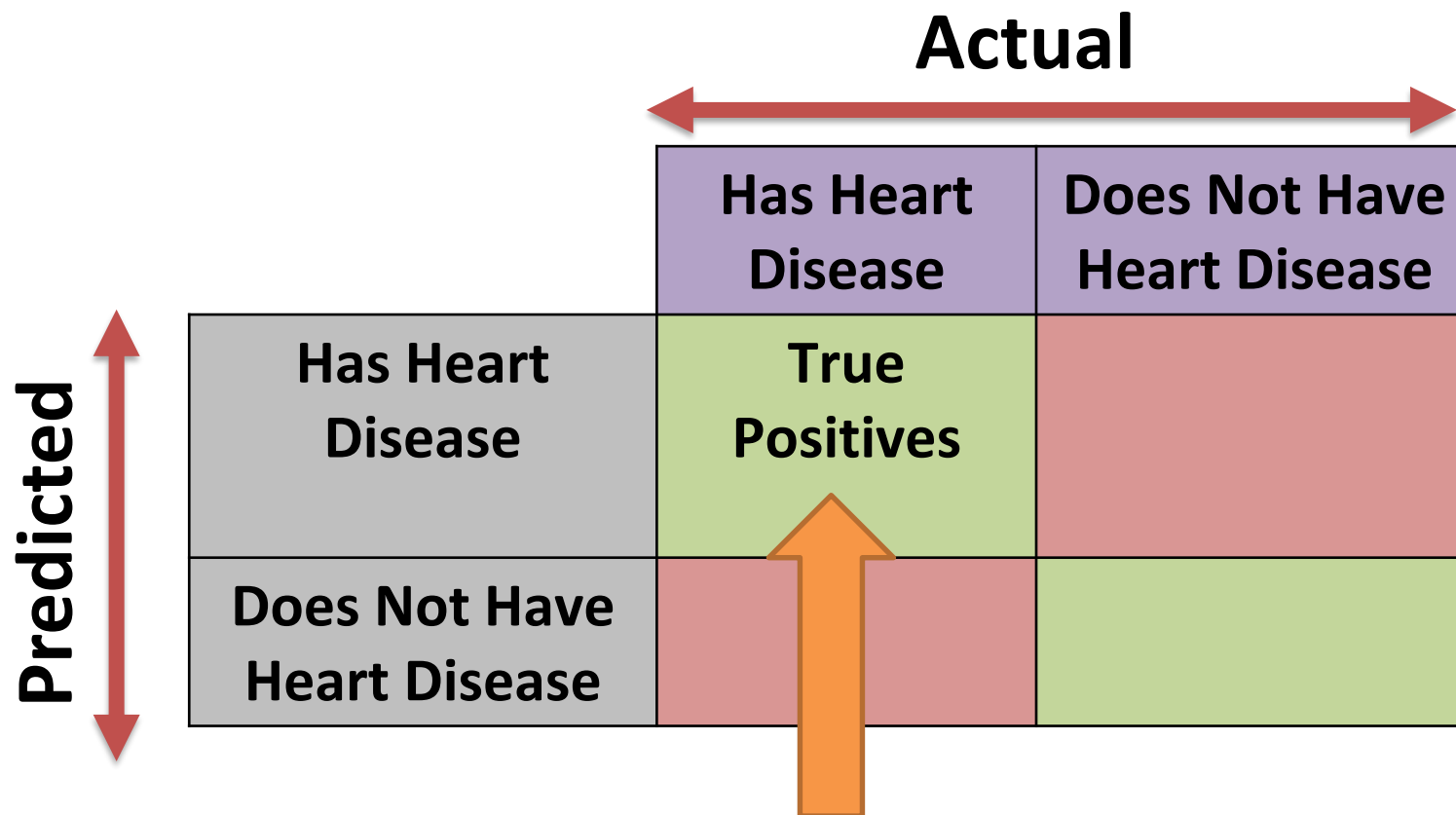
# Confusion Matrix

Actual			
Predicted	Has Heart Disease	Does Not Have Heart Disease	
	Has Heart Disease	Does Not Have Heart Disease	
Has Heart Disease			
Does Not Have Heart Disease			

The **Red Boxes** tell us how many times the samples were misclassified by the algorithm



# Confusion Matrix



The diagram illustrates a Confusion Matrix for heart disease classification. It features a 2x2 grid of cells. The columns are labeled 'Actual' at the top, with 'Has Heart Disease' and 'Does Not Have Heart Disease'. The rows are labeled 'Predicted' on the left, with 'Has Heart Disease' and 'Does Not Have Heart Disease'. The top-left cell (True Positive) is green and labeled 'True Positives'. The top-right cell (False Negative) is red. The bottom-left cell (False Positive) is red. The bottom-right cell (True Negative) is green. A red double-headed arrow is above the columns, and another red double-headed arrow is to the left of the rows. An orange arrow points upwards from the text below to the 'True Positives' cell.

Actual		
	Has Heart Disease	Does Not Have Heart Disease
Predicted	Has Heart Disease True Positives	
	Does Not Have Heart Disease	

Patients that had heart disease and that were correctly identified by the algorithm

# Confusion Matrix

Actual			
Predicted	Has Heart Disease	Does Not Have Heart Disease	
	Has Heart Disease	Does Not Have Heart Disease	
Has Heart Disease	True Positives		
Does Not Have Heart Disease		True Negatives	

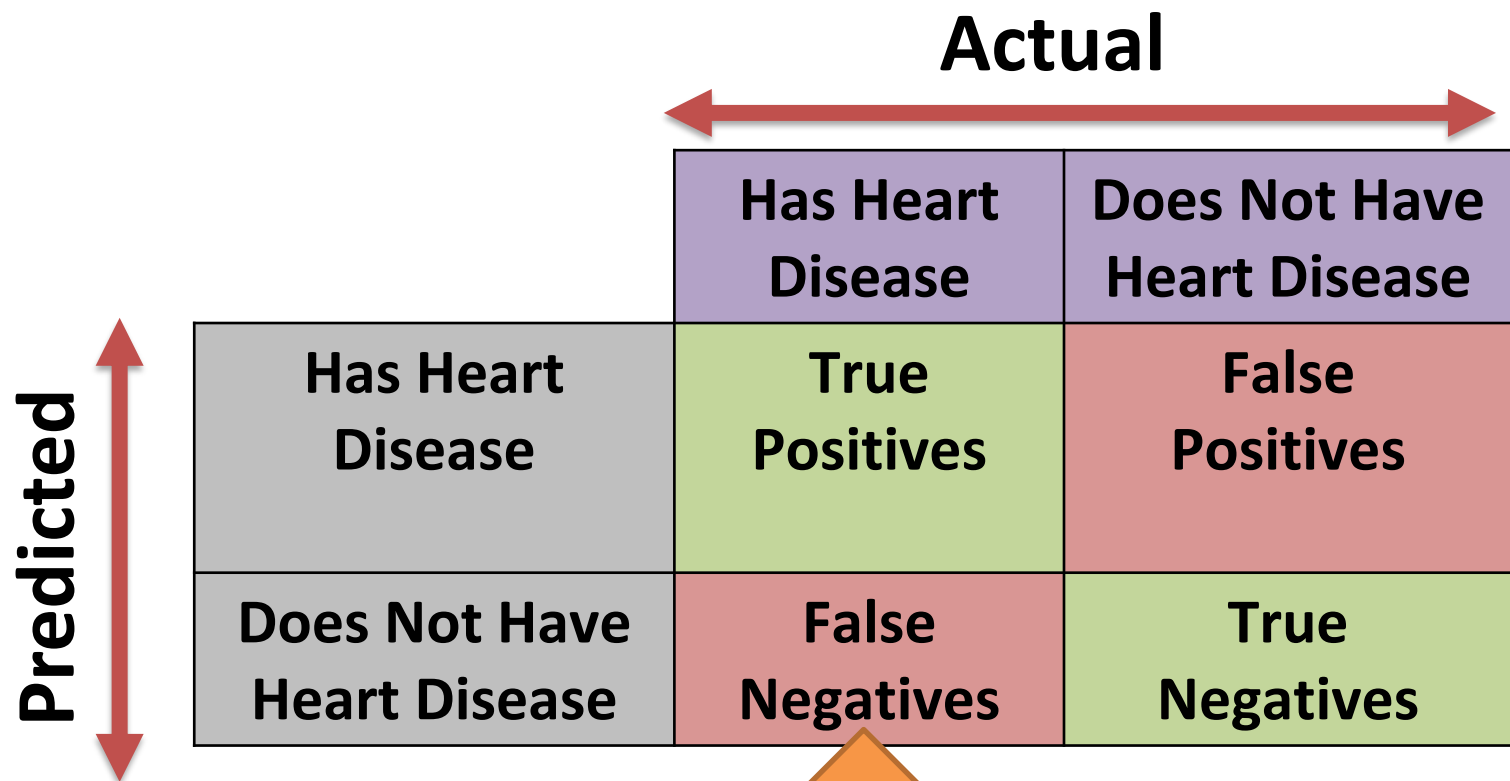
Patients that did not have heart disease and that were correctly identified by the algorithm

# Confusion Matrix

Actual			
Predicted	Has Heart Disease	Does Not Have Heart Disease	
	Has Heart Disease	Does Not Have Heart Disease	
Has Heart Disease	True Positives	False Positives	
Does Not Have Heart Disease		True Negatives	

Patients that did not have heart disease, but the algorithm says they do

# Confusion Matrix



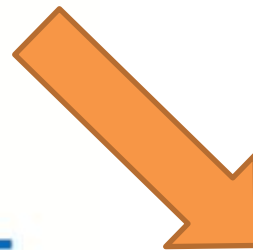
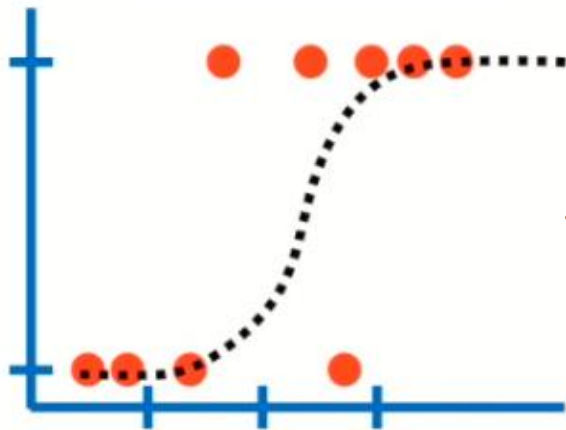
The diagram illustrates a Confusion Matrix for heart disease classification. It features a 2x2 grid of colored cells. The columns are labeled 'Actual' at the top, with a red double-headed arrow above them. The rows are labeled 'Predicted' on the left, with a red double-headed arrow next to them. The cells are: Top-Left (purple) 'Has Heart Disease', Top-Right (purple) 'Does Not Have Heart Disease', Bottom-Left (green) 'True Positives', Bottom-Right (red) 'False Positives', Bottom-Left (red) 'False Negatives', and Bottom-Right (green) 'True Negatives'. An orange arrow points from the text below to the 'False Negatives' cell.

Actual			
Predicted	Has Heart Disease	Does Not Have Heart Disease	
	Has Heart Disease	True Positives	False Positives
Does Not Have Heart Disease	False Negatives	True Negatives	

Patients that had heart disease, but the algorithm says they didn't

# Confusion Matrix

Logistic Regression



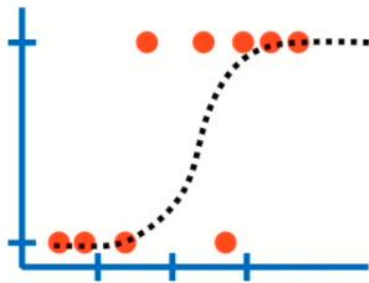
**Actual**

**Predicted**

	Actual	
	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	140	15
Does Not Have Heart Disease	18	127

# Confusion Matrix

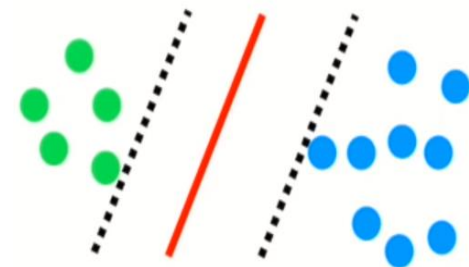
Logistic Regression



	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	140	15
Does Not Have Heart Disease	18	127

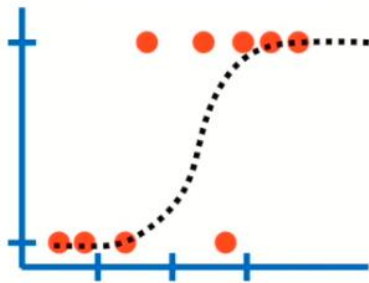
	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	151	8
Does Not Have Heart Disease	9	132

...or support vector machines (SVM)...



# Confusion Matrix

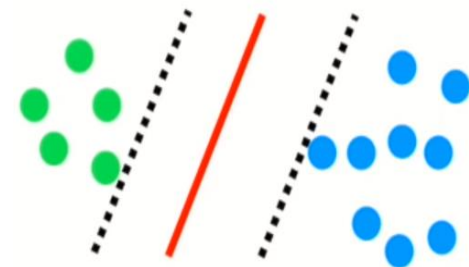
Logistic Regression



	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	140	15
Does Not Have Heart Disease	18	127

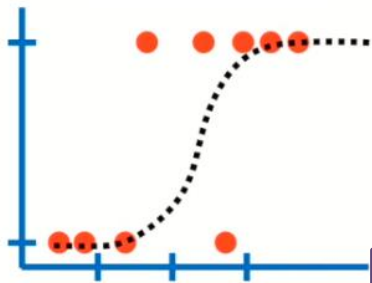
	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	151	8
Does Not Have Heart Disease	9	132

...or support vector machines (SVM)...



# Confusion Matrix

Logistic Regression

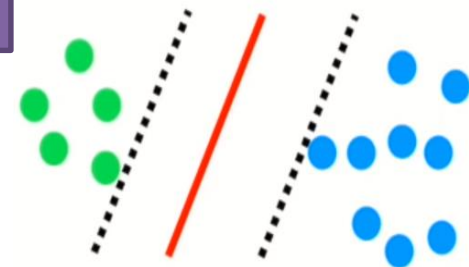


	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	140	15
Does Not Have	18	127

If we had to choose between  
Logistic regression and SVM,  
We would choose SVM !

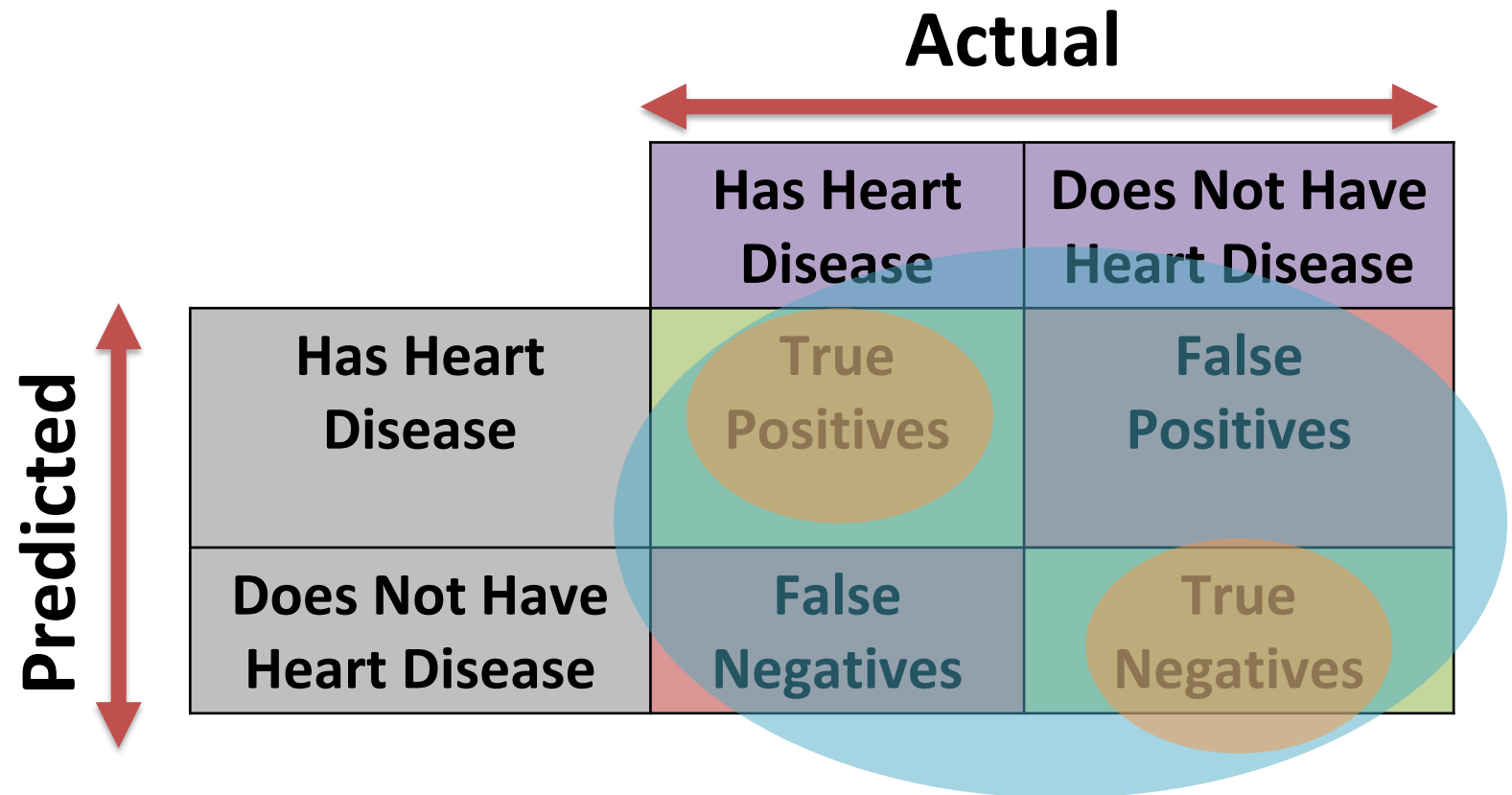
Has Heart Disease	131	9
Does Not Have Heart Disease	9	132

...or support vector machines (SVM)...



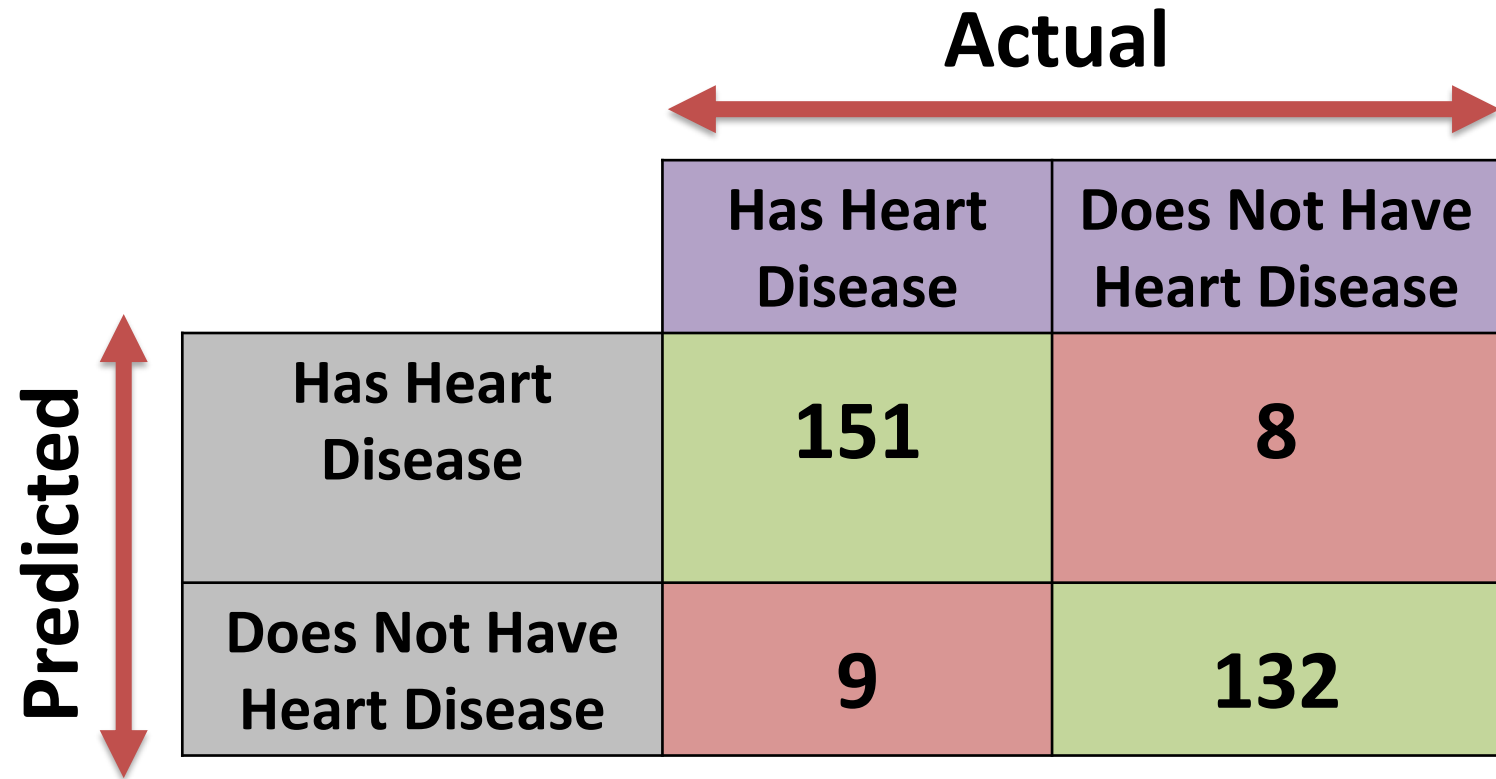


**Accuracy** is the proportion of the total number of predictions that are correct.



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

# Accuracy



A confusion matrix for heart disease prediction. The vertical axis is labeled 'Predicted' with a red double-headed arrow, and the horizontal axis is labeled 'Actual' with a red double-headed arrow. The matrix is a 2x2 grid. The top row represents 'Has Heart Disease' (Actual), and the bottom row represents 'Does Not Have Heart Disease' (Actual). The left column represents 'Has Heart Disease' (Predicted), and the right column represents 'Does Not Have Heart Disease' (Predicted). The cells contain the counts: 151 (True Positive), 8 (False Positive), 9 (False Negative), and 132 (True Negative).

	Actual	
	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	151	8
Does Not Have Heart Disease	9	132

$$Accuracy = \frac{151 + 132}{151 + 132 + 8 + 9} \approx 94.33\%$$

# Accuracy

**Actual**

## When to use Accuracy:

Accuracy is a good measure when the target variable classes in the data are nearly balanced

Have  
ease

$$Accuracy = \frac{151 + 132}{151 + 132 + 8 + 9} \approx 94.33\%$$

# Accuracy

**Actual**

## When NOT to use Accuracy:

Accuracy should NEVER be used as a measure when the target variable classes in data are a majority of 1 class

Have  
ease

$$\text{Accuracy} = \frac{151 + 132}{151 + 132 + 8 + 9} \approx 94.33\%$$

**Precision** shows correctness achieved in positive prediction

		Actual	
		Has Heart Disease	Does Not Have Heart Disease
Predicted	Has Heart Disease	True Positives	False Positives
	Does Not Have Heart Disease	False Negatives	True Negatives

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

# Precision

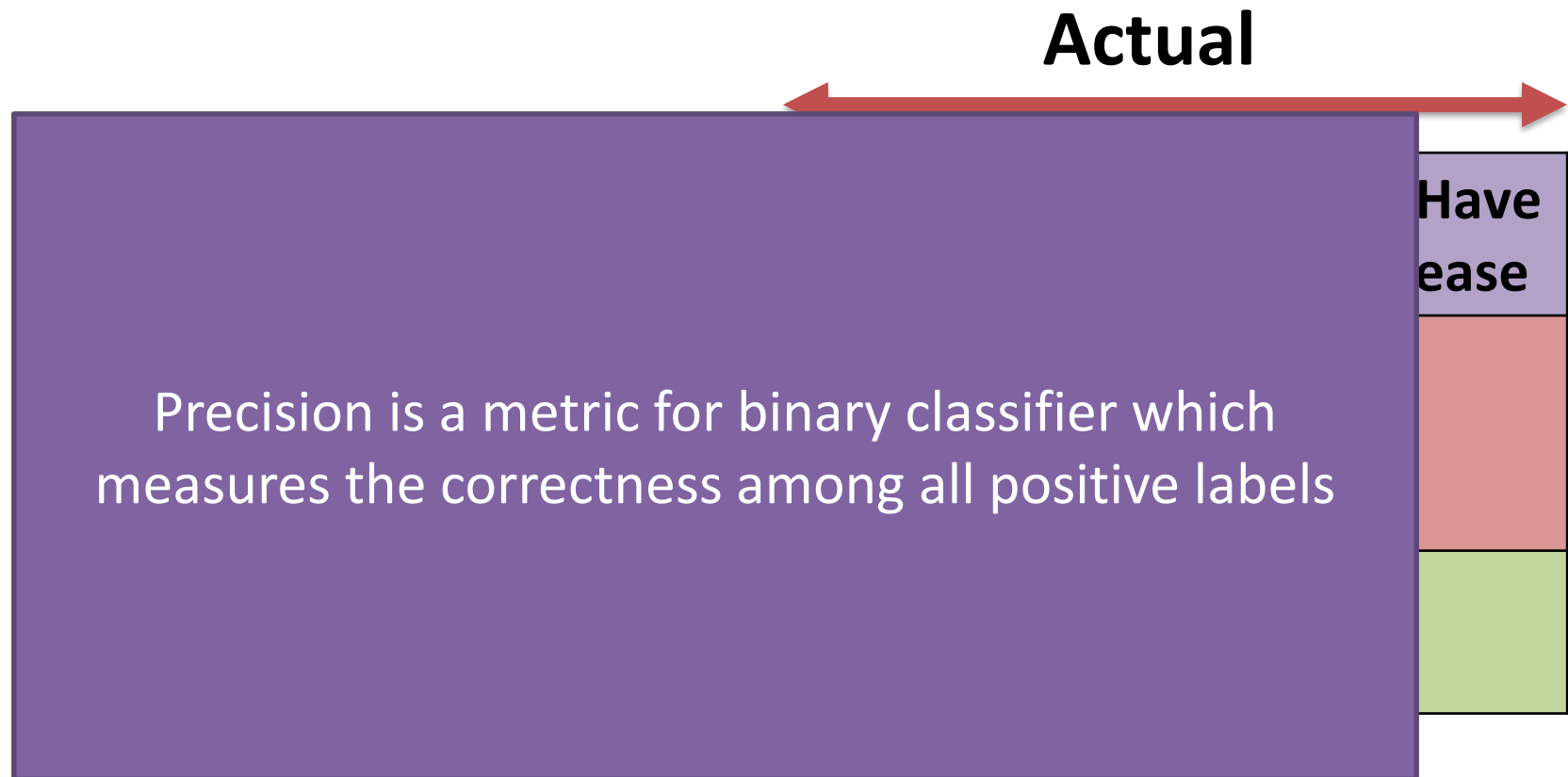
**Actual**

**Predicted**

	Has Heart Disease	Does Not Have Heart Disease
Has Heart Disease	151	8
Does Not Have Heart Disease	9	132

$$Precision = \frac{151}{151 + 8} \approx 94.97\%$$

# Precision



$$\textit{Precision} = \frac{151}{151 + 8} \approx 94.97\%$$

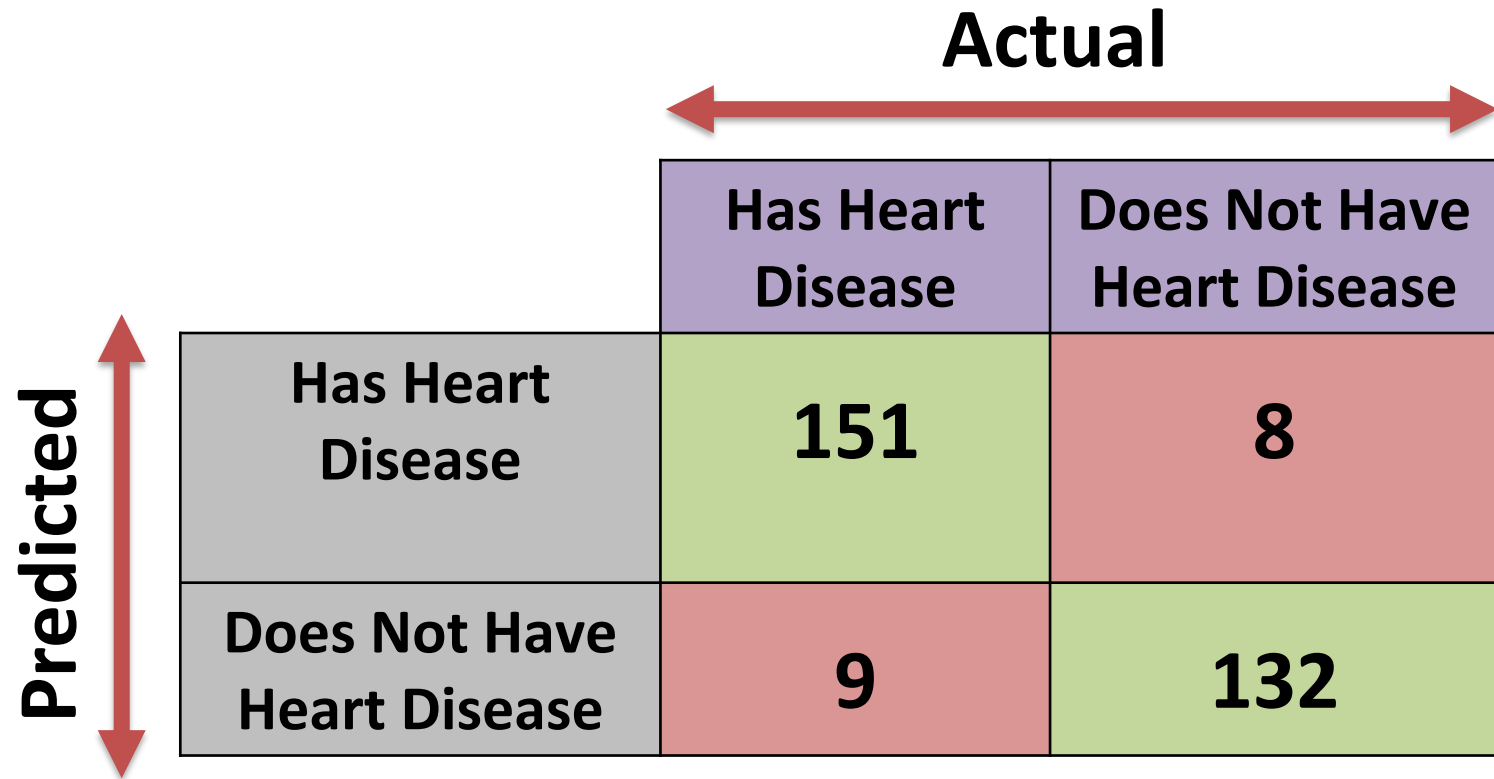
**Recall (Sensitivity)** It is measure of positive examples labeled as positive by classifier

Actual			
Predicted	Has Heart Disease	Does Not Have Heart Disease	
	Has Heart Disease	True Positives	False Positives
	Does Not Have Heart Disease	False Negatives	True Negatives

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



## Recall (Sensitivity)

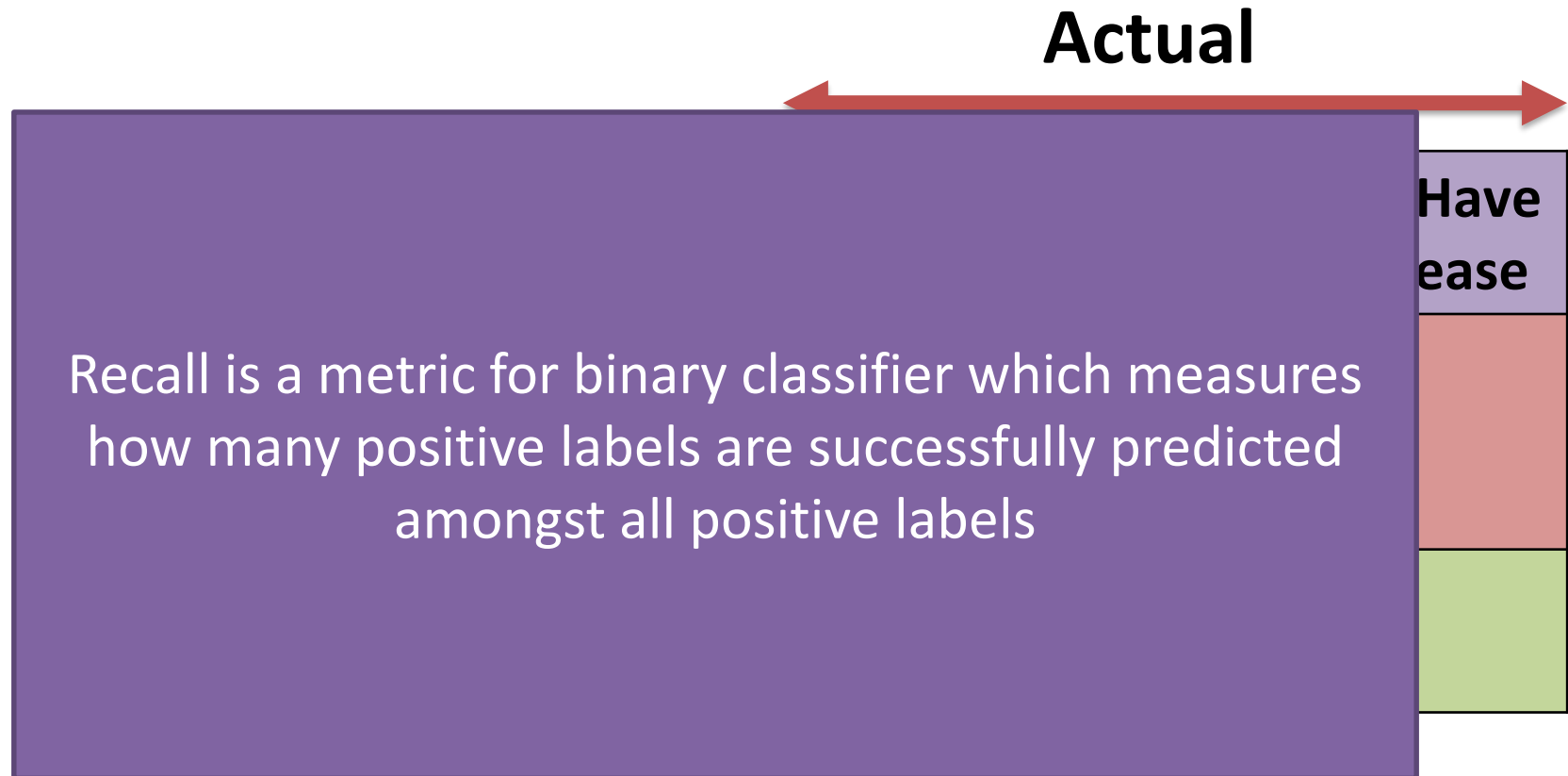


A confusion matrix for heart disease prediction. The vertical axis is labeled 'Predicted' with a red double-headed arrow, and the horizontal axis is labeled 'Actual' with a red double-headed arrow. The matrix is a 2x2 grid. The top row represents 'Actual Has Heart Disease' and the bottom row represents 'Actual Does Not Have Heart Disease'. The left column represents 'Predicted Has Heart Disease' and the right column represents 'Predicted Does Not Have Heart Disease'. The counts are: 151 (True Positives), 8 (False Negatives), 9 (False Positives), and 132 (True Negatives).

		Actual	
		Has Heart Disease	Does Not Have Heart Disease
Predicted	Has Heart Disease	151	8
	Does Not Have Heart Disease	9	132

$$Recall = \frac{151}{151 + 9} \approx 94.38\%$$

## Recall (Sensitivity)



$$\text{Recall} = \frac{151}{151 + 9} \approx 94.38\%$$

**F1 Score** is a weighted average of the recall and precision

- F1 score might be good choice when you seek to balance between Precision and Recall.

**F1 Score** is a weighted average of the recall and precision

- F1 score might be good choice when you seek to balance between Precision and Recall.

$$F1\ Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

**F1 Score** is a weighted average of the recall and precision

- F1 score might be good choice when you seek to balance between Precision and Recall.

$$F1\ Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

- It helps to compute recall and precision in one equation so that the problem to distinguish the models with low recall and high precision or vice versa could be solved

## F1 Score

$$\textit{Precision} = \frac{151}{151 + 8} \approx 94.97\%$$

$$\textit{Recall} = \frac{151}{151 + 9} \approx 94.38\%$$

$$\begin{aligned}\textit{F1 Score} &= 2 \times \frac{\textit{Precision} \times \textit{Recall}}{\textit{Precision} + \textit{Recall}} \\ &= 2 \times \frac{94.97 \times 94.38}{94.97 + 94.38} = 94.97\%\end{aligned}$$

# Key Elements of Machine Learning

## Evaluation

### Performance Metrics for Regression problems

1. Mean Squared Error (MSE)
2. Root Mean Squared Error (RMSE)
3. Mean Absolute Error (MAE)
4. R Squared ( $R^2$ )
5. Adjusted R Squared ( $R^2$ )

# Key Elements of Machine Learning

## Optimization

This is the way candidate programs are generated, also known as the **search process**

1. Combinatorial optimization
2. Convex optimization
3. Constrained optimization



# Compare Machine Learning Methods

Starting with some data ...

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...

## Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...

**We want to use some variables such as Chest pain, Good blood circulation, Blocked arteries and the Weight**

# Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...

To predict if someone has heart disease or not

## Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...

**When a new patient shows up**

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	No	No	115	

## Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...

**And predict if that person have heart disease or not**

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	No	No	115	Yes/ No ?

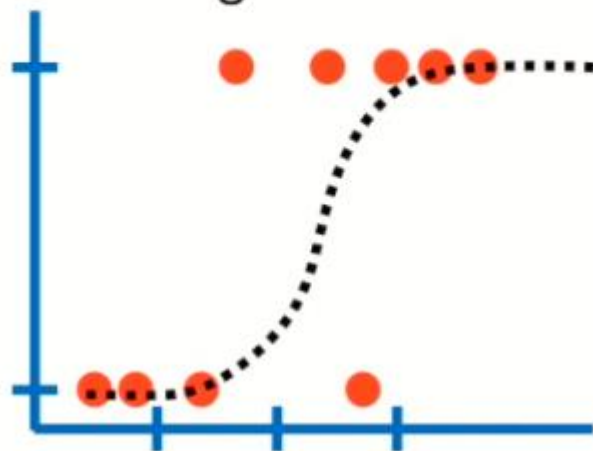
## Compare Machine Learning Methods

Chest Pain	Blood Circulation	Blocked Arteries	Weight	Heart Disease
No	No	No	80	No
Yes	Yes	Yes	125	Yes
Yes	Yes	No	140	No
...	...	...	...	...

**However, first we have to decide which Machine Learning method would be best for our actual problem ...**

# Compare Machine Learning Methods

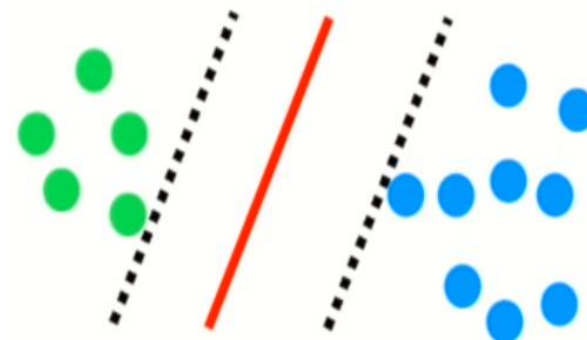
We could use Logistic Regression...



...or K-nearest neighbors...

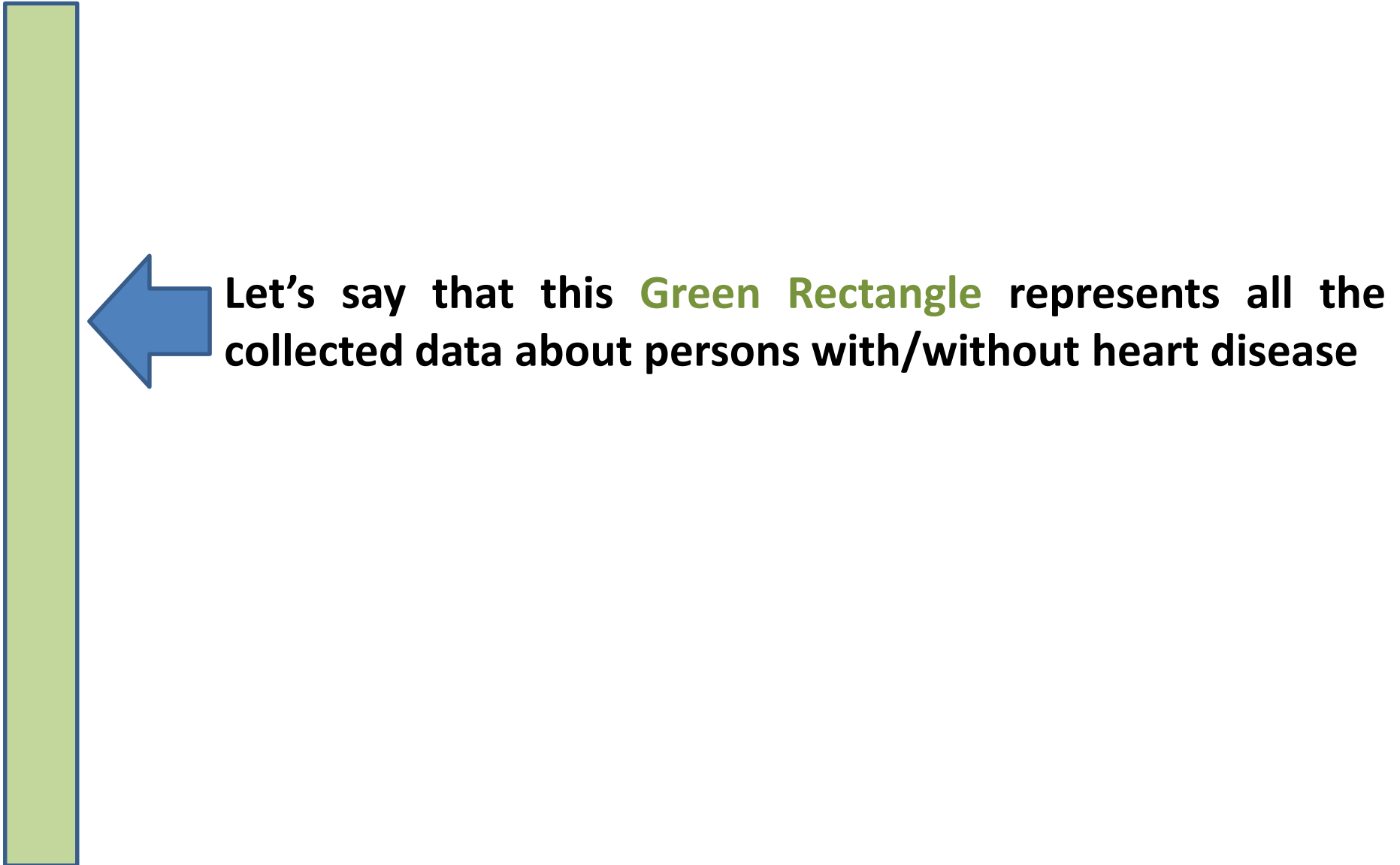


...or support vector machines (SVM)...



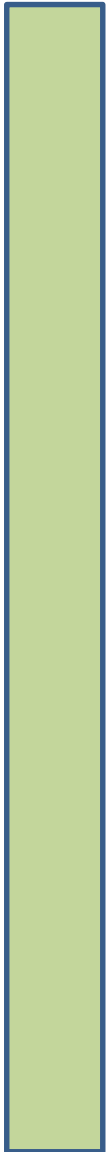
**Cross validation** allows us to compare different machine learning methods and get a sense of how they will work in practice

# Cross Validation



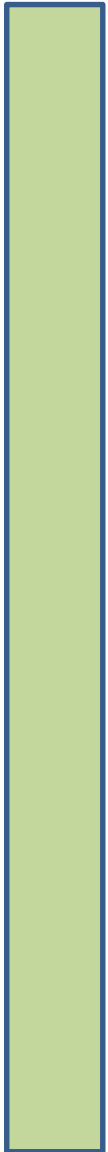


# Cross Validation



**We need to do two things with the data !**

# Cross Validation



**We need to do two things with the data !**

- ✓ **Estimate the parameters for the Machine Learning methods**

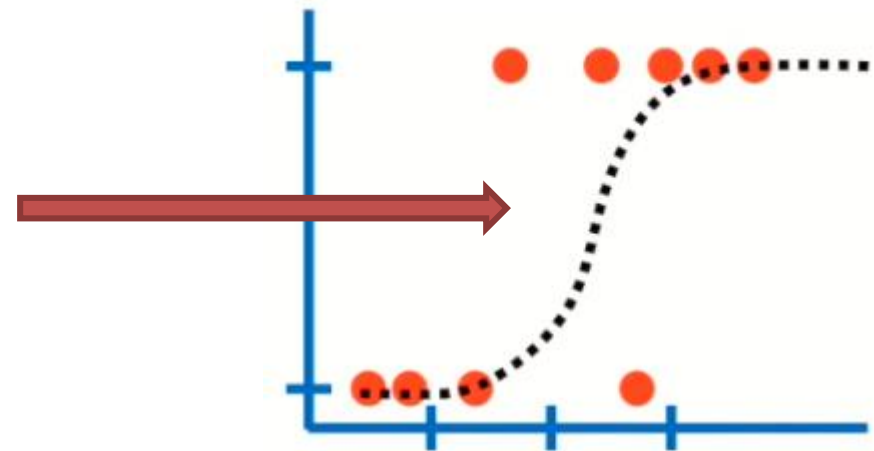
# Cross Validation



**We need to do two things with the data !**

- ✓ **Estimate the parameters for the Machine Learning methods**

If we take Logistic regression,  
we have to use some data to  
estimate the shape of this  
curve



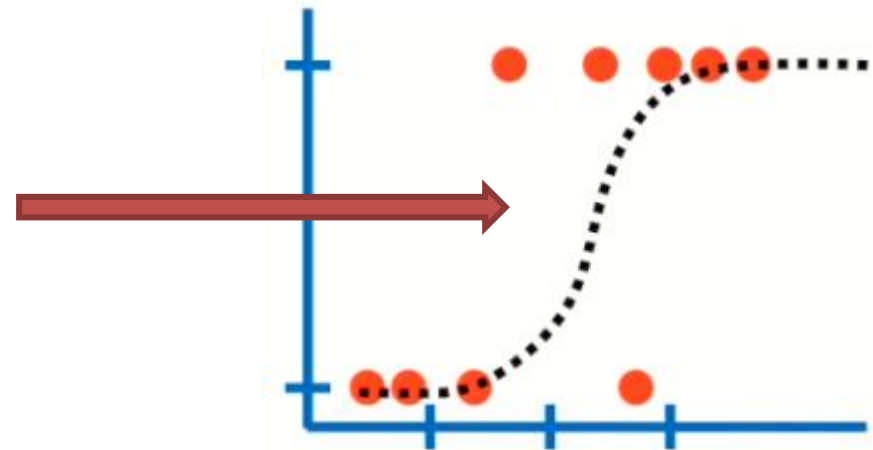
# Cross Validation



**We need to do two things with the data !**

- ✓ **Estimate the parameters for the Machine Learning methods**

In Machine Learning lingo,  
estimating parameters is called  
**training the model**



# Cross Validation



**We need to do two things with the data !**

- ✓ **Estimate the parameters for the Machine Learning methods**
- ✓ **Evaluate how well Machine Learning methods work**

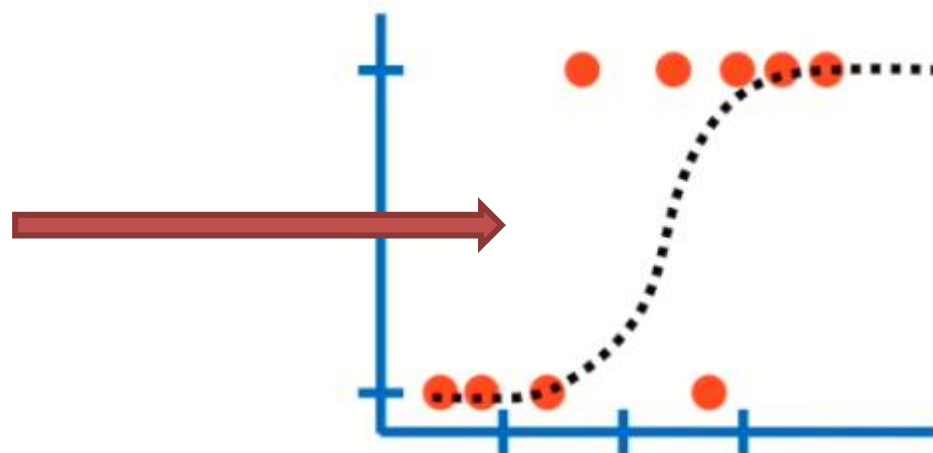
# Cross Validation



**We need to do two things with the data !**

- ✓ **Estimate the parameters for the Machine Learning methods**
- ✓ **Evaluate how well Machine Learning methods work**

Does this curve do a good job while categorizing new data ?



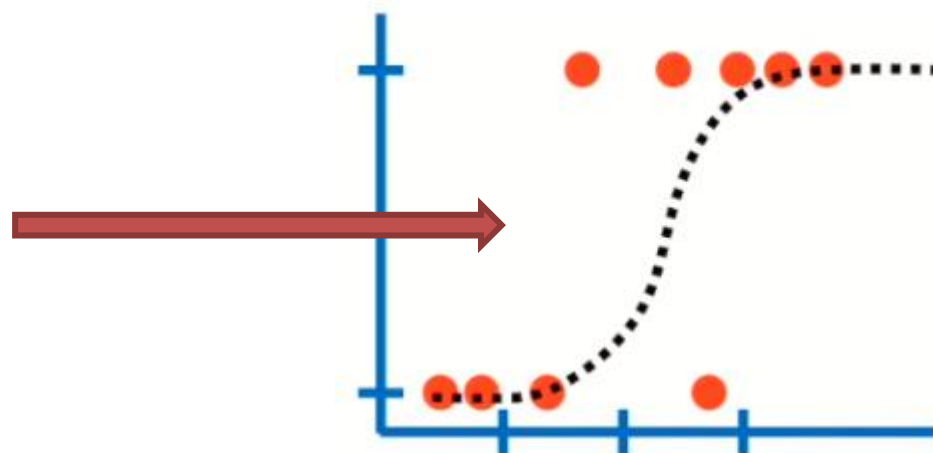
# Cross Validation



**We need to do two things with the data !**

- ✓ **Estimate the parameters for the Machine Learning methods**
- ✓ **Evaluate how well Machine Learning methods work**

In Machine Learning lingo,  
evaluating a method is called  
**testing the model**



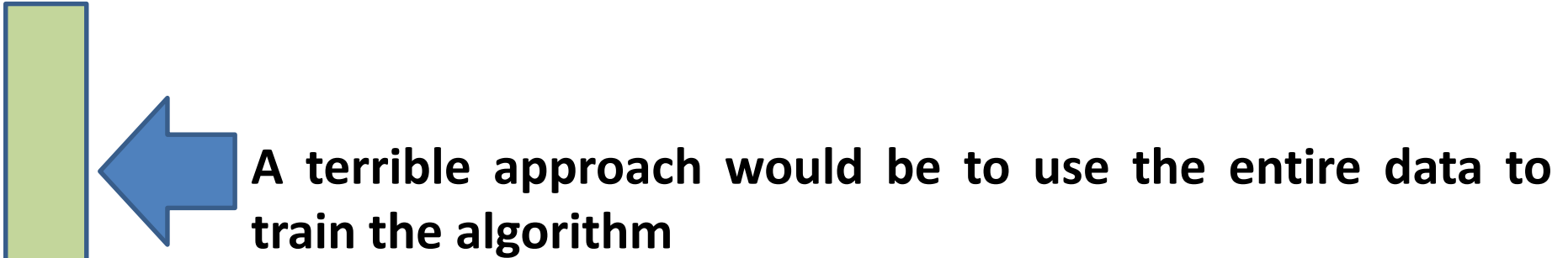
# Cross Validation

By using the Machine Learning lingo we will:

- ❖ ~~Estimate the parameters for~~ **Train** the Machine Learning methods
- ❖ ~~Evaluate how well~~ **Test** the Machine Learning methods



## Cross Validation

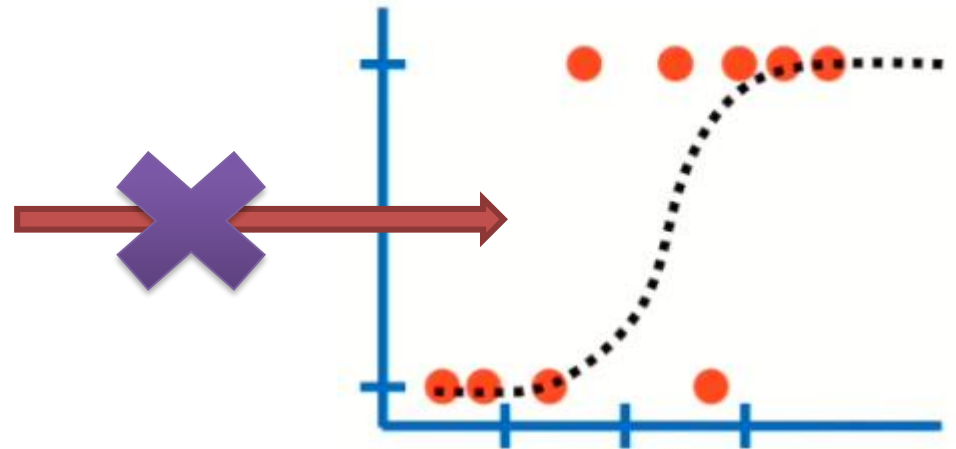


# Cross Validation



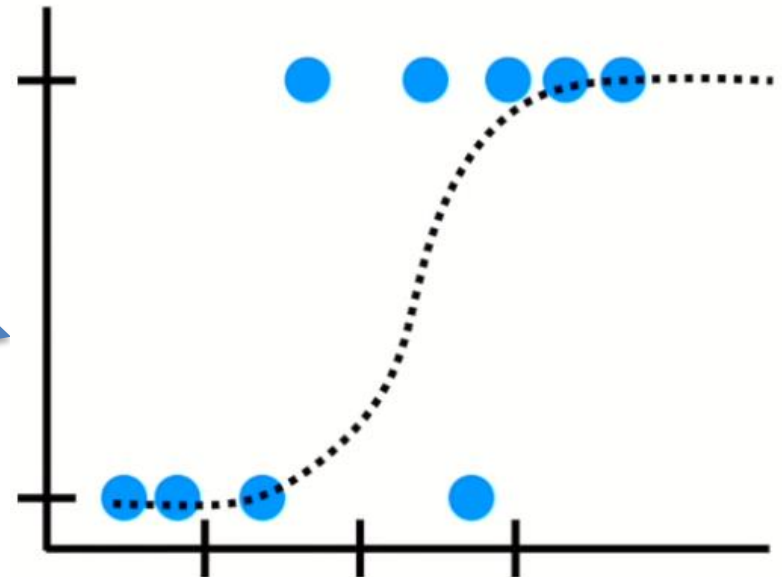
A terrible approach would be to use the entire data to train the algorithm

Reusing the same data for both training and testing is a **bad idea !**



# Cross Validation

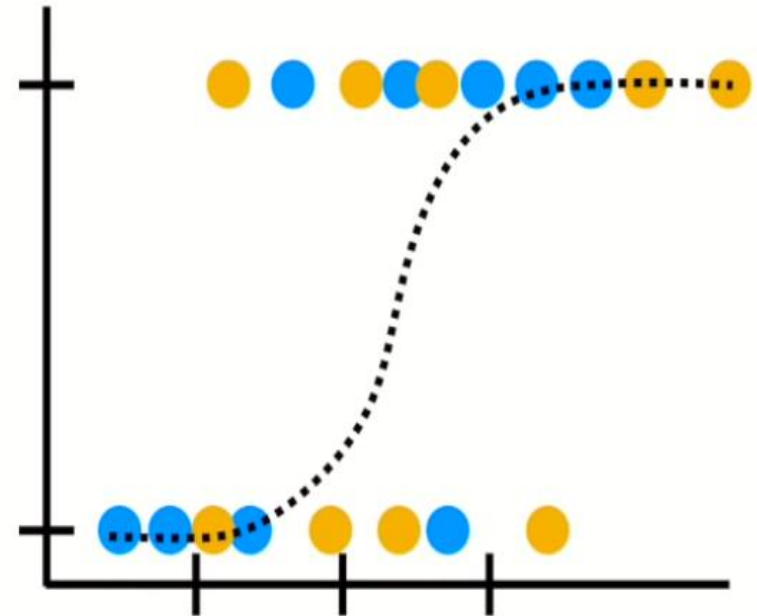
**A better approach would be to use the first 75% of the data for training phase**



# Cross Validation

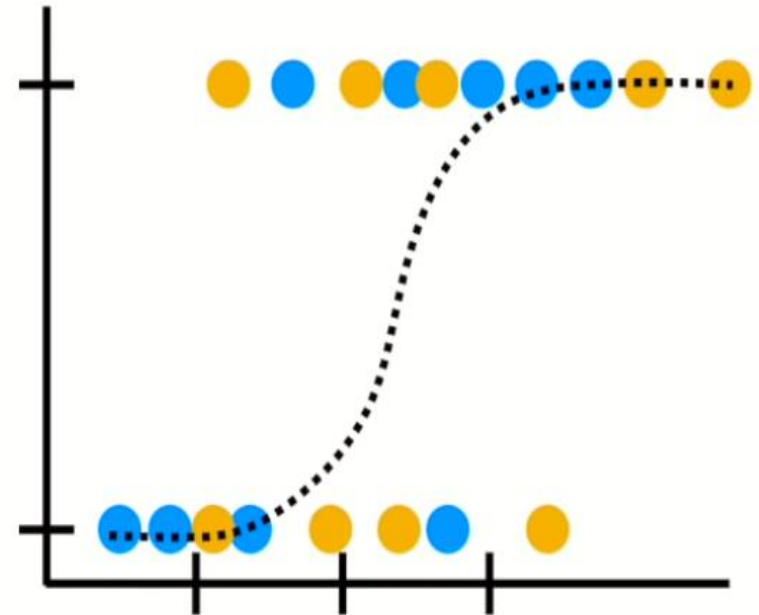
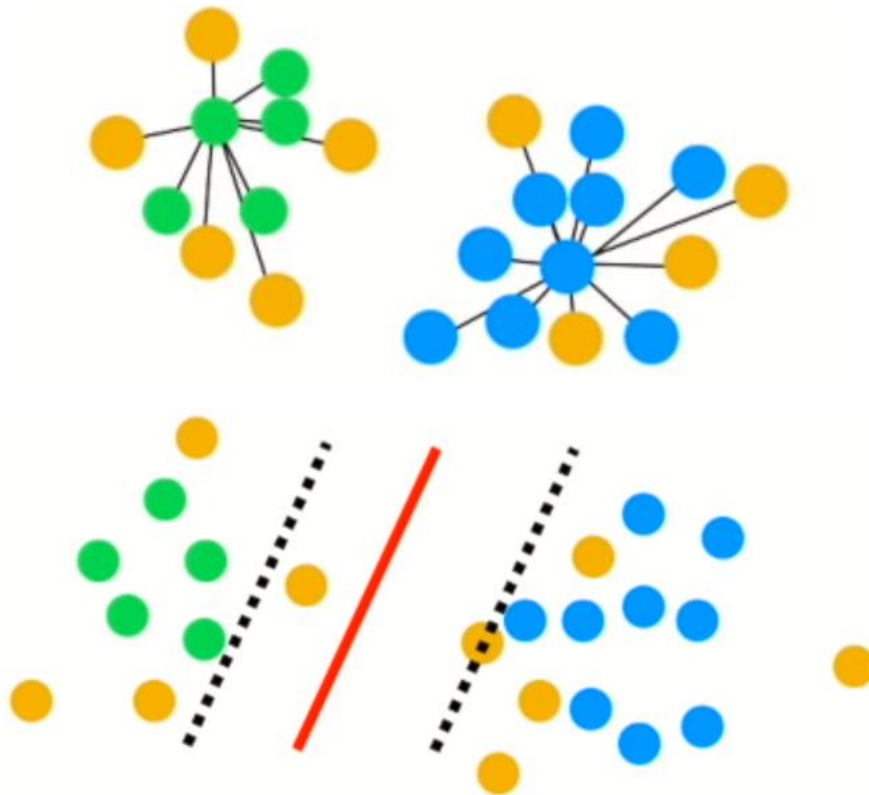


And use the last 25% of the data for testing



# Cross Validation

We can now compare methods by seeing how well each one categorized the test data



# Cross Validation



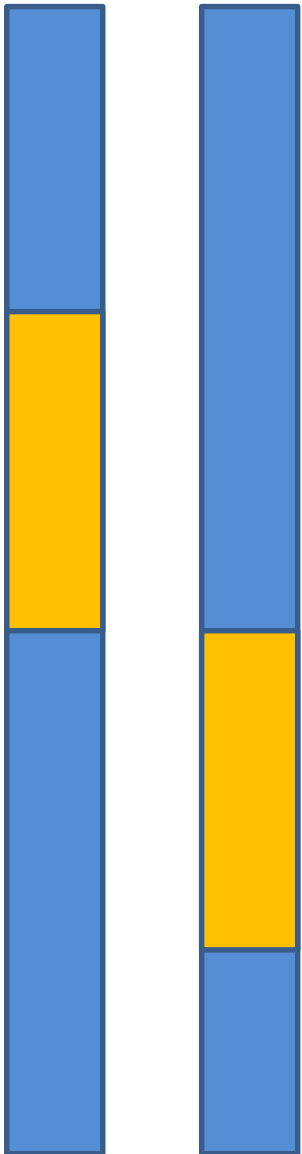
**How do we know that using the first 75% as training data and the last 25% as testing data is the best way to divide the data ?**

# Cross Validation



**What if we inverse the order ?**

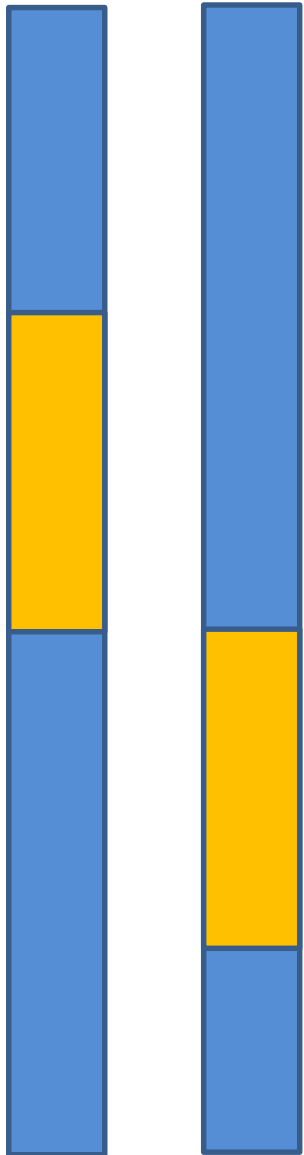
## Cross Validation



**What if we take one of the middle rectangles ?**

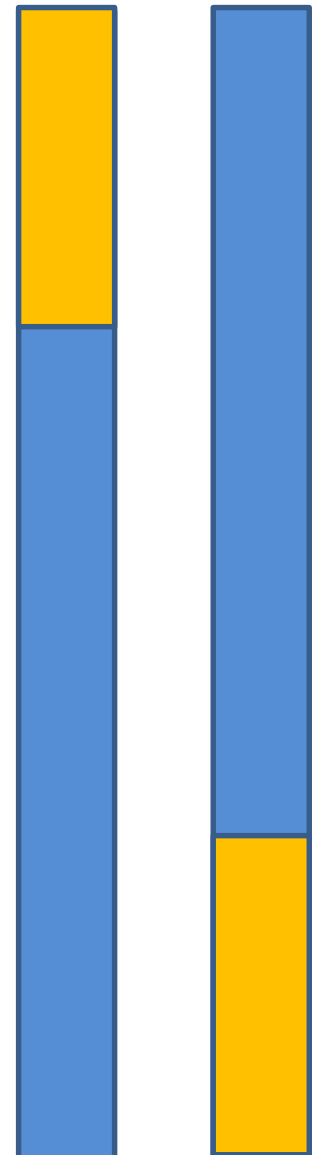


# Cross Validation



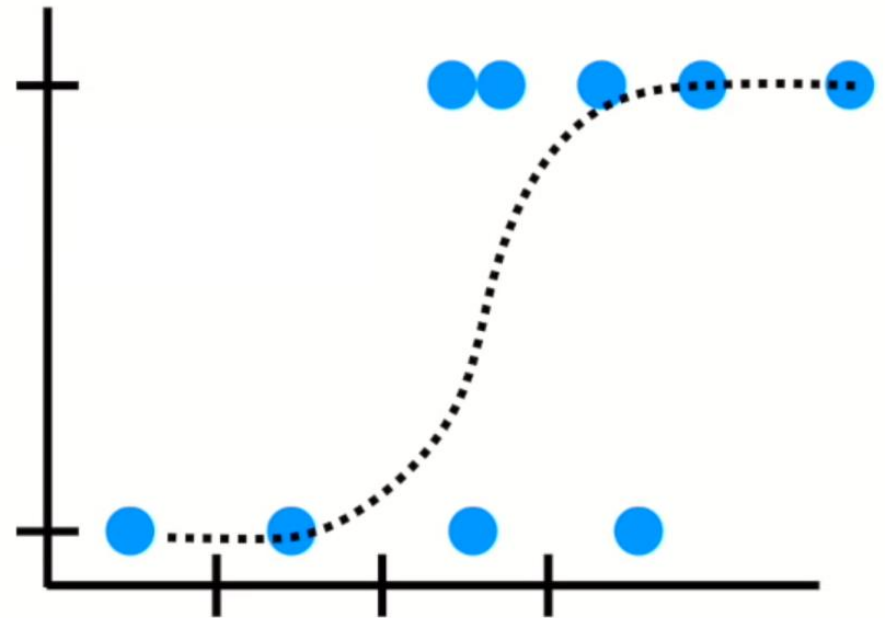
Don't worry **Cross Validation**  
comes the rescue !

It uses all the possible splits, and  
summarizes the results at the end

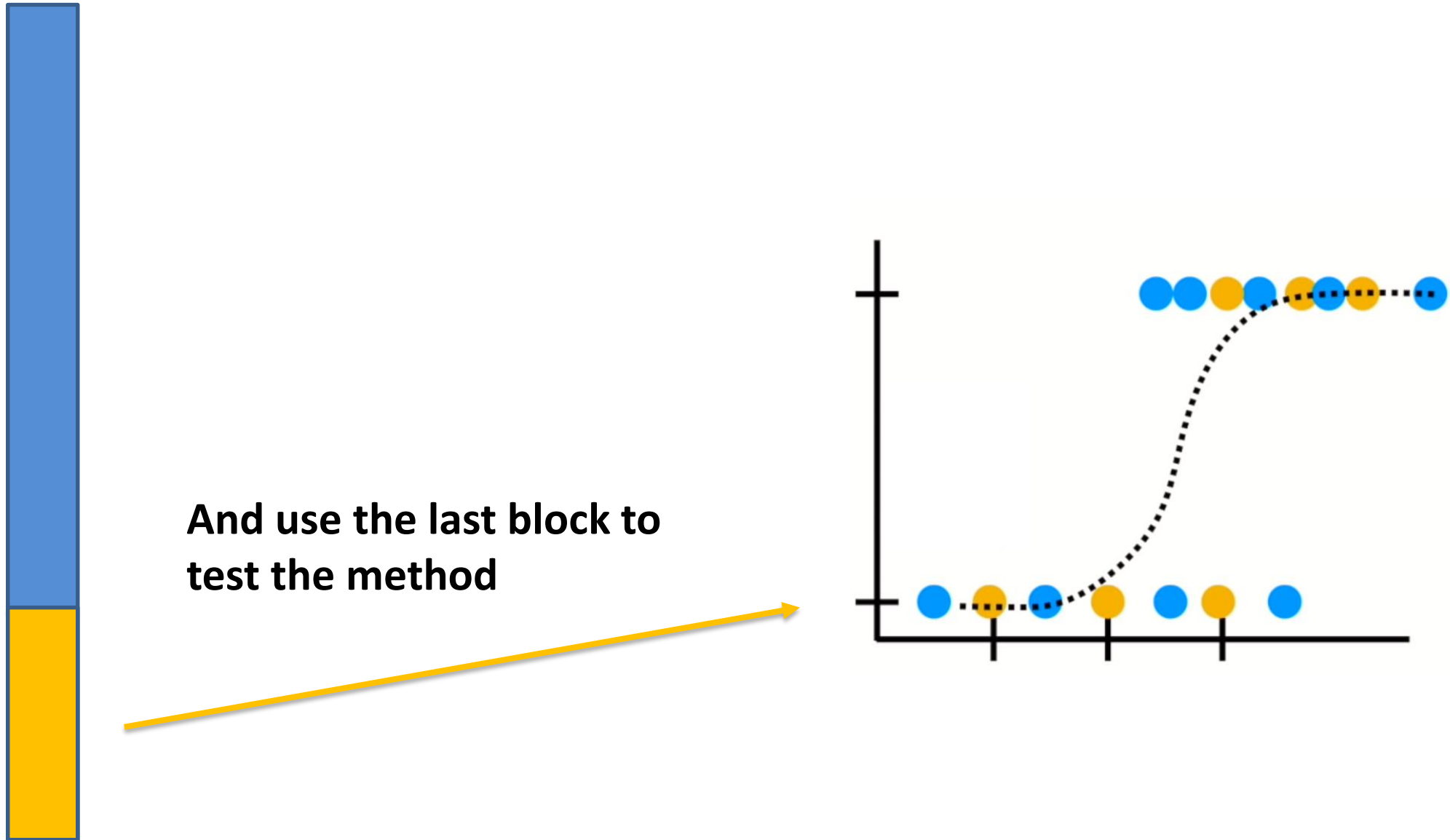


# Cross Validation

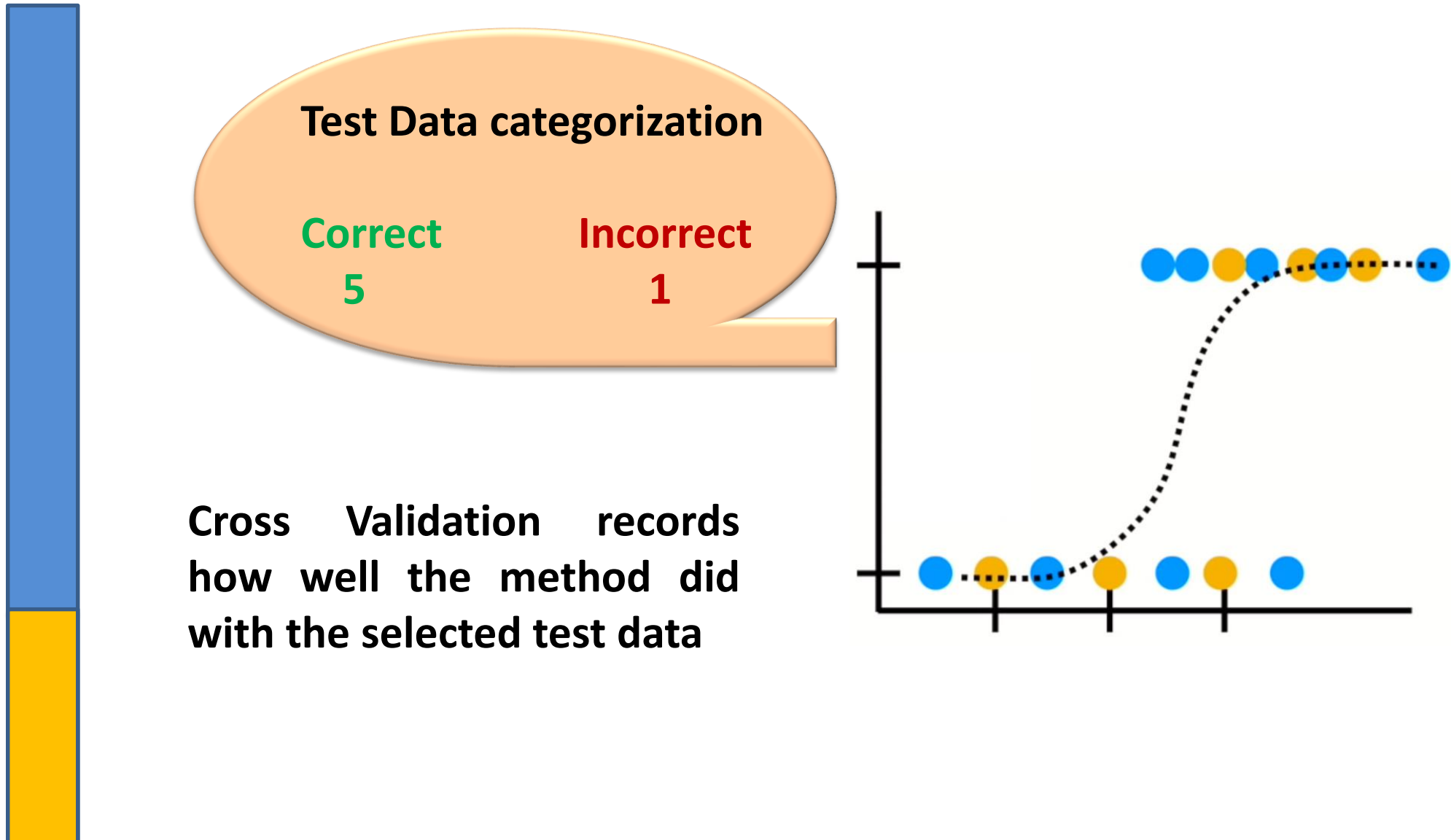
Start with the first 3 block  
to train the method



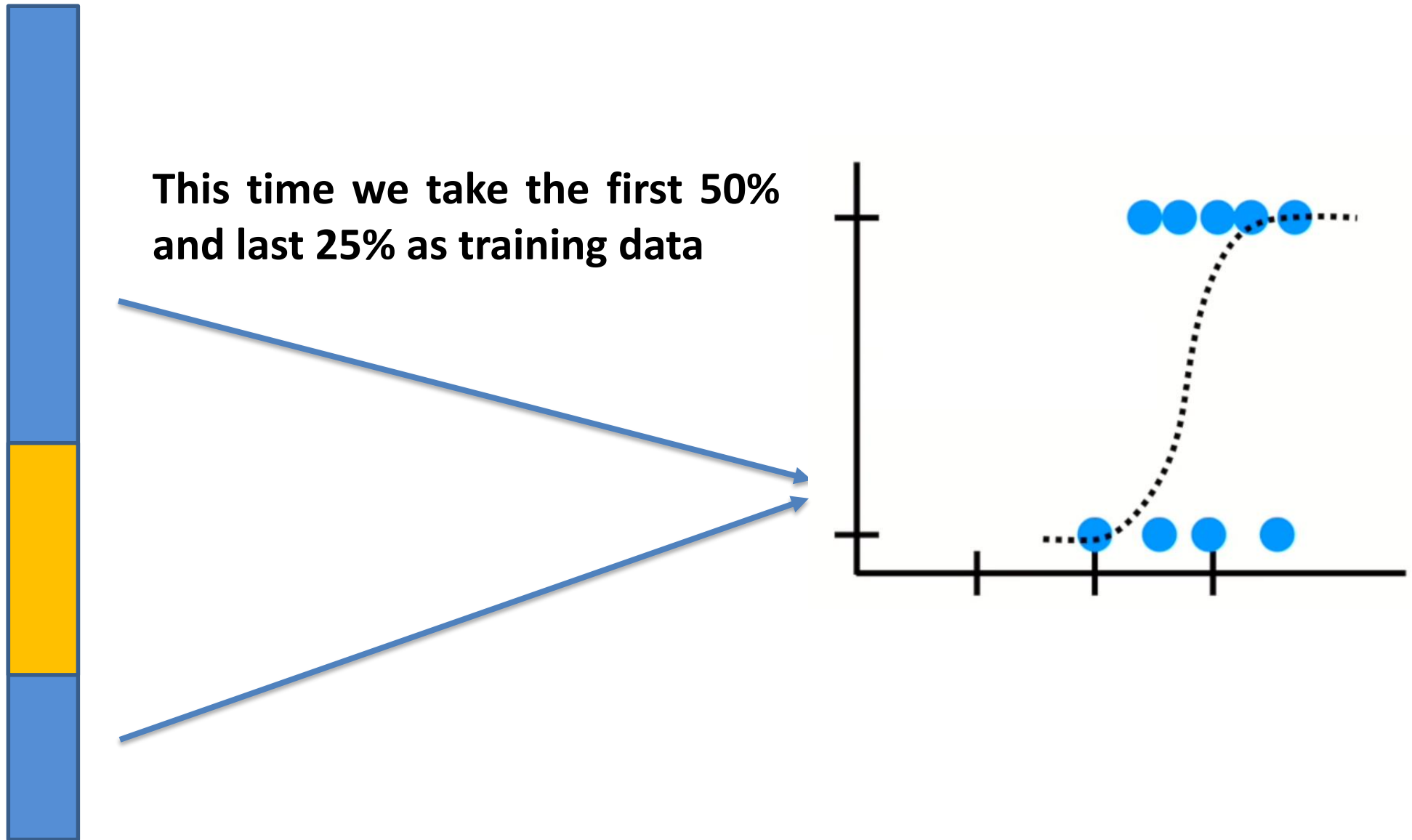
# Cross Validation



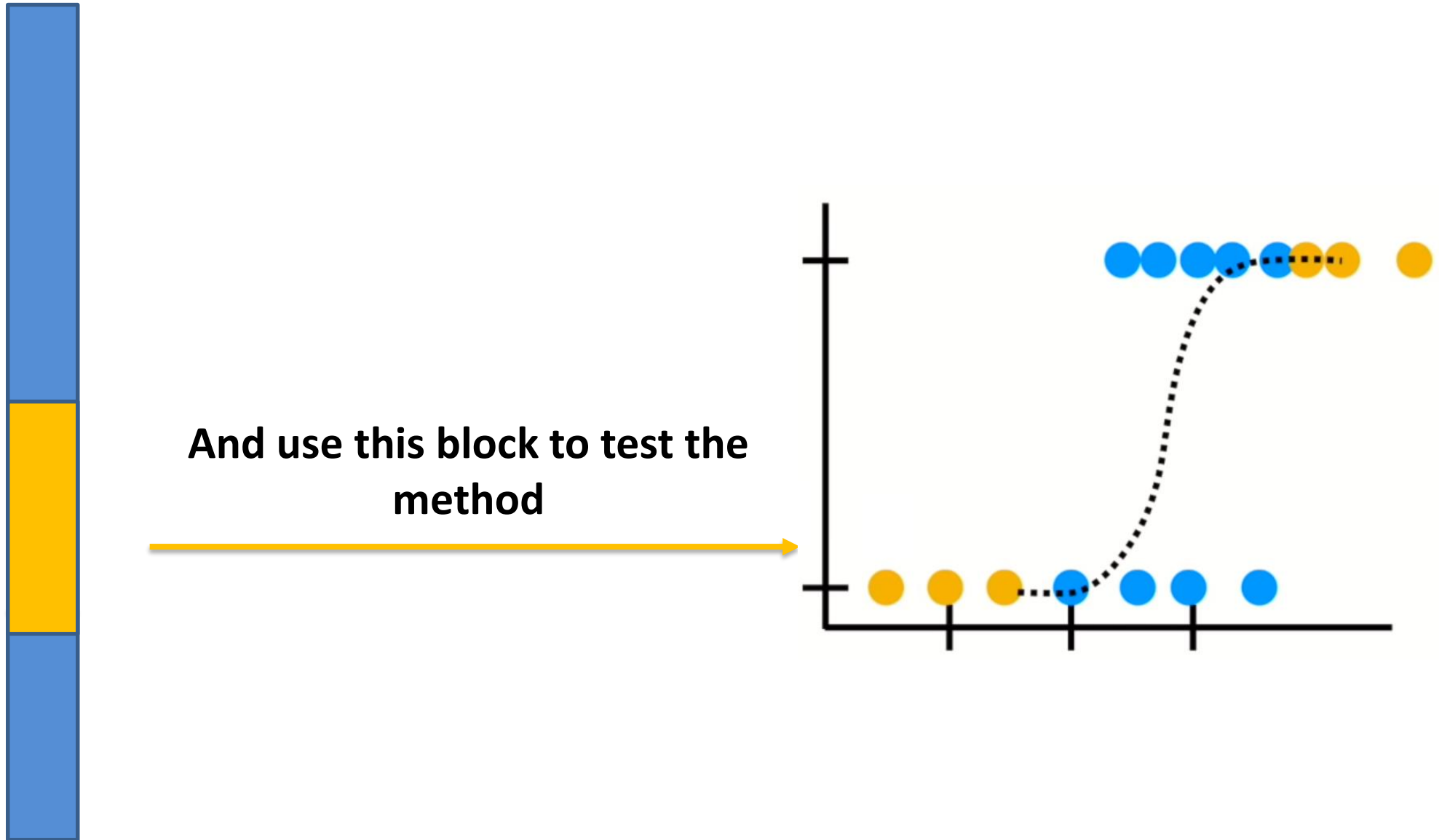
# Cross Validation



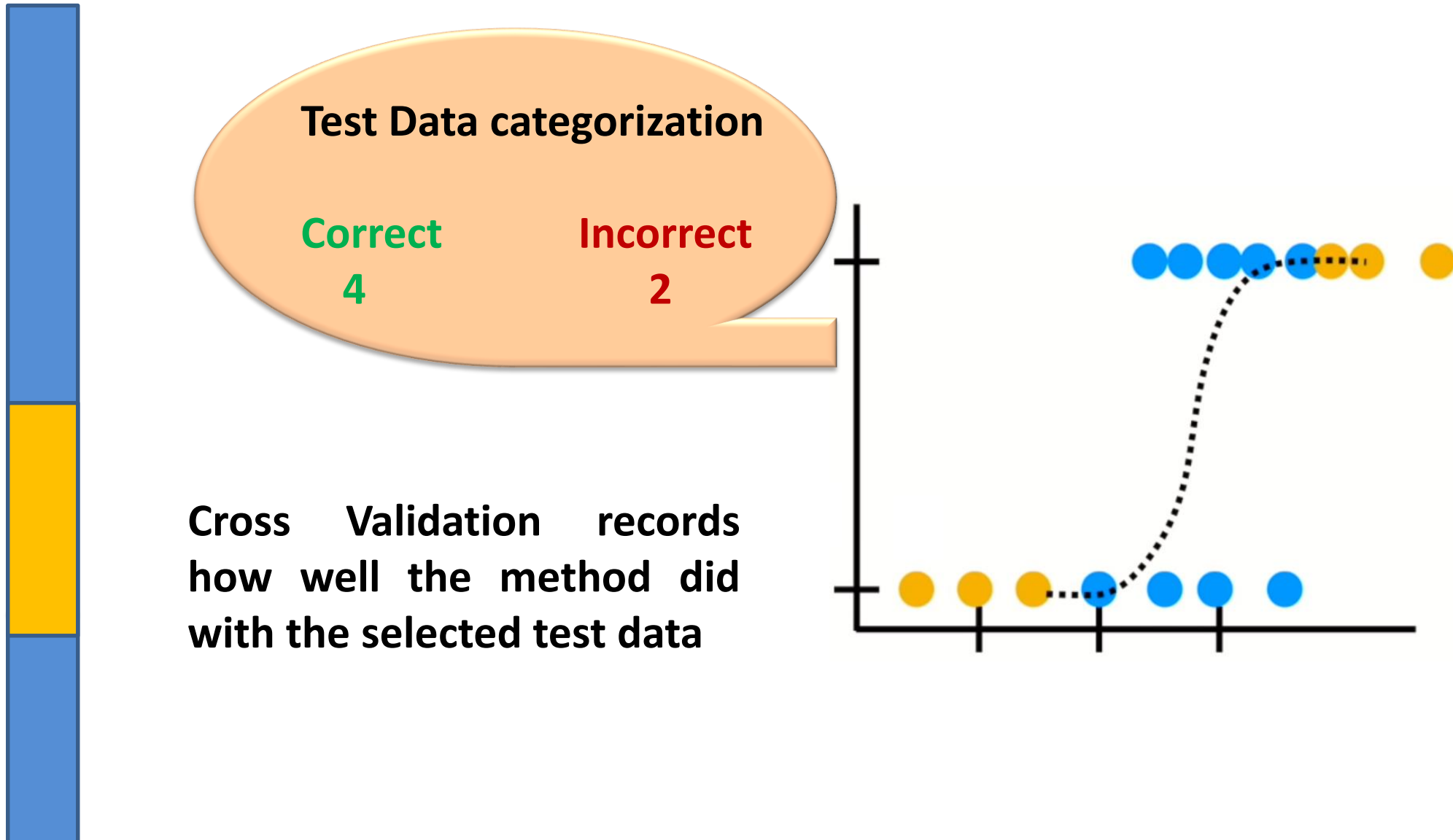
## Cross Validation



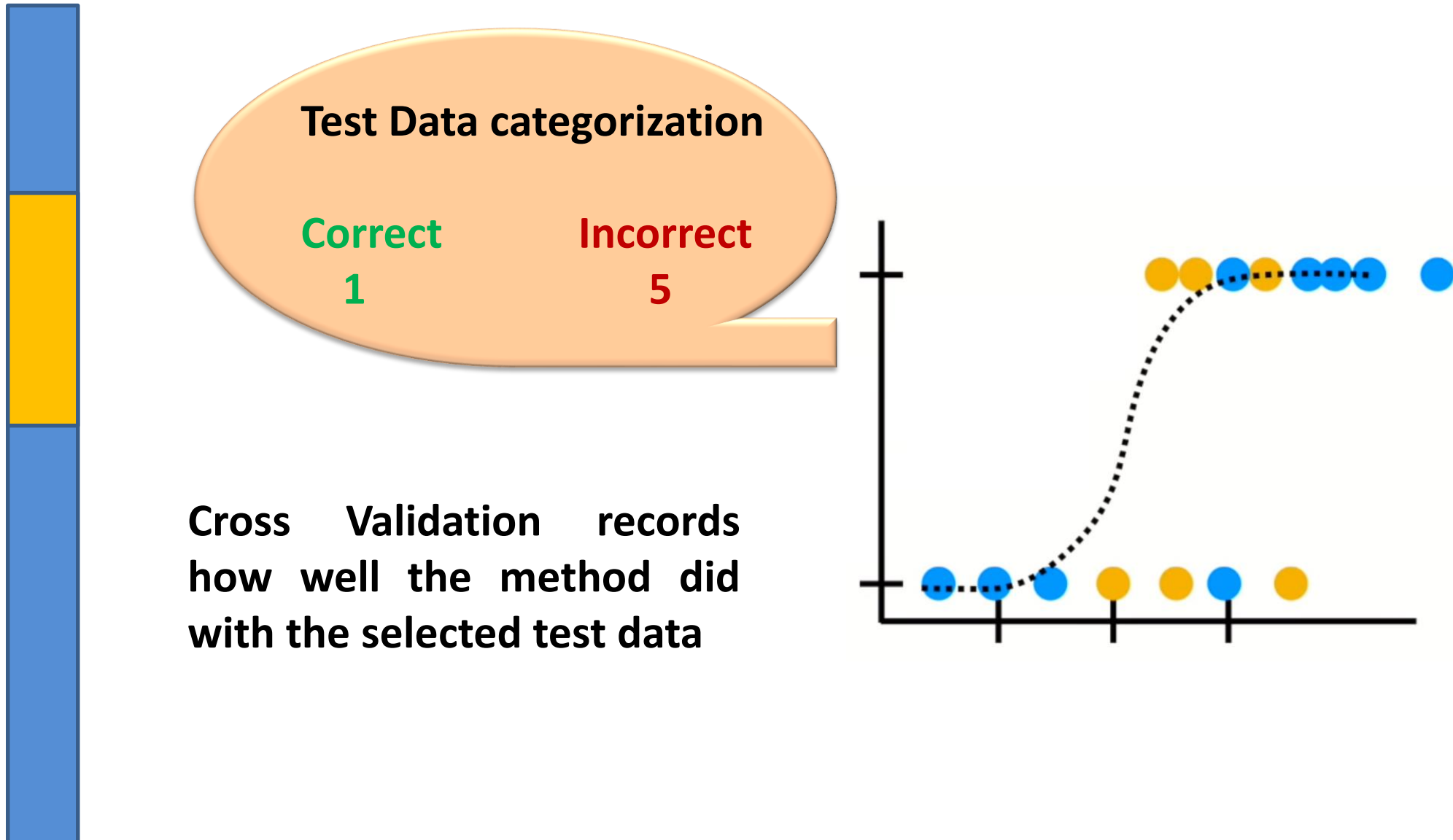
# Cross Validation



# Cross Validation

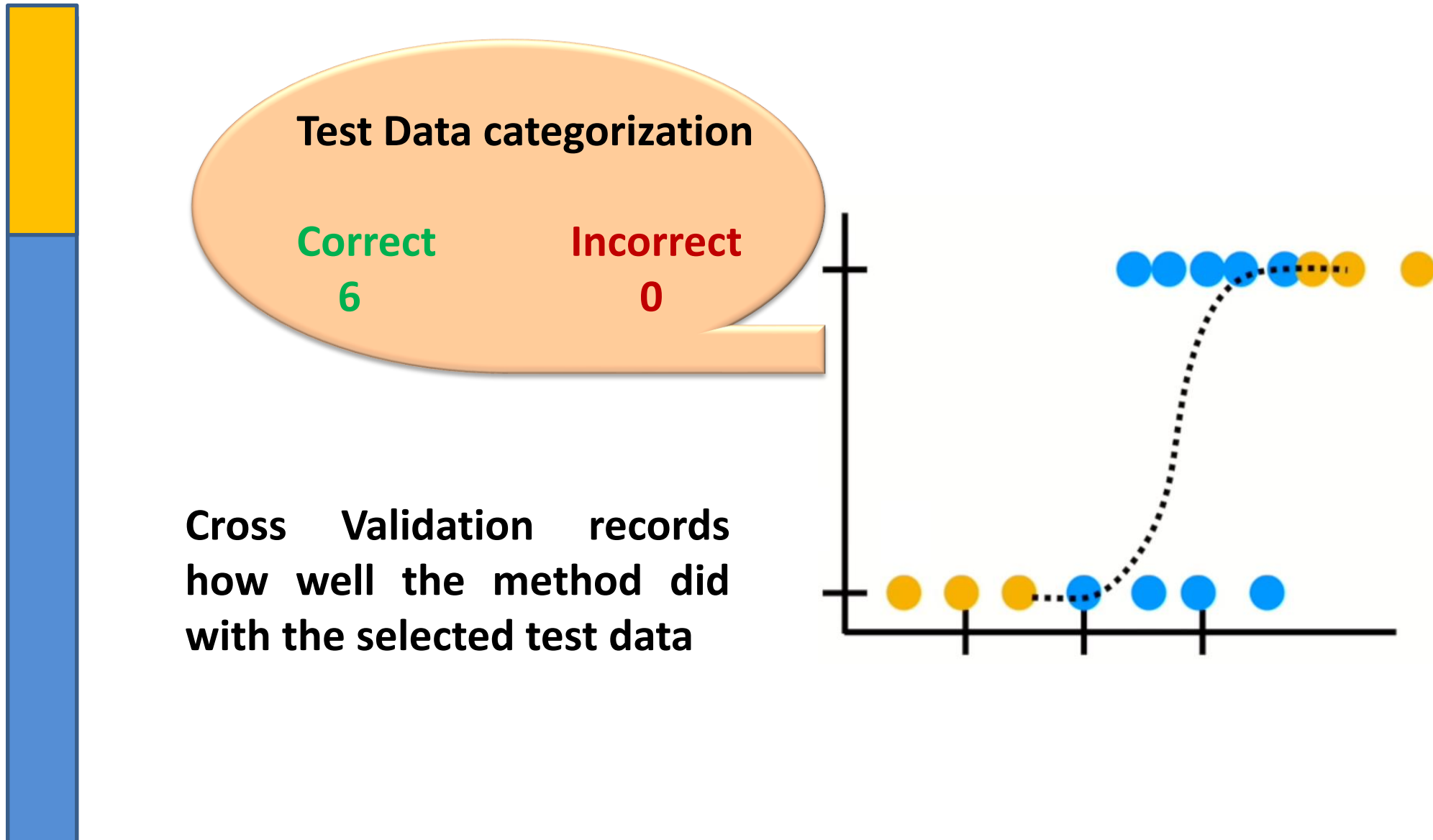


# Cross Validation



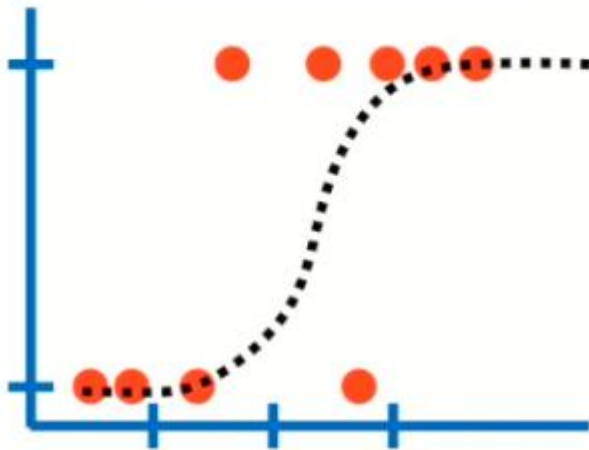


# Cross Validation



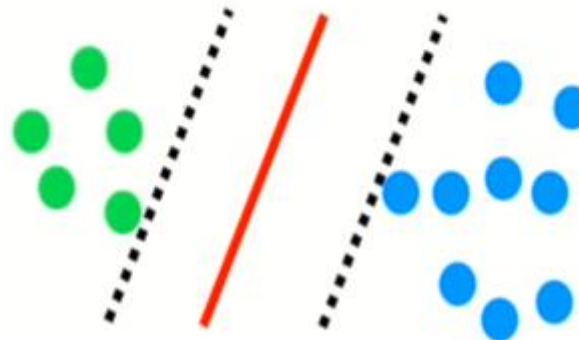
# Cross Validation

Logistic Regression



Correct	Incorrect
16	8

support vector machines (SVM)



Correct	Incorrect
18	6

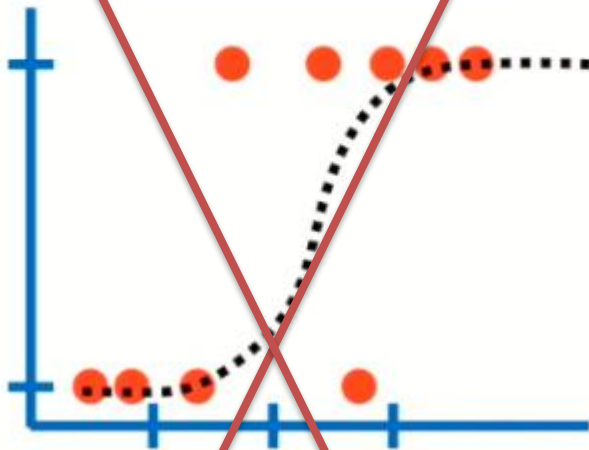
K-nearest neighbors



Correct	Incorrect
10	12

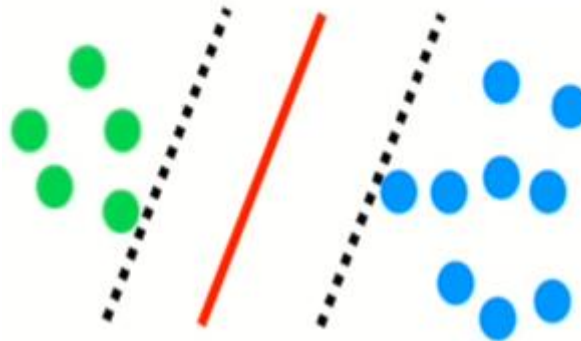
## Cross Validation

Logistic Regression



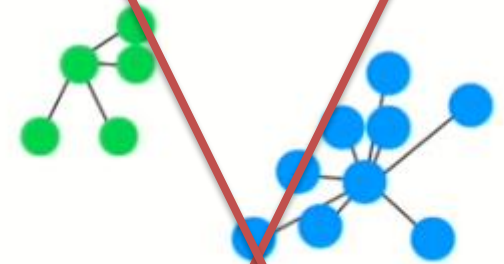
Correct	Incorrect
16	8

support vector machines (SVM)



Correct	Incorrect
18	6

K-nearest neighbors



Correct	Incorrect
10	12

# Cross Validation



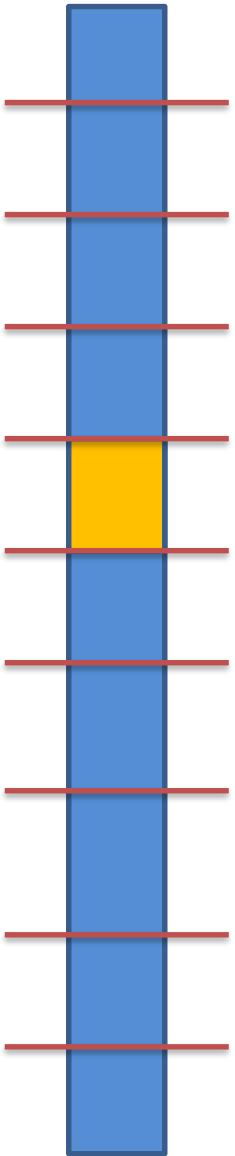
For this example we divided the data the data into 4 blocks. This is called **Four-Fold Cross Validation**

# Cross Validation



In an extreme special cases we could consider each sample as a block. This is called **Leave One Out Cross Validation**

# Cross Validation



It is very common to divide the data into 10 blocks.  
This is called **Ten-Fold Cross Validation**