# Introduction / Background

- Consists of 100,000 records of patients diagnosed with diabetes collected from 130 US hospitals over 10 years.
- Some of the questions we intend to answer from this dataset is whether we could predict readmission of a patient given certain attributes and whether there are any associations between the diabetic medications given to patients.
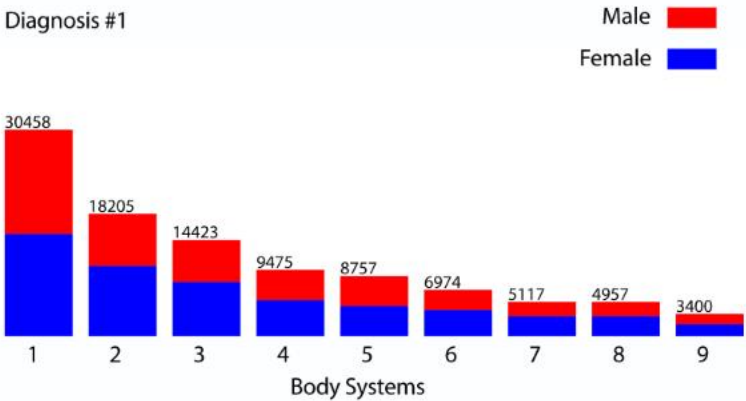
# Exploratory Analysis

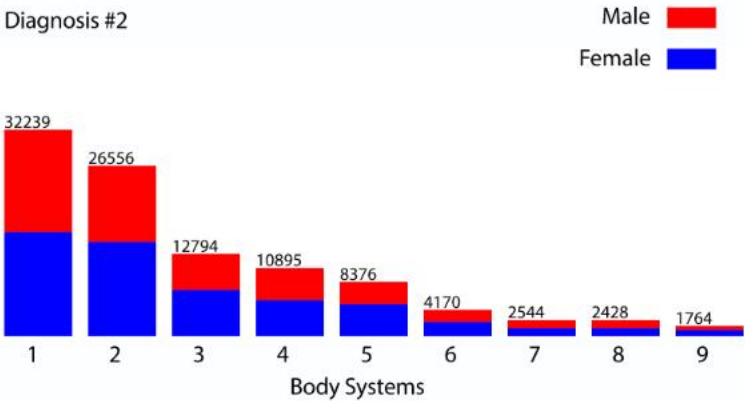- There are 101, 766 patient encounters
- 55 features for each encounter

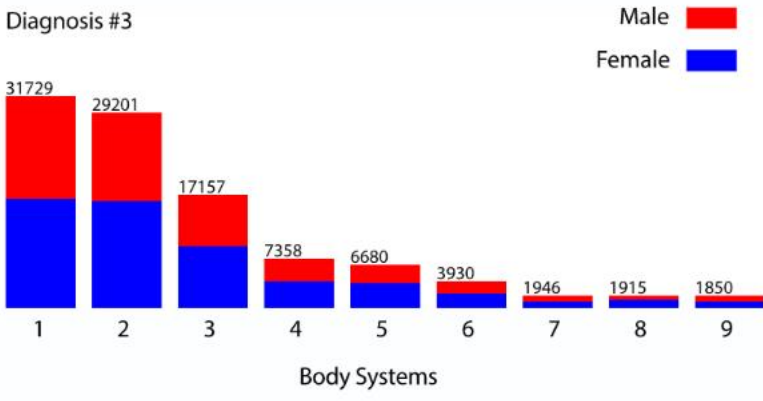**ICD9 values to Body systems**

**Preprocessing**

# ICD9 values to Body systems



**Diagnosis #1**

Male ■ (red)
Female ■ (blue)

30458, 18205, 14423, 9475, 8757, 6974, 5117, 4957, 3400

Body Systems (1–9)

1. Circulatory
2. Other
3. Respiratory
4. Digestive
5. Diabetes
6. Injury
7. Genitourinary
8. Musculoskeletal
9. Neoplasms

**Diagnosis #2**

Male ■ (red)
Female ■ (blue)

32239, 26556, 12794, 10895, 8376, 4170, 2544, 2428, 1764

Body Systems (1–9)

1. Circulatory
2. Other
3. Diabetes
4. Respiratory
5. Genitourinary
6. Digestive
7. Neoplasms
8. Injury
9. Musculoskeletal

**Diagnosis #3**

Male ■ (red)
Female ■ (blue)

31729, 29201, 17157, 7358, 6680, 3930, 1946, 1915, 1850

Body Systems (1–9)

1. Circulatory
2. Other
3. Diabetes
4. Respiratory
5. Genitourinary
6. Digestive
7. Injury
8. Musculoskeletal
9. Neoplasms

# Preprocessing

As part of preprocessing we identified that the following features were unuseful and hence dropped them.

- Encounter Id
- Patient number
- Weight
- Payer code
- Medical specialty

# FP-Growth

## Attributes:

- metformin
- repaglinide
- nateglinide
- chlorpropamide
- glimepiride
- acetohexamide
- glipizide
- glyburide
- tolbutamide
- pioglitazone
- rosiglitazone
- acarbose
- miglitol
- troglitazone
- tolazamide
- examide
- citoglipton
- insulin

We found 2 rules:

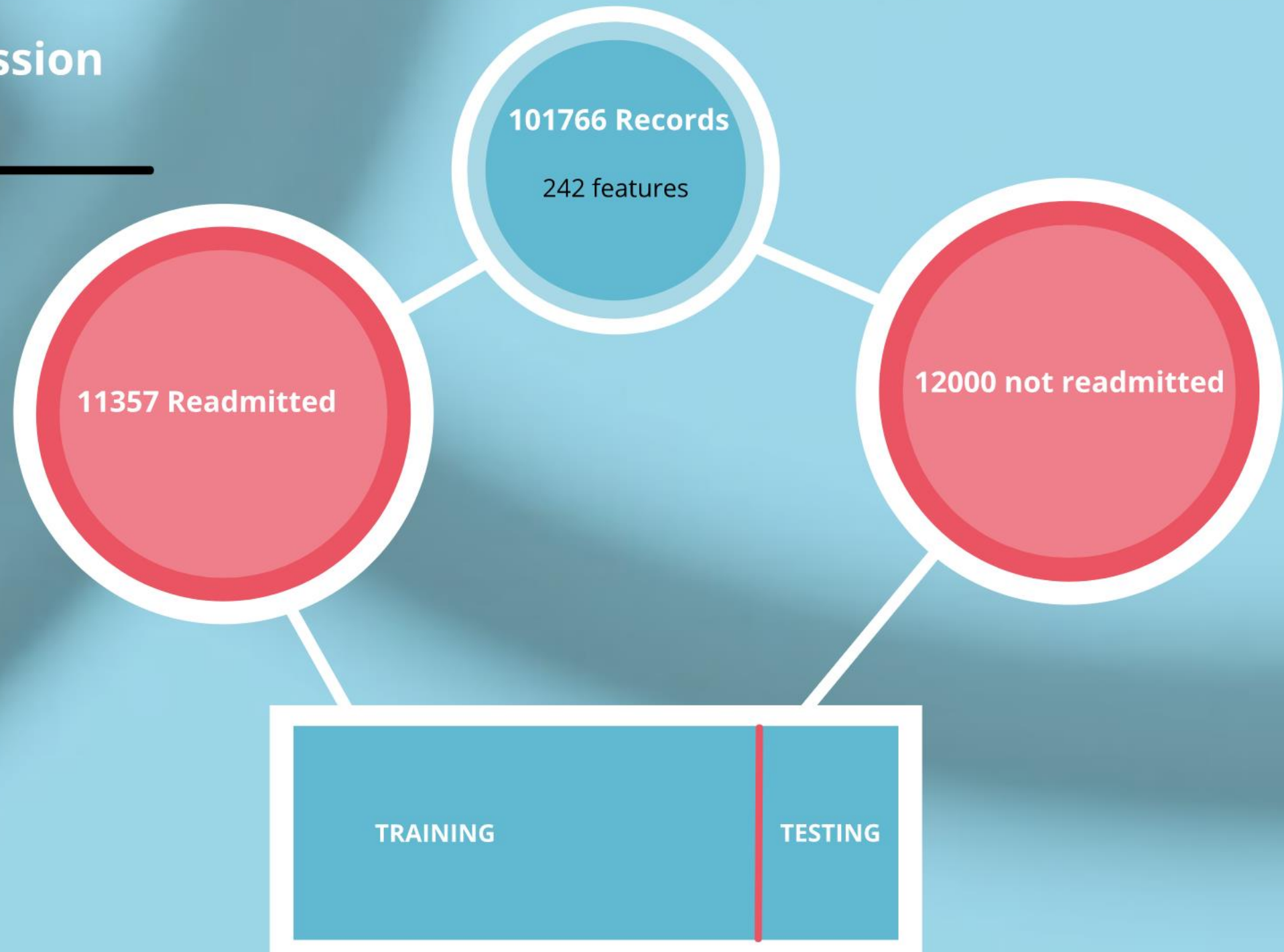metformin ==> insulin  with Confidence 0.5 and Lift 0.72

insulin ==> metformin with Confidence 0.18 and Lift 0.72

# Results

Training Data: 242 features
Model Accuracy: 56.57%


Applied PCA to reduce dimensionality to 6 (48% variance in original data)
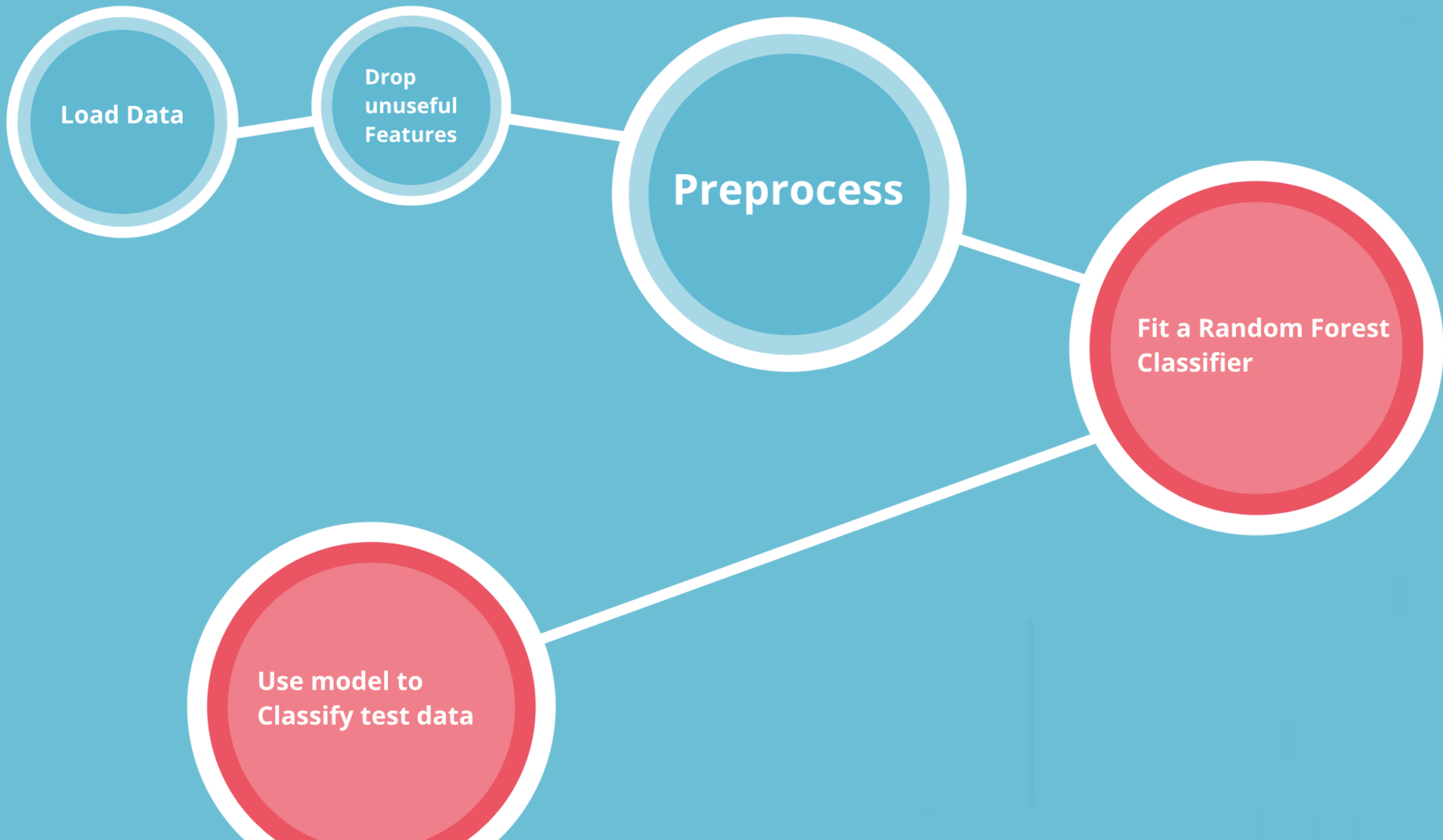Model Accuracy: 60.64%

# Random Forest

- Ensemble learning method for classification
- It constructs a large number of decision trees on the training data to create a classifier model, this model is later used on the test data for prediction.

**Preprocess**

**Results**

# Results

Training Data: 44 features
Model Accuracy: 51.19%

## Confusion Matrix

# Conclusion

Unsupervised

- FP-Growth
- PCA

Supervised

- Logistic Regression
- Random Forests