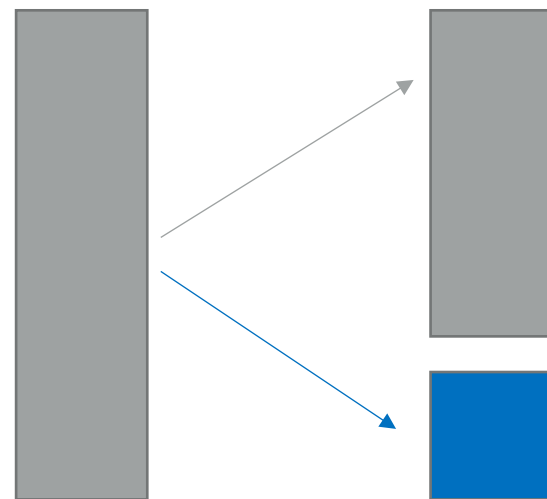# Model Basics
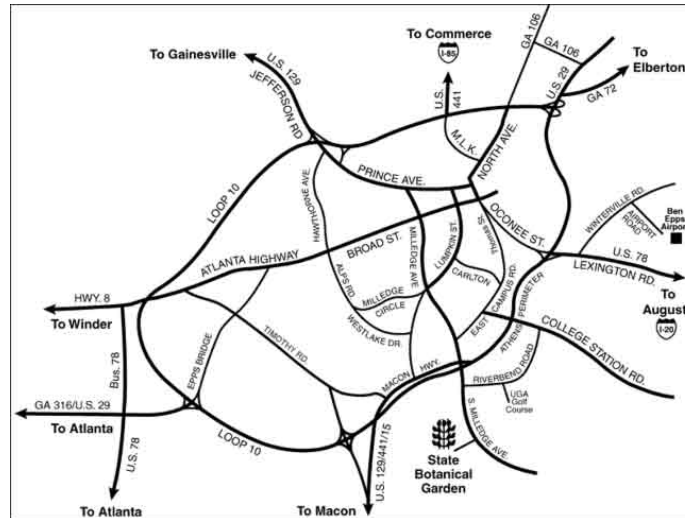
John Rios

*Business Intelligence and Analytics*
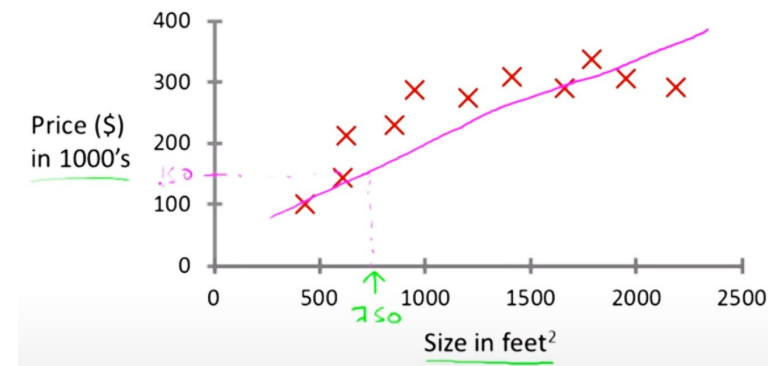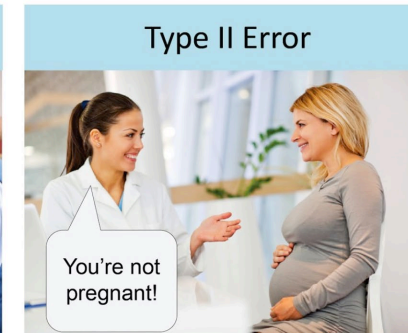
Type I Error — You're pregnant!

Type II Error — You're not pregnant!

Price ($) in 1000's vs Size in feet²

Data

Training Data

Testing Data

| | Attributes | | | Target attribute |
|---|---|---|---|---|
| Name | Balance | Age | Employed | Write-off |
| Mike | $200,000 | 42 | no | yes |
| Mary | $35,000 | 33 | yes | no |
| Claudio | $115,000 | 40 | no | no |
| Robert | $29,000 | 23 | yes | yes |
| Dora | $72,000 | 31 | no | no |

This is one row (example).
Feature vector is: **<Claudio,115000,40,no>**
Class label (value of Target attribute) is **no**

# Where We Are



✓ Connect Tableau to various Datasets: Excel and CSV files
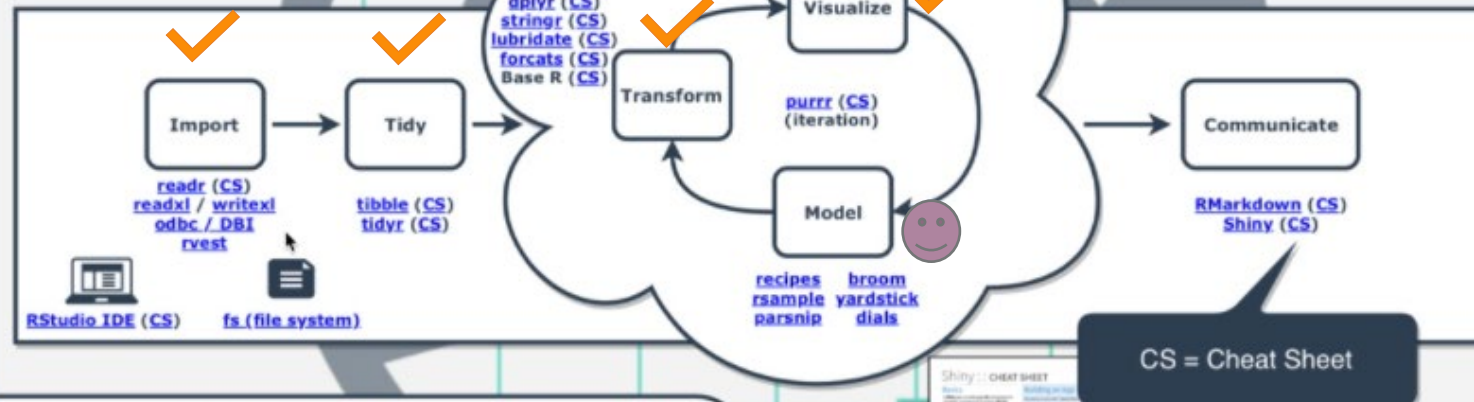
**Import**

✓ Create Area Charts   ✓ Create Barcharts
✓ Create Scatterplots   ✓ Create Maps
✓ Create Treemaps
✓ Create Storylines
✓ Create Interactive Dashboards

ggplot2

**Explore**

```
> # how did we do with test set? confusion matrix
> confusionMatrix(data = churn_test$pred_churn,
+                 reference = churn_test$Churn,
+                 mode = "prec_recall",
+                 positive = "Yes")
Confusion Matrix and Statistics

          Reference
Prediction  No Yes
       No  906 186
       Yes 126 187

               Accuracy : 0.7779
                 95% CI : (0.7553, 0.7994)
    No Information Rate : 0.7345
    P-Value [Acc > NIR] : 9.989e-05

                  Kappa : 0.3998

 Mcnemar's Test P-Value : 0.0008371

              Precision : 0.5974
                 Recall : 0.5013
                     F1 : 0.5452
             Prevalence : 0.2655
         Detection Rate : 0.1331
   Detection Prevalence : 0.2228
      Balanced Accuracy : 0.6896

       'Positive' Class : Yes
```

Root Mean Squared Error (RMSE)

**Evaluate**

**Question**

left_join()
right_join()
inner_join()
full_join()
semi_join()
anti_join()

**Wrangle**

✓ Understand Types of Joins and how they work

CLASSICAL MACHINE LEARNING

SUPERVISED   UNSUPERVISED
CLASSIFICATION   CLUSTERING
REGRESSION   ASSOCIATION
DIMENSION REDUCTION (generalization)

```
case_when(condition-1 ~ output-1    # IF
         ,condition-2 ~ output-2    # ELIF
         ,condition-3 ~ output-3    # ELIF
         ,TRUE ~ output-else        # ELSE
         )
```

**Model**

actionable

**Test**

All customers → Predicted as likely to **churn**

A. Run retention campaign
A/B Test → Any significant difference in churn levels afterwards?
B. Business as usual

A   B
CONTROL   VARIATION

Deployment

| Database | | Statistics | | Machine Learning |

# Model Defined

A **simplified\* representation** of reality created for a **specific purpose**

- *\*based on some assumptions** about what is and is not important, or sometimes based on constraints on information or tractability

# **Model Goal**

Not to uncover truth, but to discover a simple approximation that is still useful

- **i.e., capture true "signals"** (or patterns generated by the phenomenon of interest) **and ignore "noise"** (or random variation that you're not interested in)
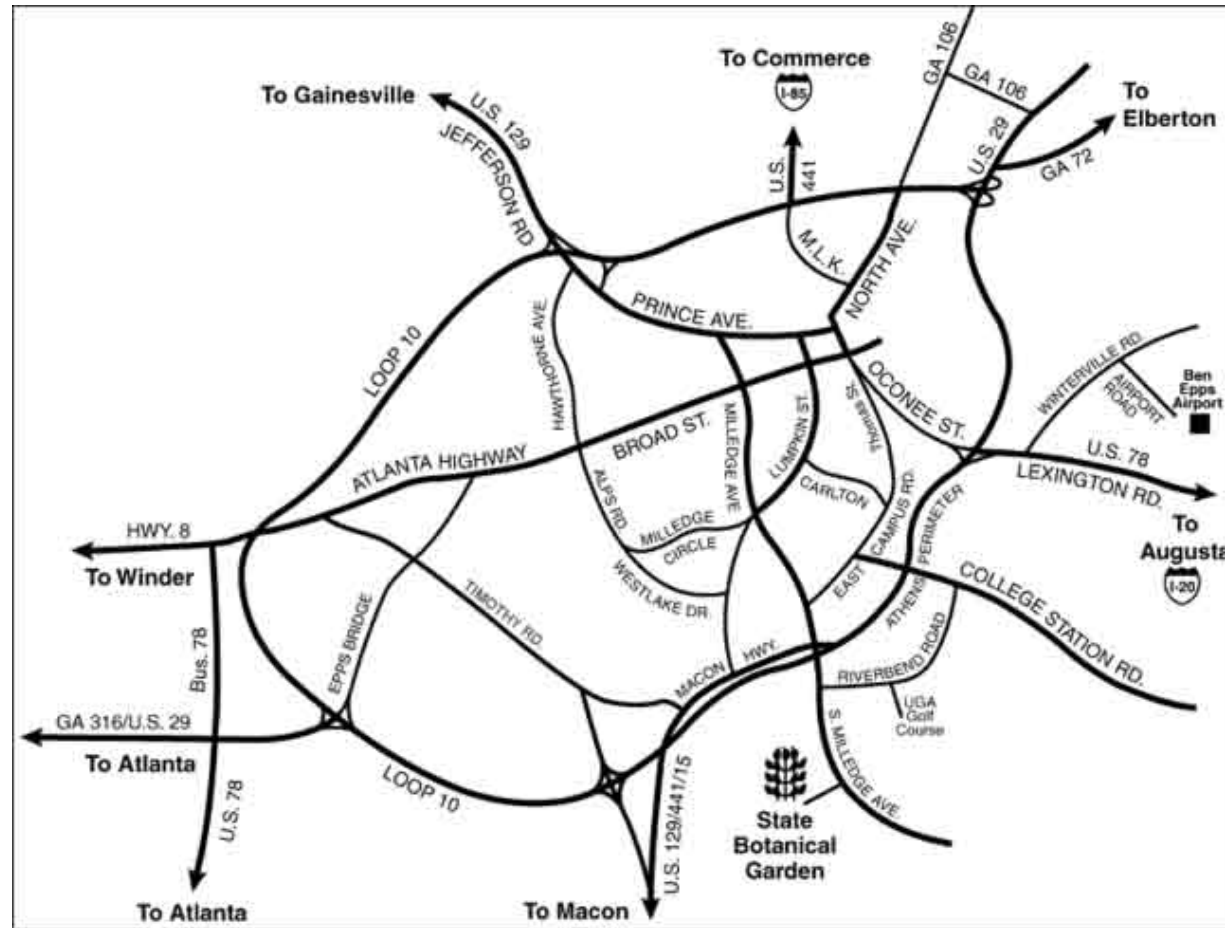
# *"All models are wrong, but some are useful"*

George Box

# Model Example

# Predictive (or Supervised Learning) Model

A formula for estimating the unknown value of interest:

- **The target!**

  - The formula could be mathematical or a logical statement, such as a rule. Often, it is a hybrid

# More Terminology

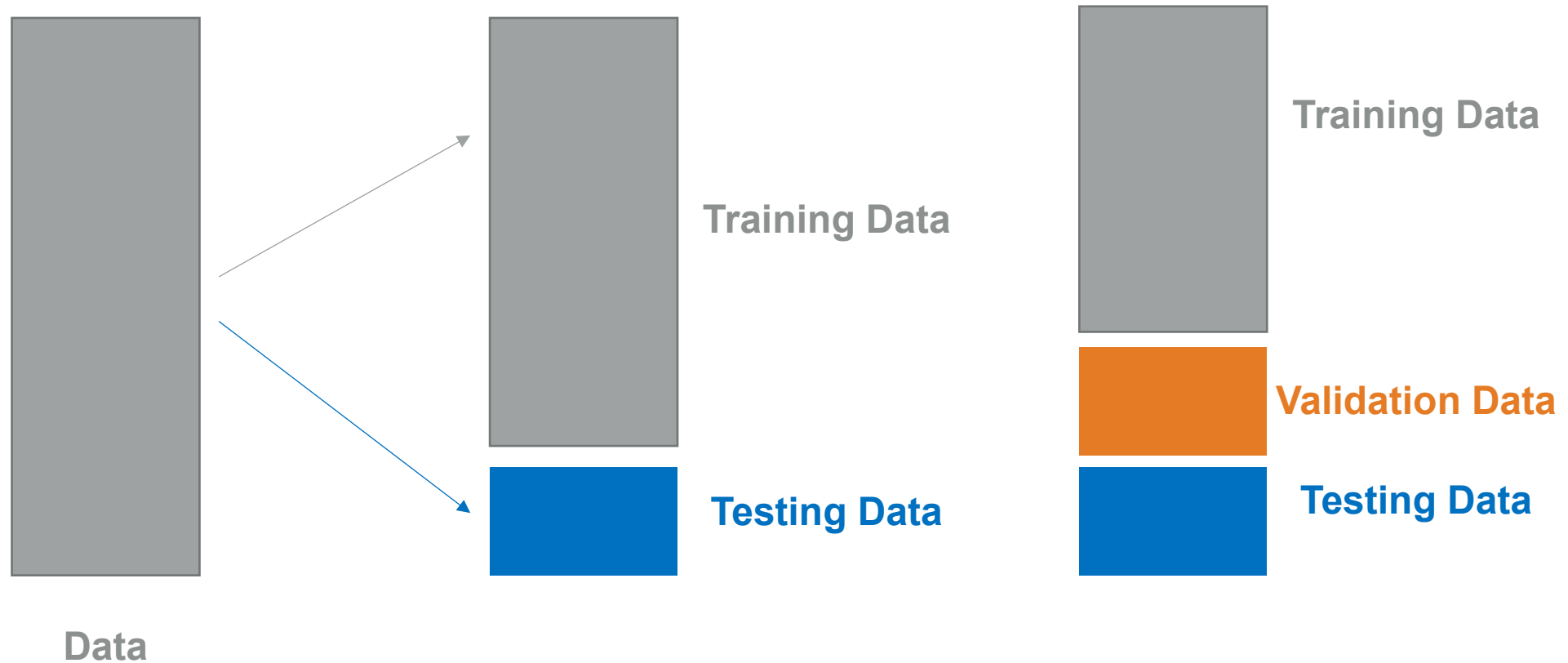**Instance / example** = a fact or data point described by a set of attributes (also known as variables, columns, or features)

**Model induction** = the creation of models from data

**Training data** = the input data used for model induction

**Testing data** = the input data used for model testing

# Predictive Model: Data Splitting



Data

Training Data

Testing Data

Training Data
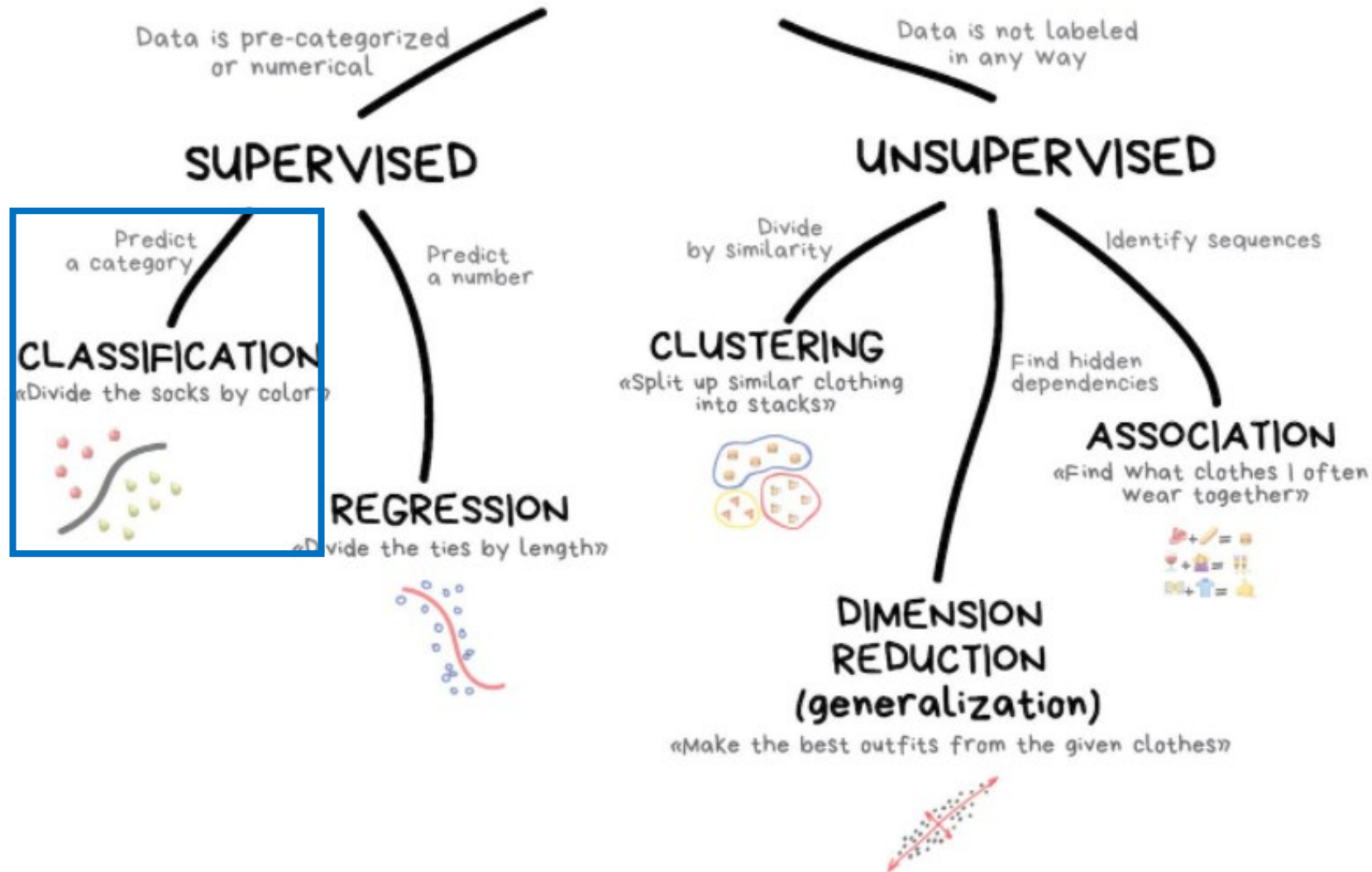
Validation Data

Testing Data

# Classification or Regression?

- Is this email spam?
- What will Google's stock price be tomorrow?
- Will Google's stock price go up tomorrow?
- Is this X/Twitter account a bot?
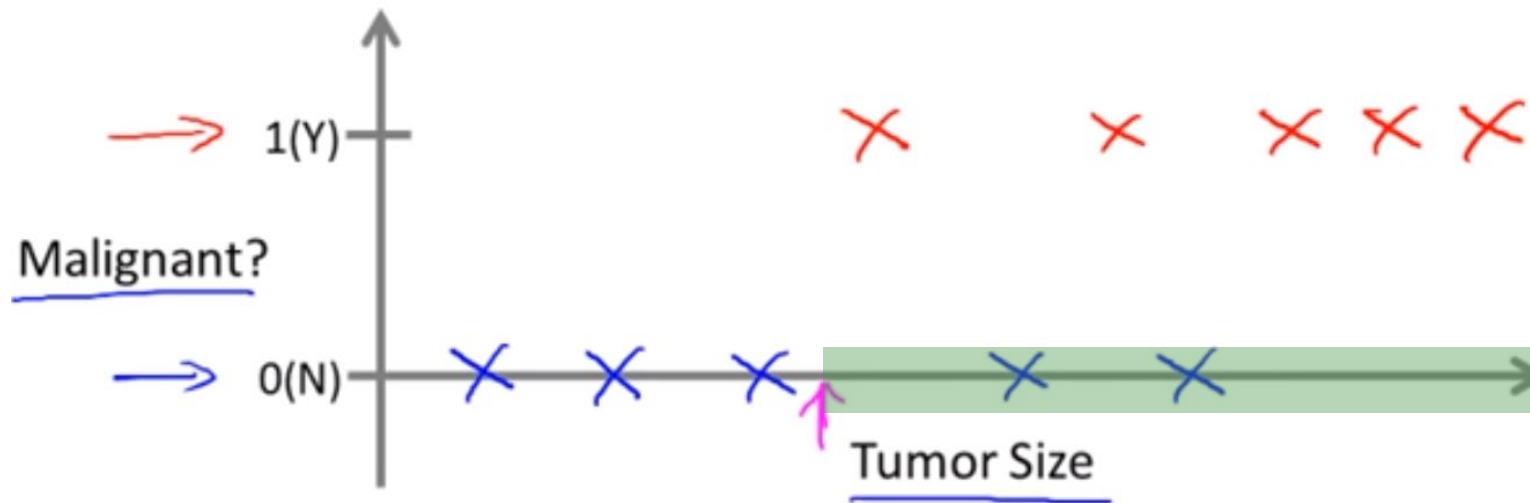- What will the total electricity demand be tomorrow in Athens, GA?

CLASSICAL MACHINE LEARNING

Data is pre-categorized or numerical

Data is not labeled in any way

**SUPERVISED**

**UNSUPERVISED**

Predict a category

Predict a number

Divide by similarity

Identify sequences

**CLASSIFICATION**
«Divide the socks by color»

**REGRESSION**
«Divide the ties by length»

**CLUSTERING**
«Split up similar clothing into stacks»

Find hidden dependencies

**ASSOCIATION**
«Find what clothes I often wear together»

**DIMENSION REDUCTION (generalization)**
«Make the best outfits from the given clothes»

# Predictive Model: Classification

Suppose you want to predict whether someone's breast cancer is malignant

**Estimate class probability (e.g., with a logistic regression)**

# Predictive Model: Classification Performance

|           |          | **Actual** | |
|-----------|----------|----------------|----------------|
|           |          | **Positive**   | **Negative**   |
| **Predicted** | **Positive** | True Positives | False Positives |
|           | **Negative** | False Negatives | True Negatives |

- An email sent from your friend identified as spam (+: spam)
- A fraud successfully caught by the system monitor (+: fraud)
- An internet intrusion passed as a normal activity (+: intrusion)
- A recovered patient approved to be discharged (+: still sick)

# Predictive Model: Classification Performance

Additional metrics

|  |  | Actual | |
|---|---|---|---|
|  |  | **Positive** | **Negative** |
| **Predicted** | **Positive** | True Positives | False Positives |
|  | **Negative** | False Negatives | True Negatives |

- Precision = true positives / (true positives + false positives)
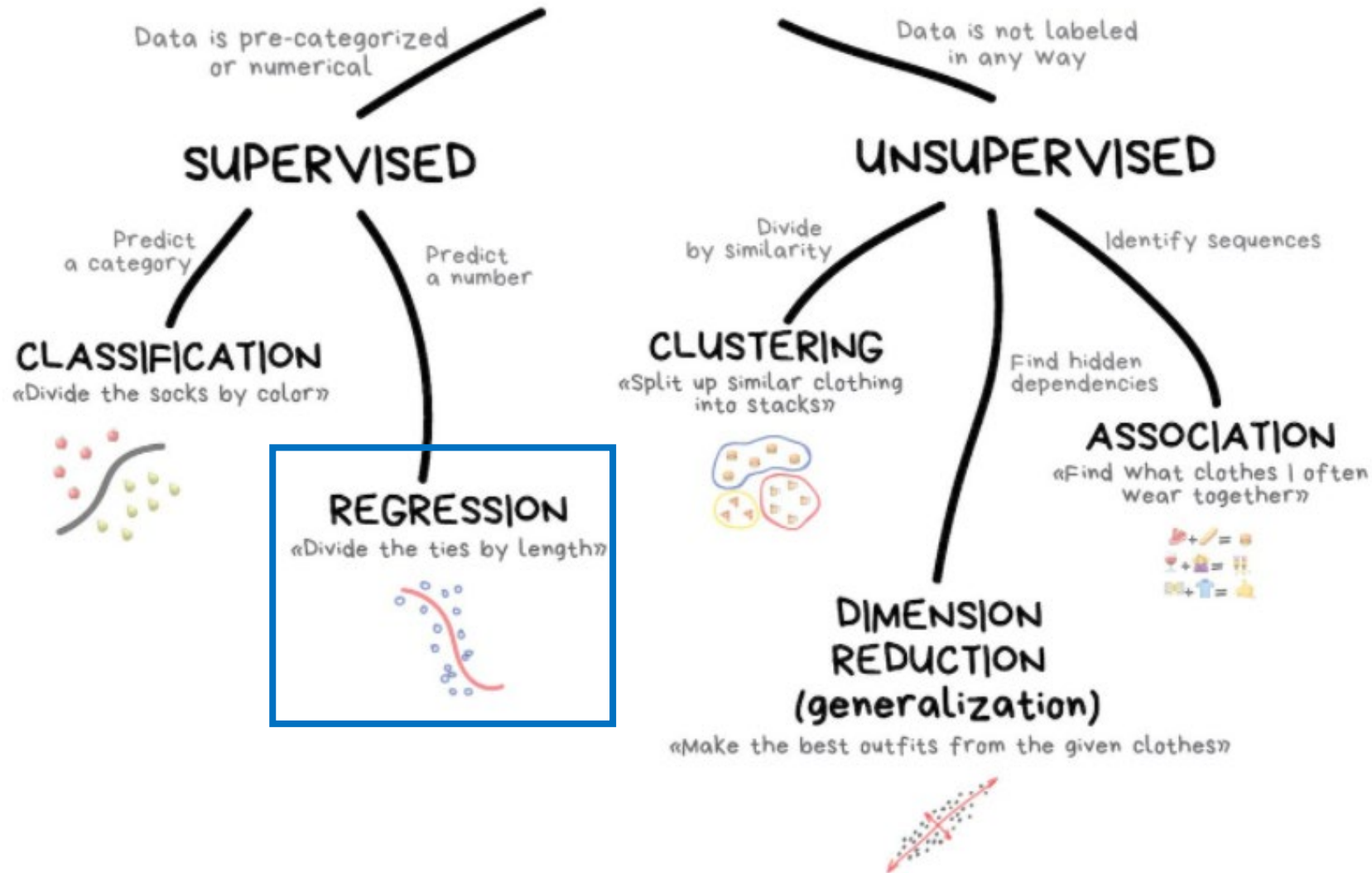- Recall = true positives / (true positives + false negatives)

# **Exercise**

Model X is developed to help doctors diagnose breast cancer. Its performance was evaluated over a test data set of 1000 patients.

- What is the accuracy of Model X?
- What is the precision and recall?

Actual
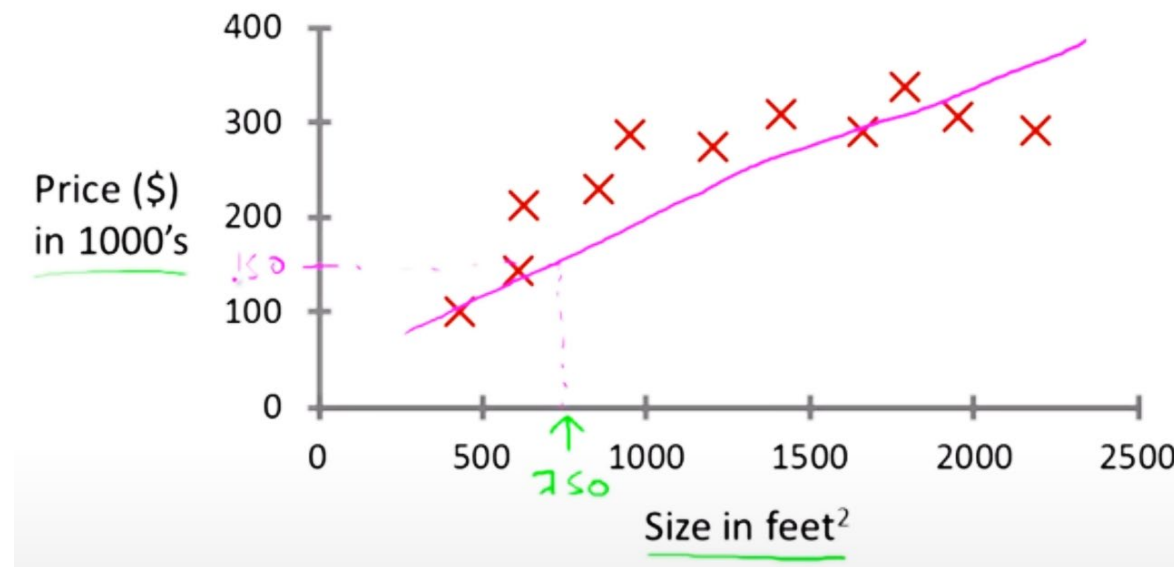
Predicted

|   | + | - |
|---|---|---|
| **+** | 16 | 4 |
| **-** | 3 | 977 |

# CLASSICAL MACHINE LEARNING

Data is pre-categorized or numerical

Data is not labeled in any way

## SUPERVISED

## UNSUPERVISED

Predict a category

Predict a number

Divide by similarity

Identify sequences

### CLASSIFICATION
«Divide the socks by color»

### REGRESSION
«Divide the ties by length»

### CLUSTERING
«Split up similar clothing into stacks»

Find hidden dependencies

### ASSOCIATION
«Find what clothes I often wear together»

### DIMENSION REDUCTION (generalization)
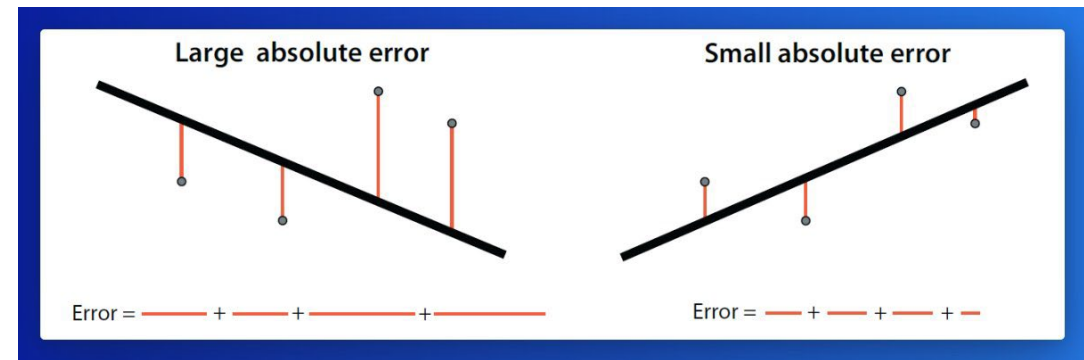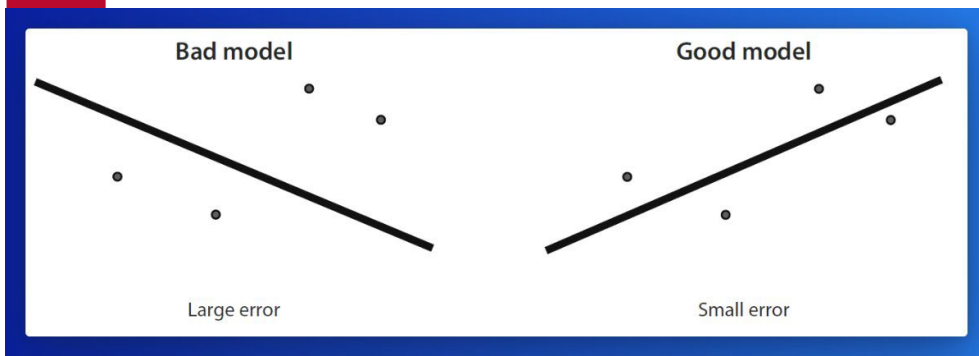«Make the best outfits from the given clothes»

# Predictive Model: Regression

Suppose you want to predict house prices, and you have some data about the price of a house (in thousands of $) over size (sqft)

**Estimate numeric value (e.g., with a linear regression)**

# Model Evaluation

Error for data record = predicted (p) minus actual (a)

**RMSE: Root Mean Squared Error**: $\sqrt{\frac{1}{n}\sum_{1}^{n}(Y_i - \hat{Y}_i)^2}$

MAE: Mean Absolute Error: $\frac{1}{n}\sum_{1}^{n}|(Y_i - \hat{Y}_i)|$

MAPE: Mean Absolute Percentage Error: $\frac{100}{n}\sum_{1}^{n}\left|\frac{Y_i - \hat{Y}_i}{Y_i}\right|$

Total SSE: Total Sum of Squared Errors: $\sum_{1}^{n}(Y_i - \hat{Y}_i)^2$

# RMSE

Error for data record = predicted (p) minus actual (a)

RMSE = how much the *p*'s diverge from the *a*'s, on average

Assume the regression equation is y = 1.74x. What is the root mean squared error for the sample dataset?

| x | a | p | (p − a)^2 |
|---|---|---|---|
| 1 | 2 | 1.74 | 0.0676 |
| 2 | 5 | 3.48 | 2.3104 |
| -1 | -2 | -1.74 | 0.0676 |

**RMSE =**

$$\sqrt{(0.0676 + 2.3104 + 0.0676) / 3}$$
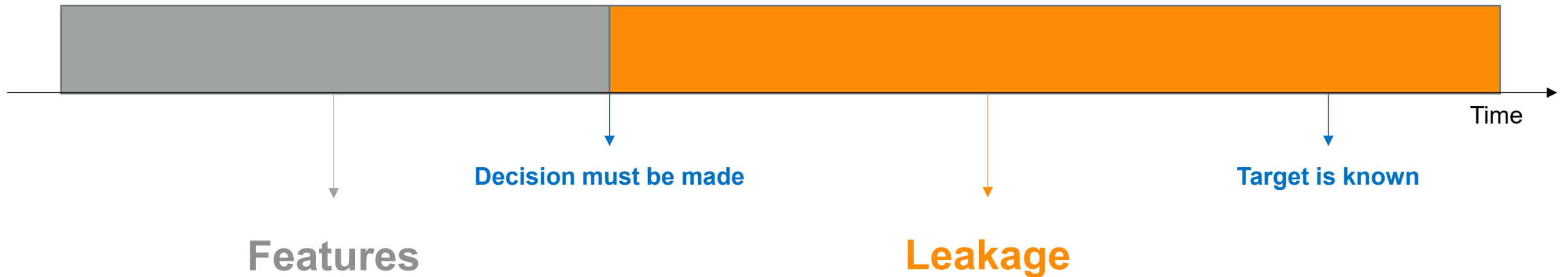
$$= \sqrt{0.8152}$$

$$= .903$$

# Things to Consider

- Is there a specific, quantifiable target that you are interested in predicting?
    - If yes, is it a **class or a number**?
        - Think about the decision

- Do you have data on the target?
    - Do you have **enough** data?
        - If the target is a class, a min of ~500 for each class type is needed

# Another Thing to Consider

- Do you have relevant data **prior** to the decision?
    - Think about the timing of decision and action leading up to it



Features      Decision must be made      Leakage      Target is known      Time
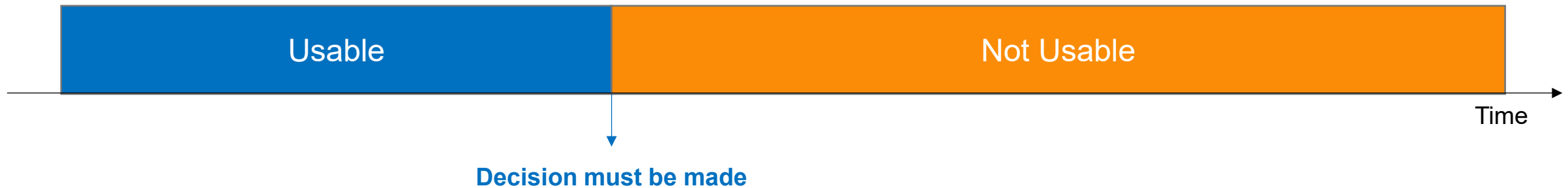
# **Predicting Loan Default**

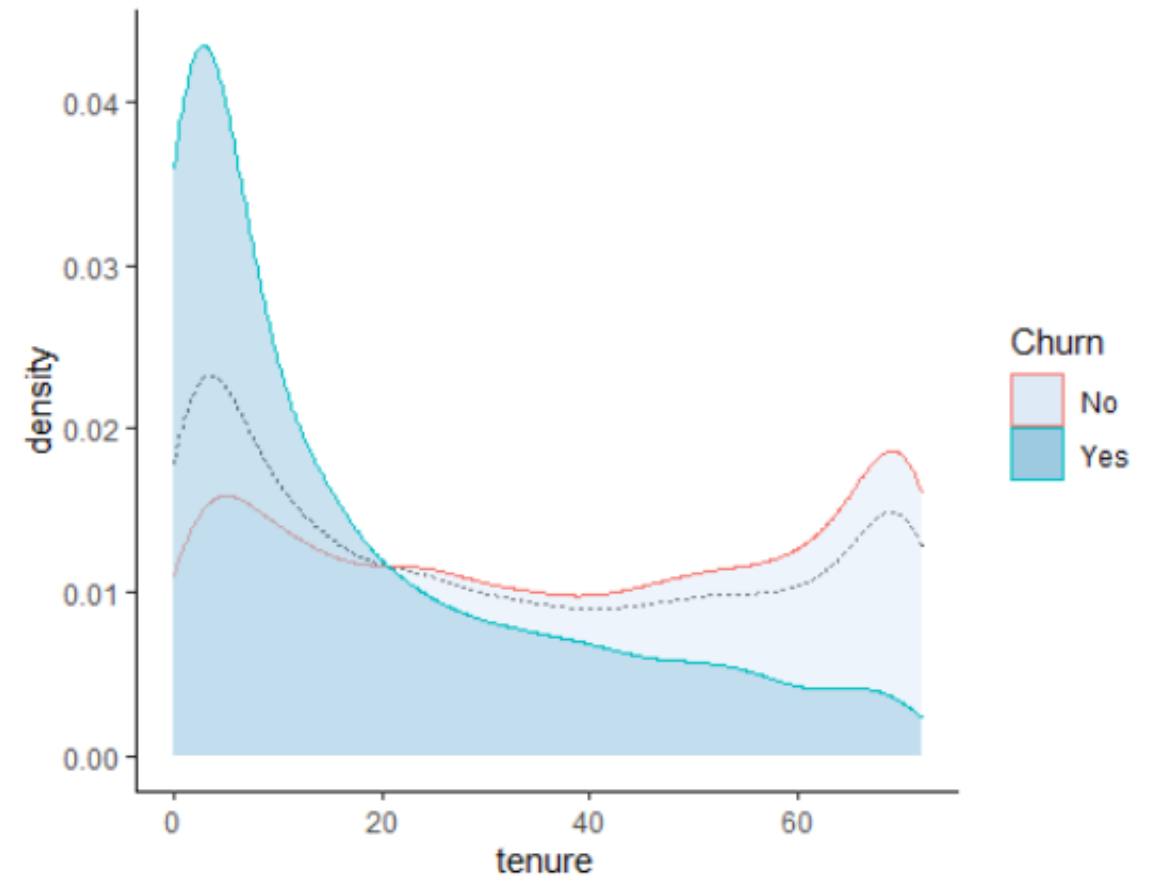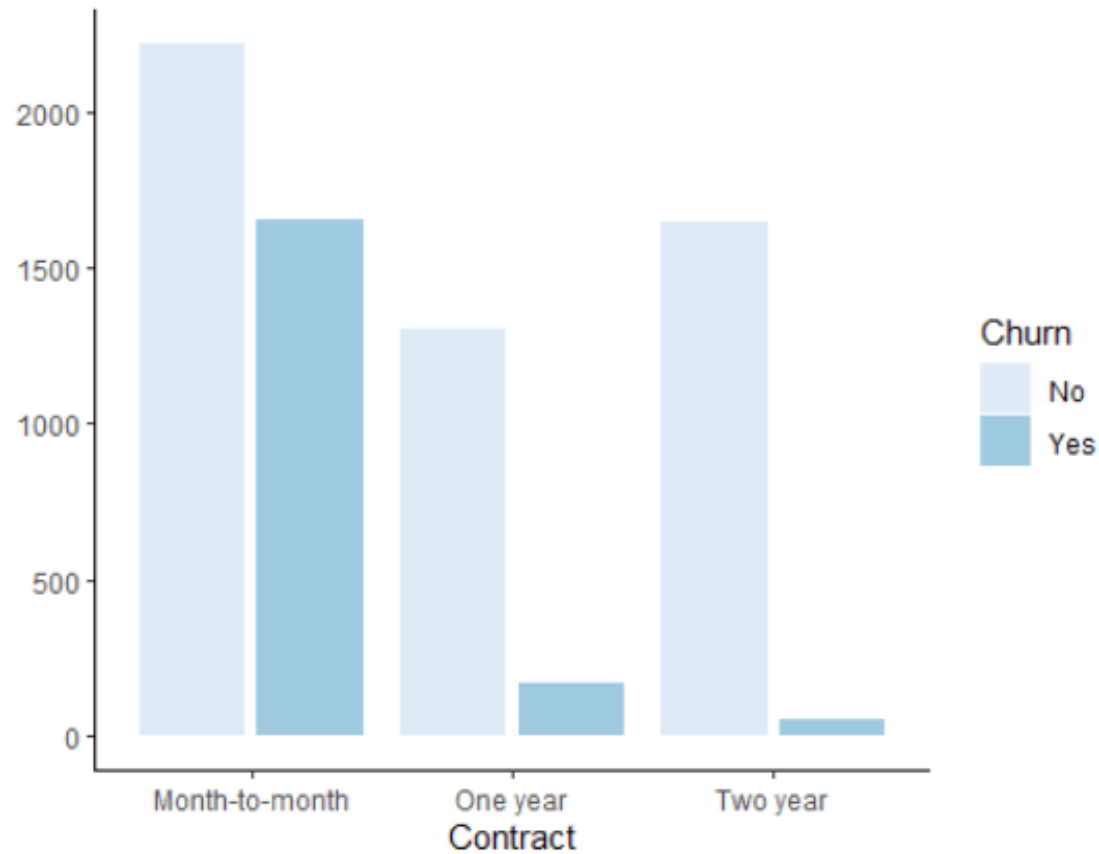What information about its customers can a bank use to predict loan default?

- 

- 

-

# Avoiding Leakage

- Do you have relevant data **prior** to the decision?
  - Think about the timing of decision and action leading up to it

| Usable | Not Usable |
|---|---|

Decision must be made

Time

# Exploratory analysis

# Class Activity 1

An online clothing store wants to create a supervised model that will offer personalized clothing recommendations to customers. This means that the model will recommend certain clothes to Janeth and different clothes to Joe. The model will use past purchasing behavior to generate training data. Mention five attributes that are useful for the model and justify your answer.

# Class Activity 1

An online clothing store wants to create a supervised model that will offer personalized clothing recommendations to customers. This means that the model will recommend certain clothes to Janeth and different clothes to Joe. The model will use past purchasing behavior to generate training data. Mention five attributes that are useful for the model and justify your answer.

- Size

- User clicks on the product description

- Clothes beauty