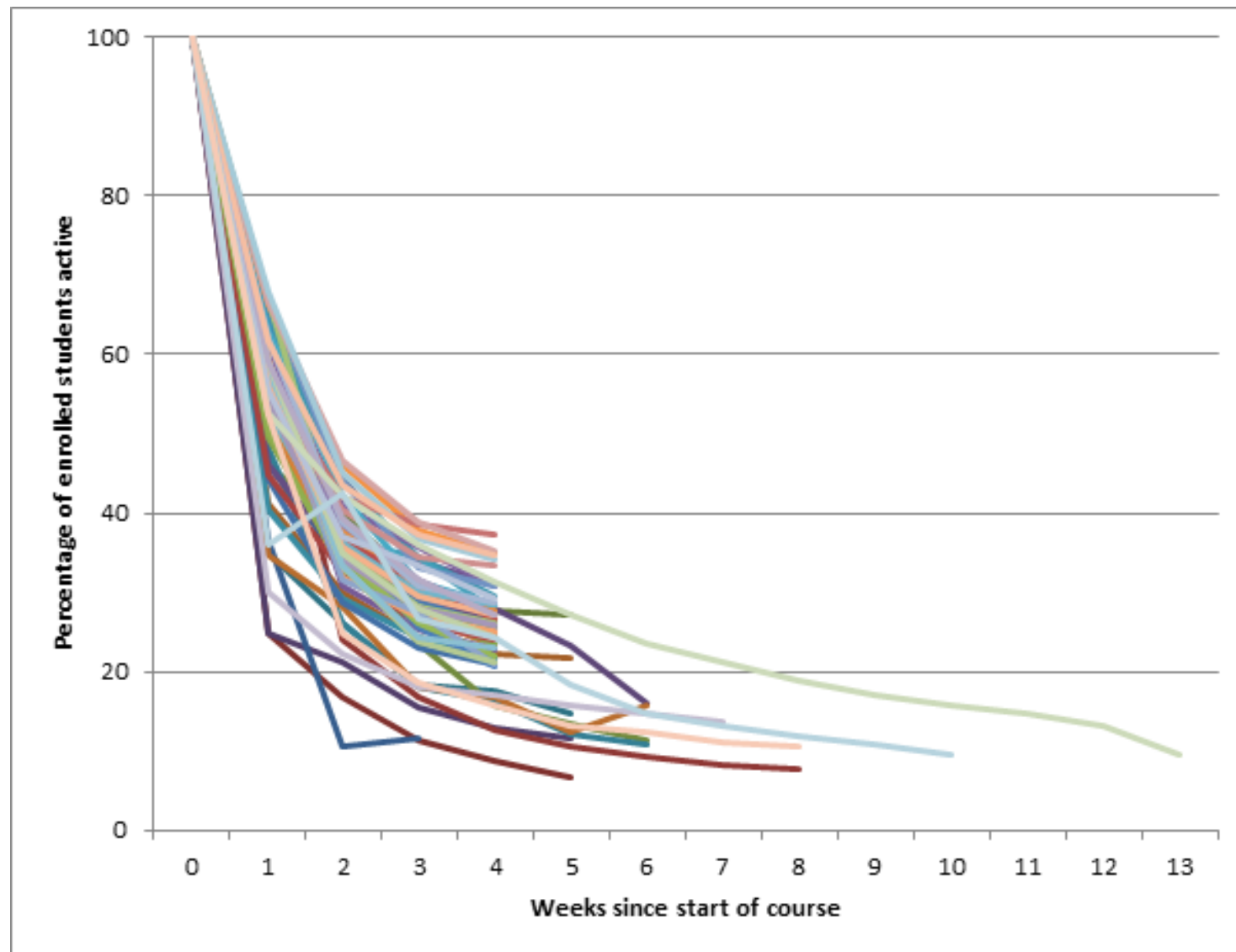# Prediction

- There are times when we want to automate a process in education

- Want to be able to predict what a student might do in the future: next question, next move in a game, drop out
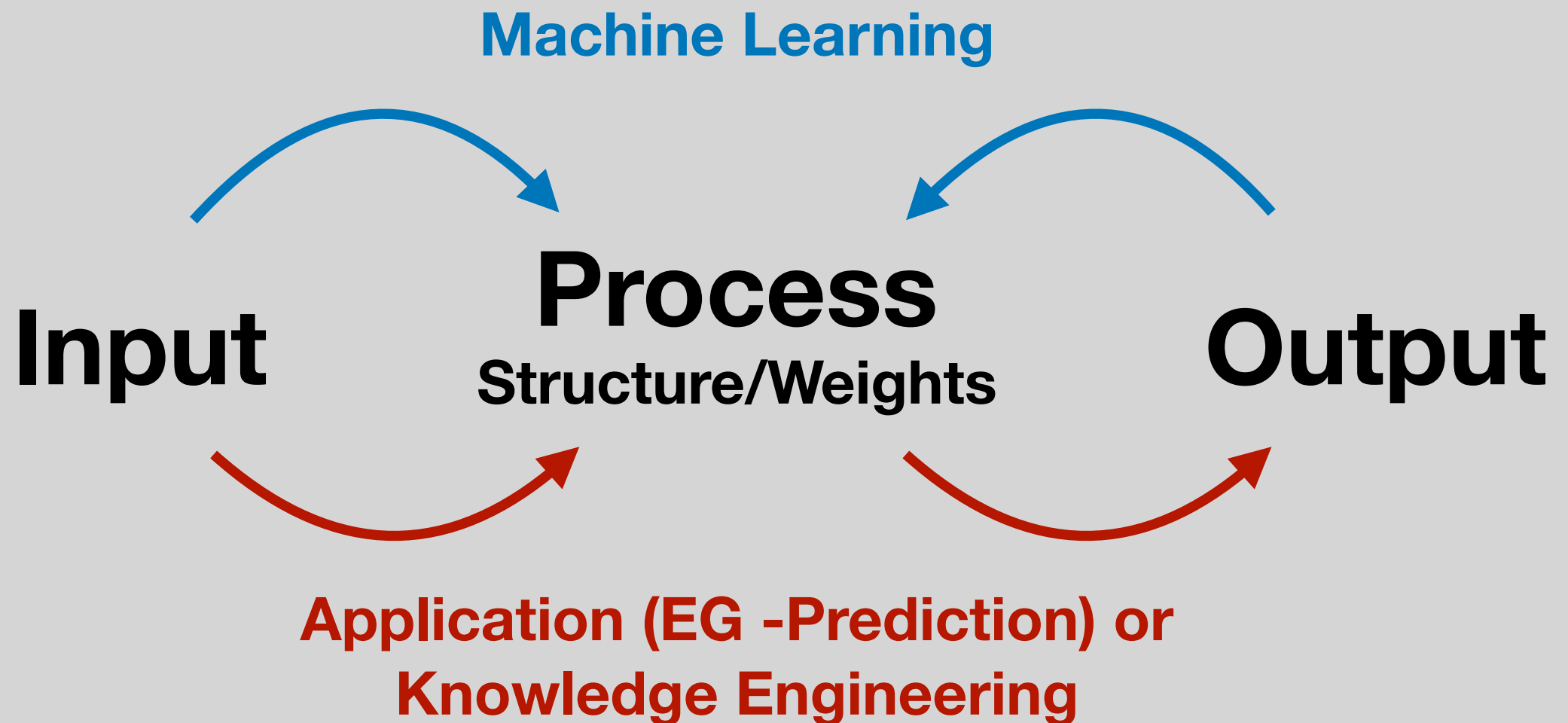
# Prediction



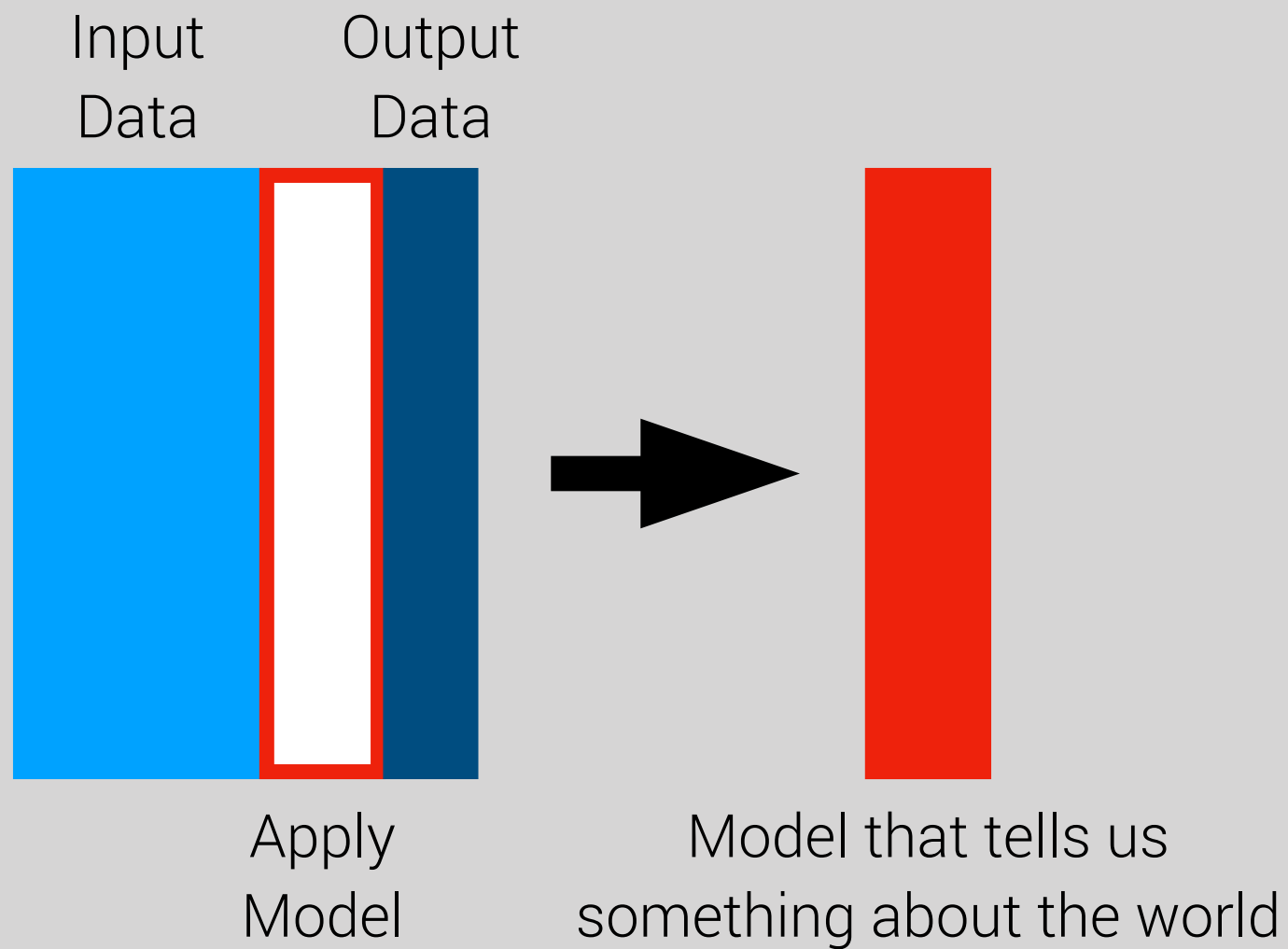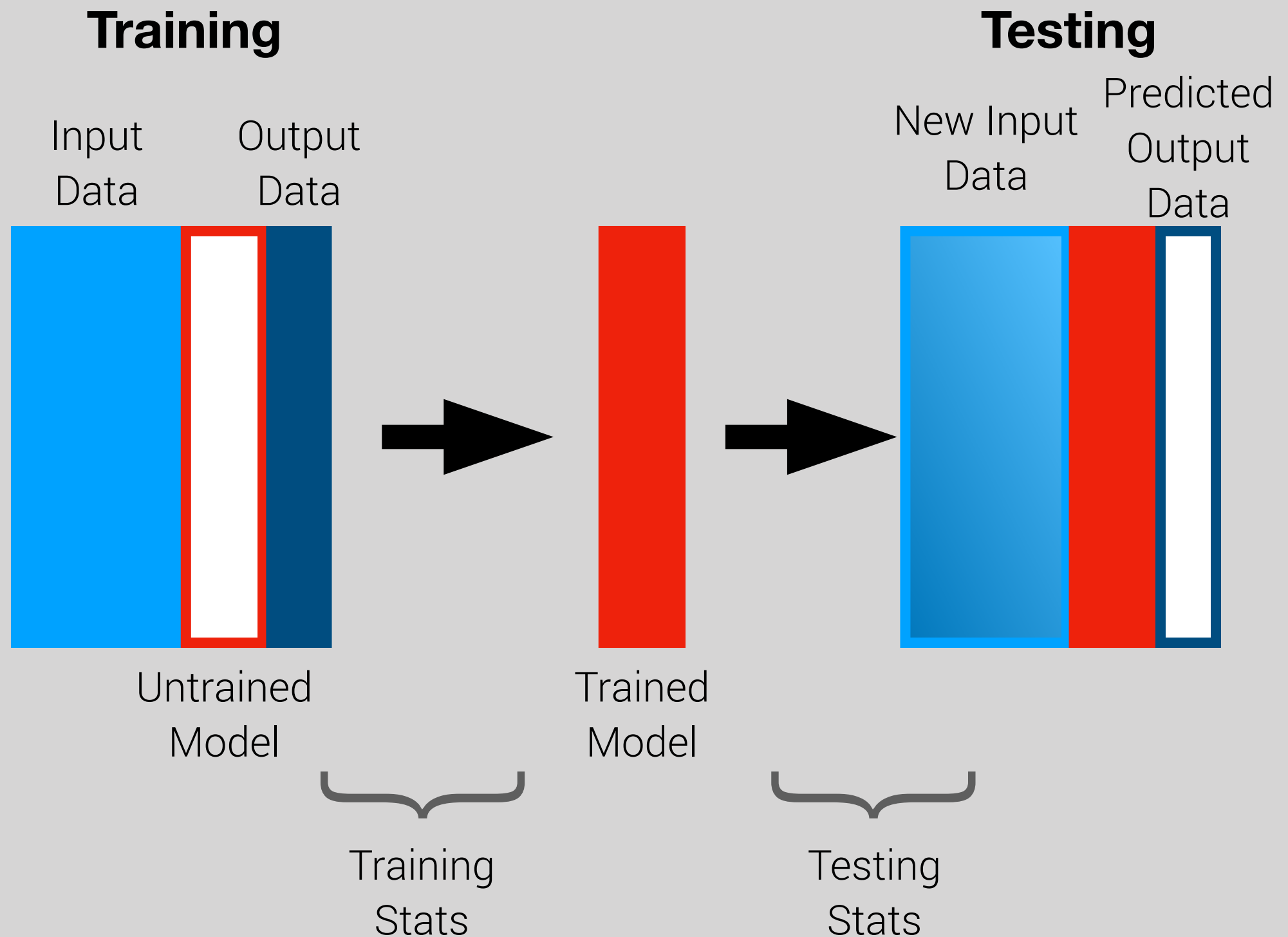K Jordan, Open University, 2013

# Machine Learning

**Machine Learning**

**Input**

**Process**
**Structure/Weights**

**Output**

**Application (EG -Prediction) or**
**Knowledge Engineering**

# Educational Statistics



Input Data

Output Data

Apply Model

Model that tells us something about the world

# Classification Confusion Matrix

|  |  | Actual Class (Observations) | |
|---|---|---|---|
|  |  | **P** | **N** |
| **Predicted Class** (Predictions) | **P** | TP | FP |
|  | **N** | FN | TN |

# Classification Confusion Matrix

|              |     | Actual Class |     |
| :----------: | :-: | :----------: | :-: |
|              |     | **P**        | **N** |
| **Predicted Class** | **P** | TP           | FP  |
|              | **N** | FN           | TN  |

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

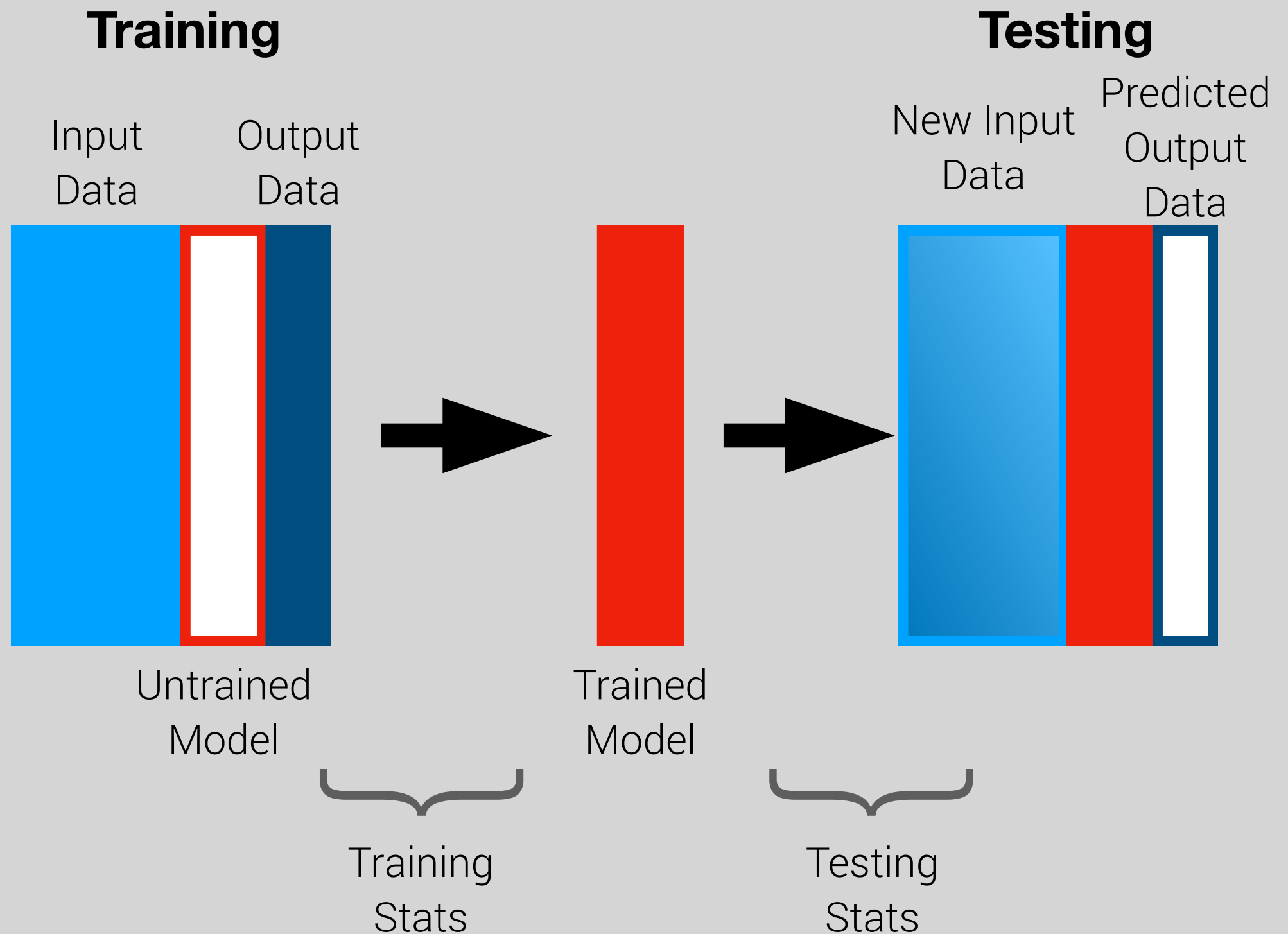$$\text{Sensitivity/Recall/TPR} = \frac{TP}{TP + FN}$$

$$\text{Specificity/Selectivity/TNR} = \frac{TN}{TN + FP}$$

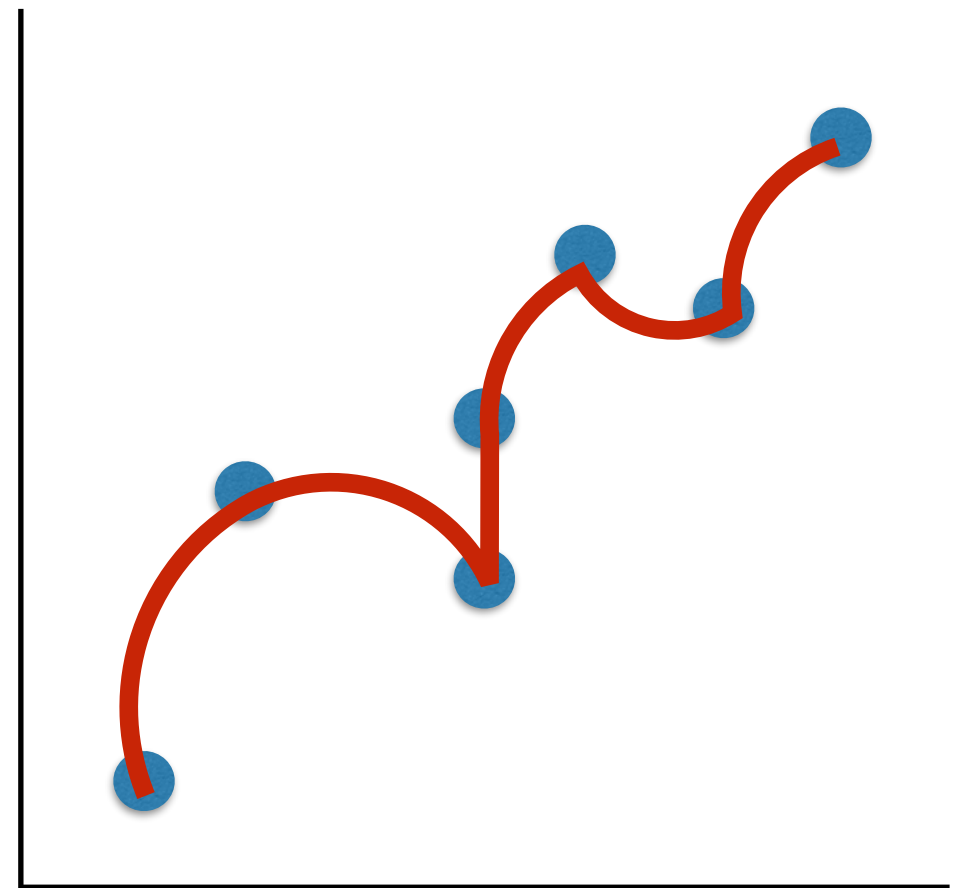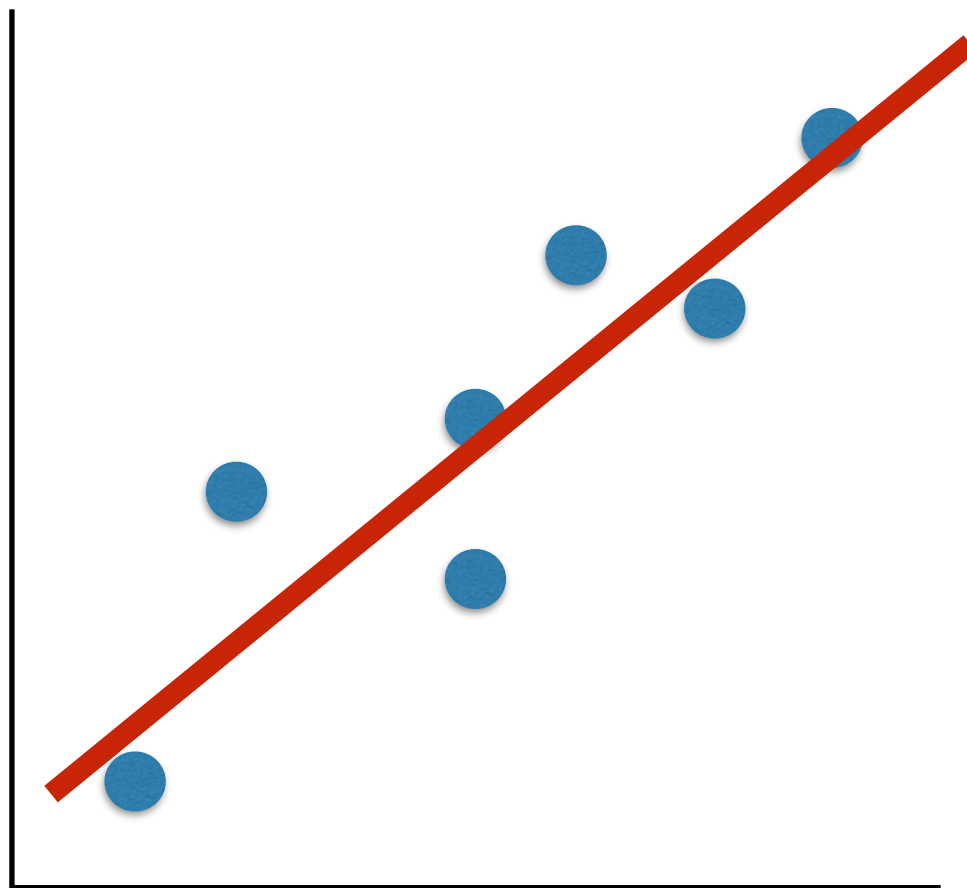$$\text{Precision/Positive Predictive Value (PPV)} = \frac{TP}{TP + FP}$$

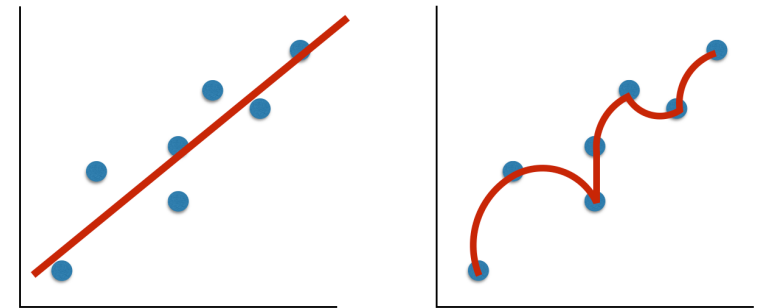$$F1 = \frac{2TP}{2TP + FP + FN}$$
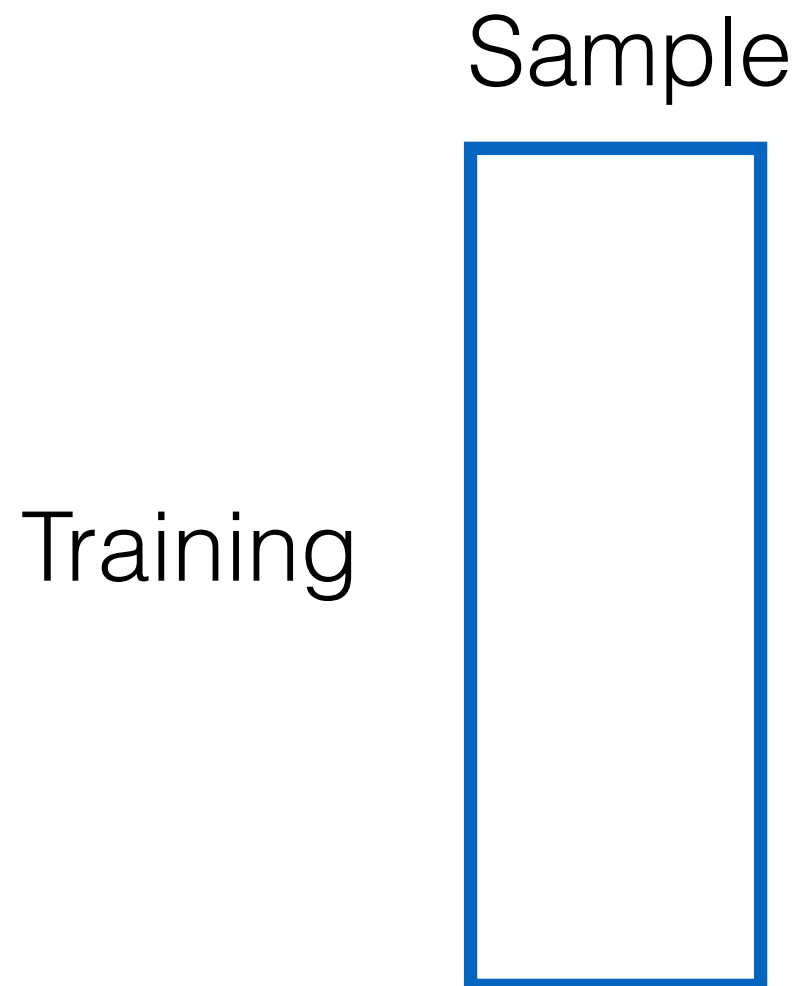
Which is more "accurate"?

Which is more "useful"?

How can we tell?
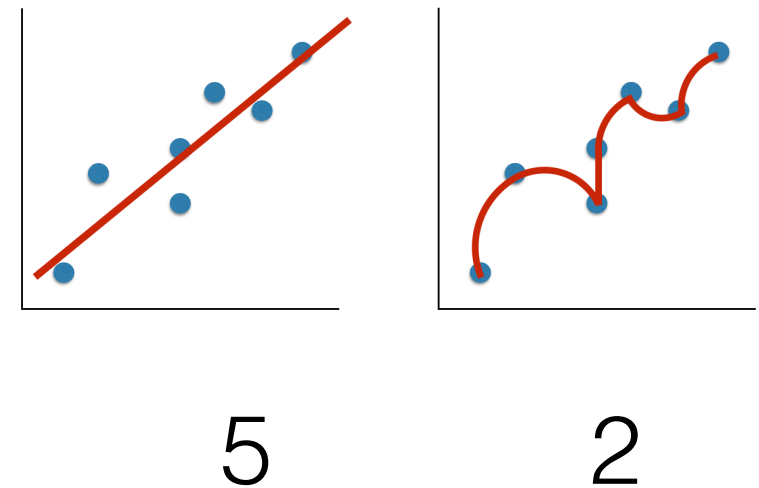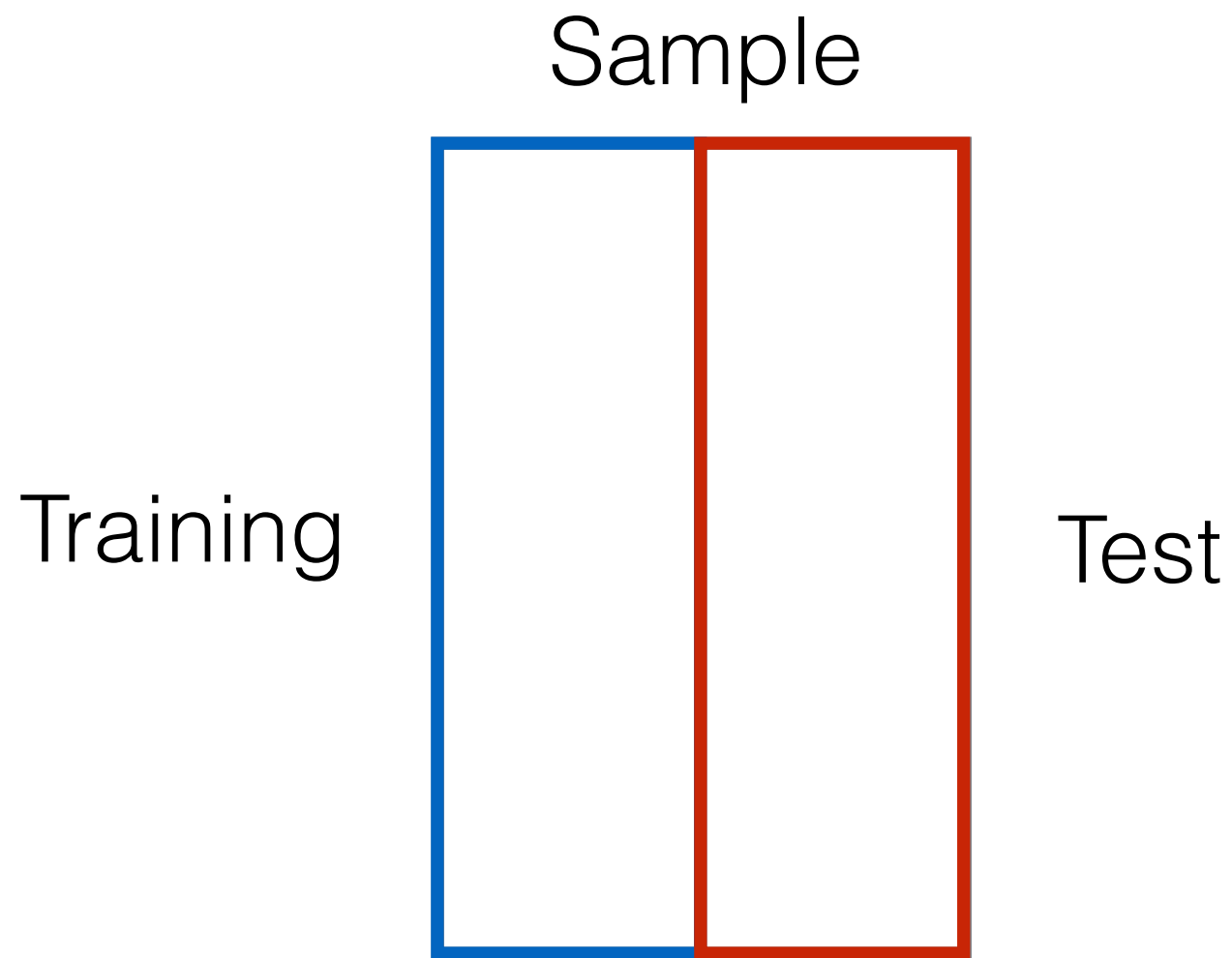
# Cross Validation

- Estimate how accurately a predictive model will perform in practice

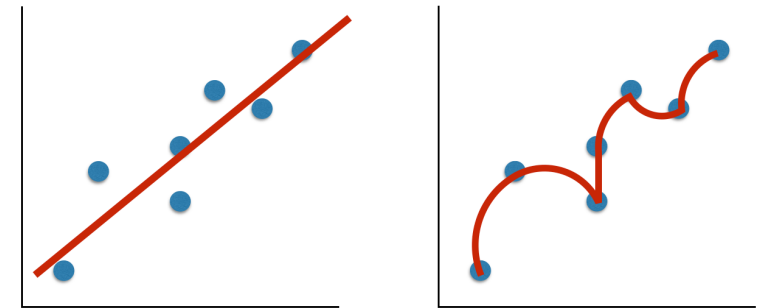- Give an insight on how the model will generalize to an independent dataset

# No Validation



**Problem**: Can't compare generalizability of models
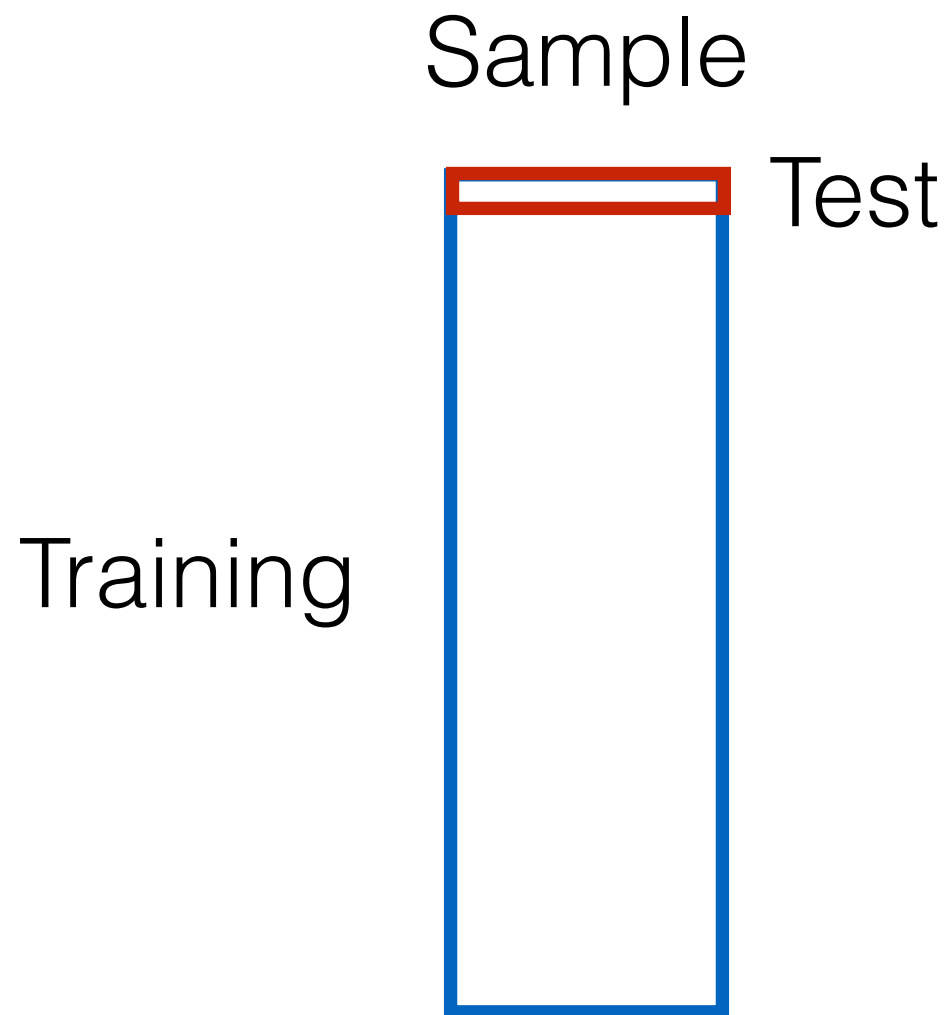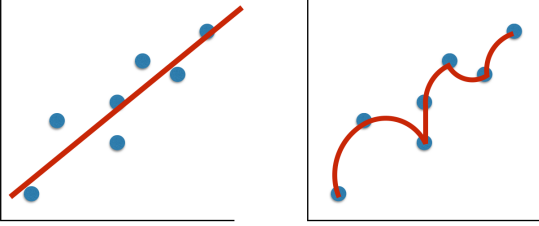
# Hold-out Validation

Sample

Training

Test

5                    2

**Problem**: very dependent on which data are in each group

# Hold One-out Validation

Sample

Test

Training

**Problem**: very computationally expensive

# K-Fold Cross Validation

Sample

| Test 1 | | | Training 1 | 5 | 2 |
| Test 2 | | | Training 2 | 4 | 2 |
| Test 3 | | | Training 3 | 3 | 1 |
| Test 4 | | | Training 4 | 5 | 4 |
| Test 5 | | | Training 5 | 4 | 2 |
| | | | | 4.2 | 2.2 |

Calculate how accurate we are in each "fold"
and average the answer

# Five Tribes

- Symbolists
- Connectionists
- Evolutionaries
- Bayesians
- Analogizers

# Classification Tree

- Decision tree

- Map observations (branches) onto classes (leaves)
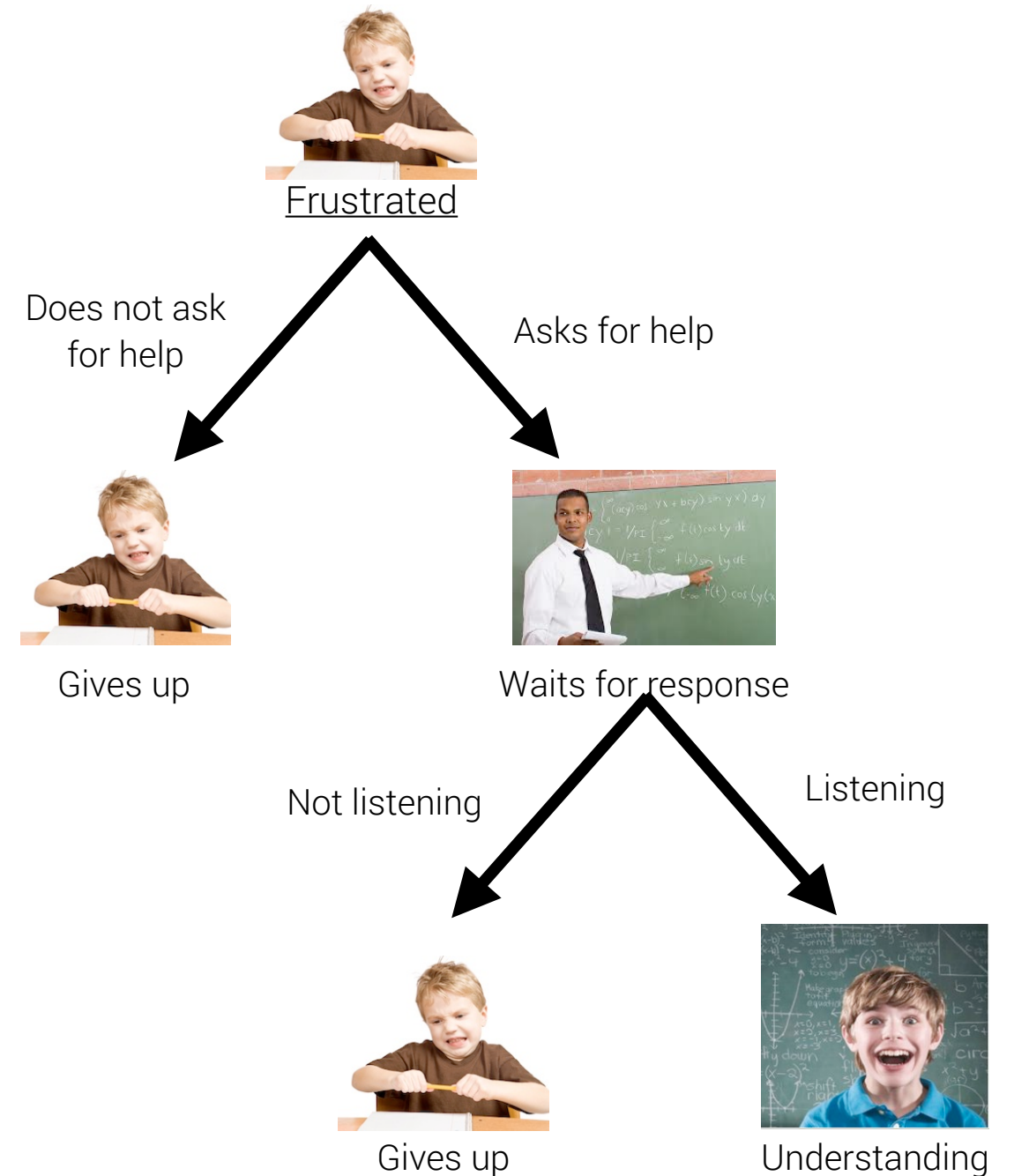
- Tree describes the data but can be used as classification

- EG: student states = leaves, student actions = branches



Frustrated

Does not ask for help

Asks for help

Gives up

Waits for response

Not listening

Listening

Gives up

Understanding

# Machine Learning

## Input

## Process
### Structure/Weights
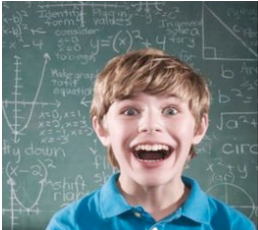
## Output

Does not ask
for help

Asks for help

Not listening

Listening


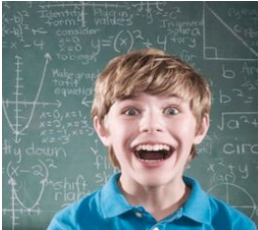Frustrated

Does not ask
for help

Asks for help


Gives up


Waits for response

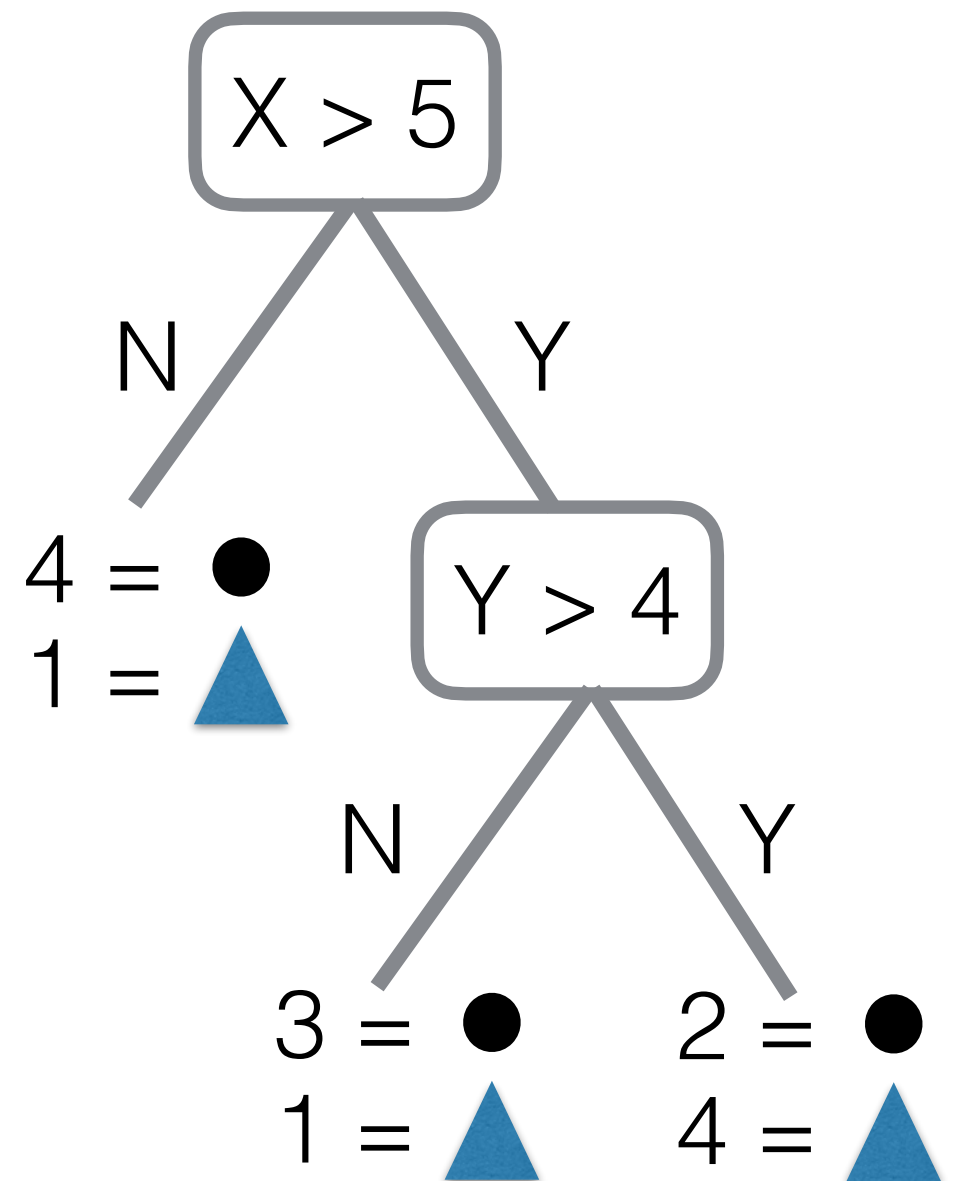Not listening
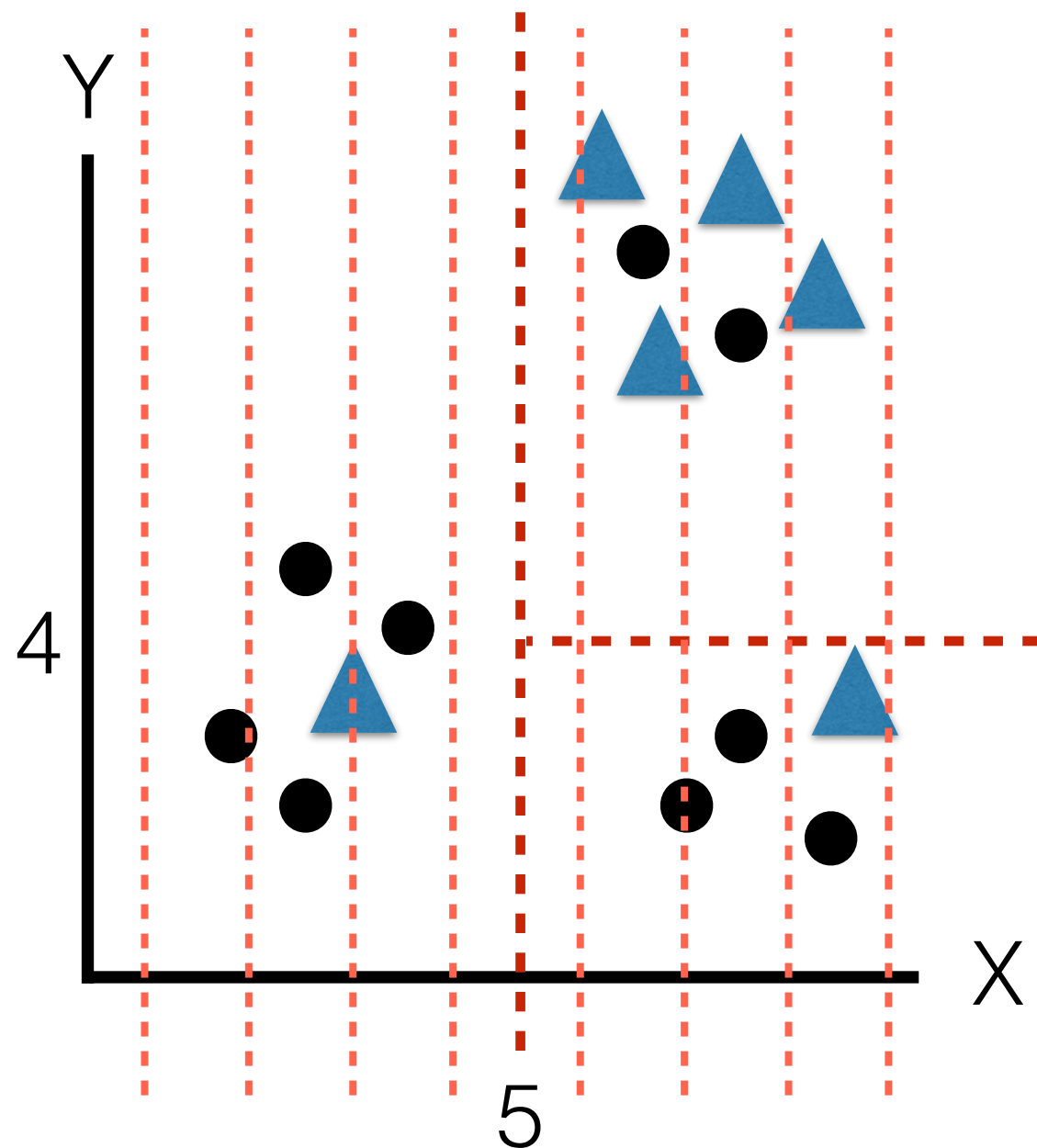
Listening


Gives up


Understanding


Gives up


Understanding

# Binary Classification Tree

* Minimize the error

# caret

- Standard syntax for comparing many models

- Generate training and testing data sets

- Run several model types

- Run resampling algorithms and alter parameters to generate the best model

- Compare using the same diagnostic metrics

- https://topepo.github.io/caret/

# caret

## Generate Training/Test Data Sets

```
trainData <- createDataPartition(
  y = data$thing, ## the outcome data are needed
  p = .75, ## The percentage of data in the
training set
  list = FALSE)

#Generates a list of index numbers for the sample

training <- DATA[ trainData,]
testing  <-DATA[-trainData,]
```

# caret

K-Fold Cross Validation

```
ctrl <- trainControl(method = "cv", repeats = 3)
```

# caret

## Train Model

```
fit1 <- train(
  thing ~ .,
  data = training,
  method = "model",
  preProc = c("center", "scale")## Center and scale
the predictors for the training set and all future
samples.
  trControl = ctrl #add cross validation specs
  metric = "cp"
)
```

# caret

## Test Model

```
pred1 <- predict(fit1, newdata = testing)

confusionMatrix(data = pred1, DATA$thing)
```

# Project

Train and test three tree-based models (CART, Conditional Inference Trees and C50) using data from the University of Michigan Open Data Set.