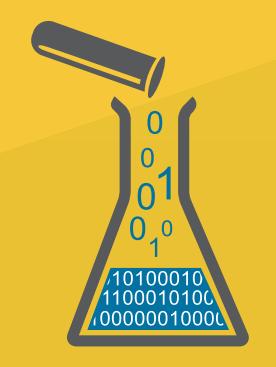


YELLOWBRICK: STEERING SCIKIT-LEARN WITH VISUAL TRANSFORMERS

By Benjamin Bengfort and Rebecca Bilbro

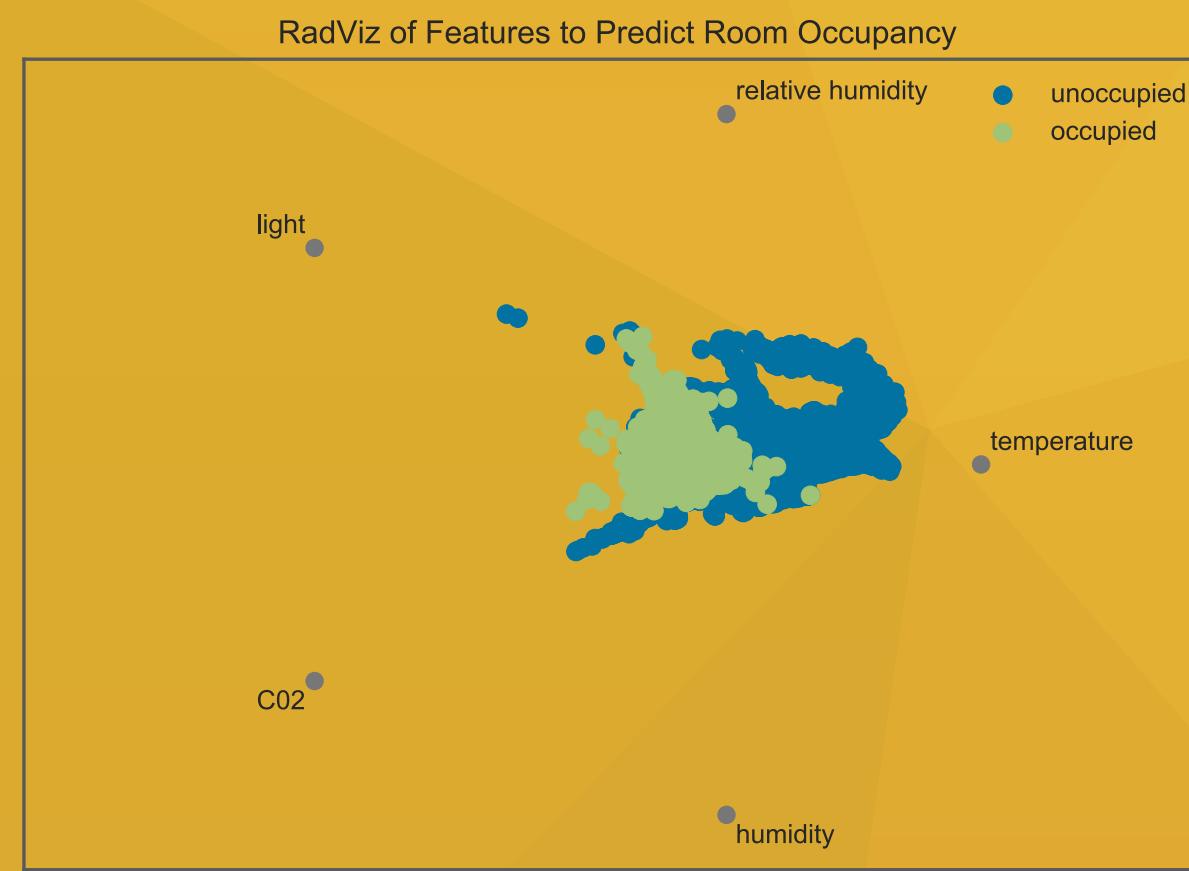


DistrictDataLabs
pycon.districtdatalabs.com

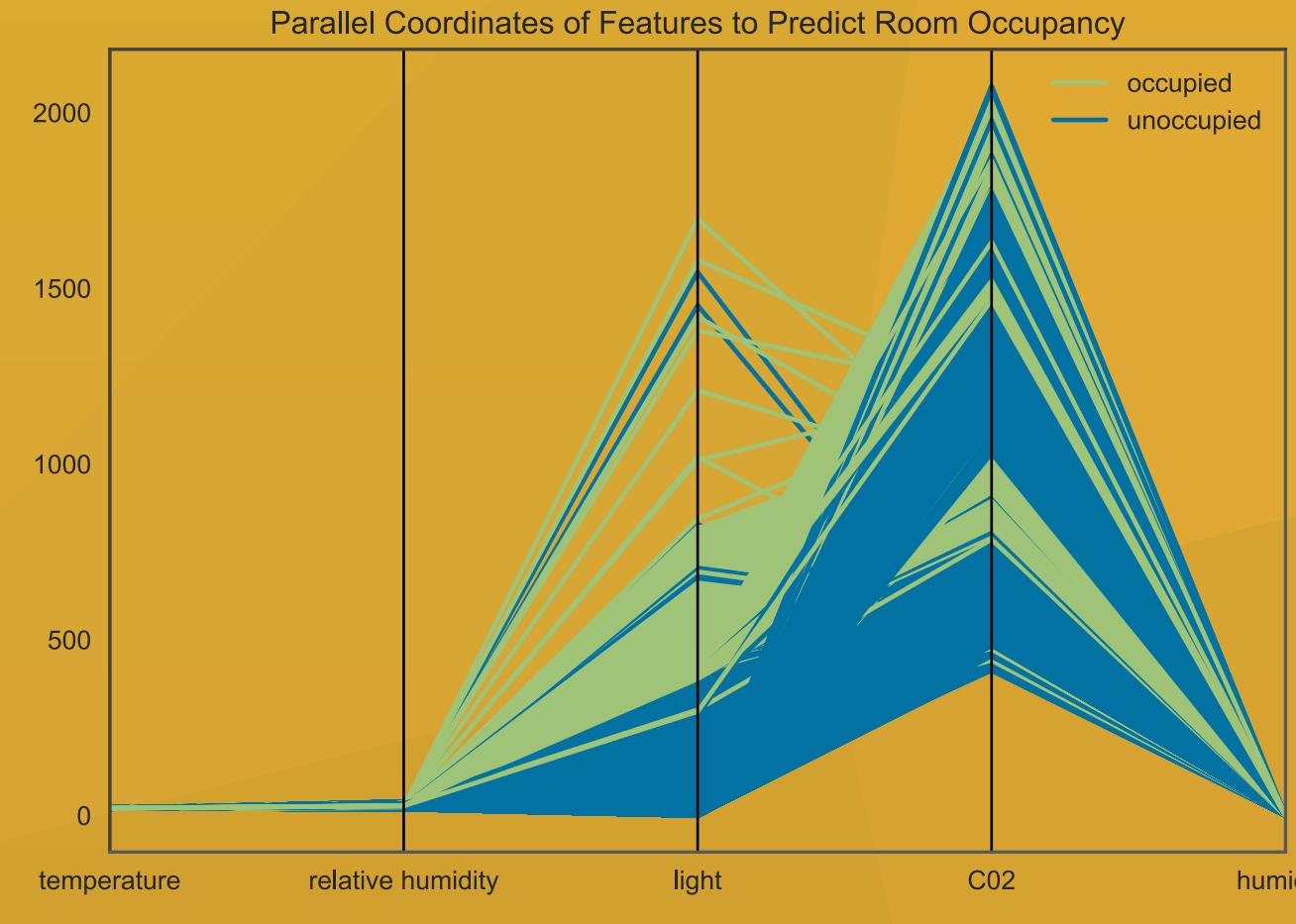
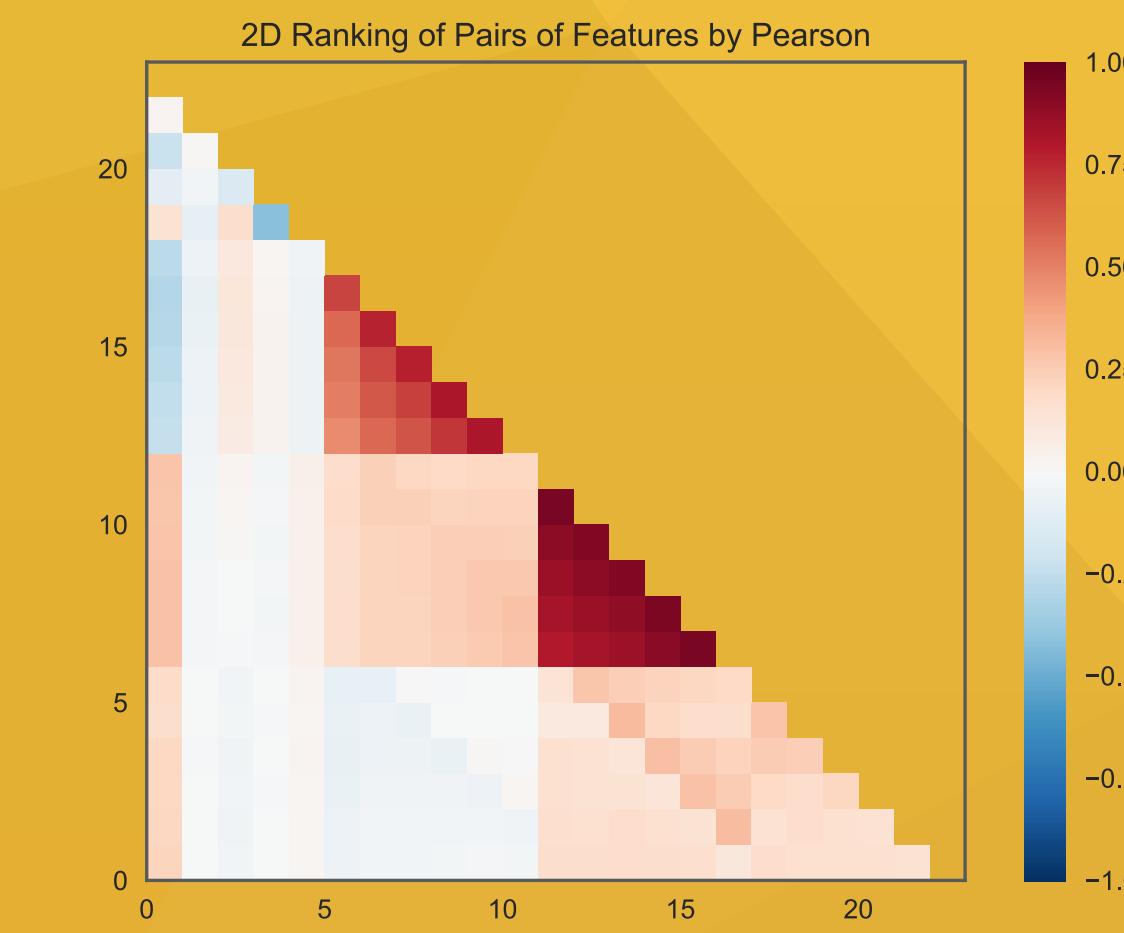
WHICH FEATURES DO I USE?

Given labeled data about rooms...

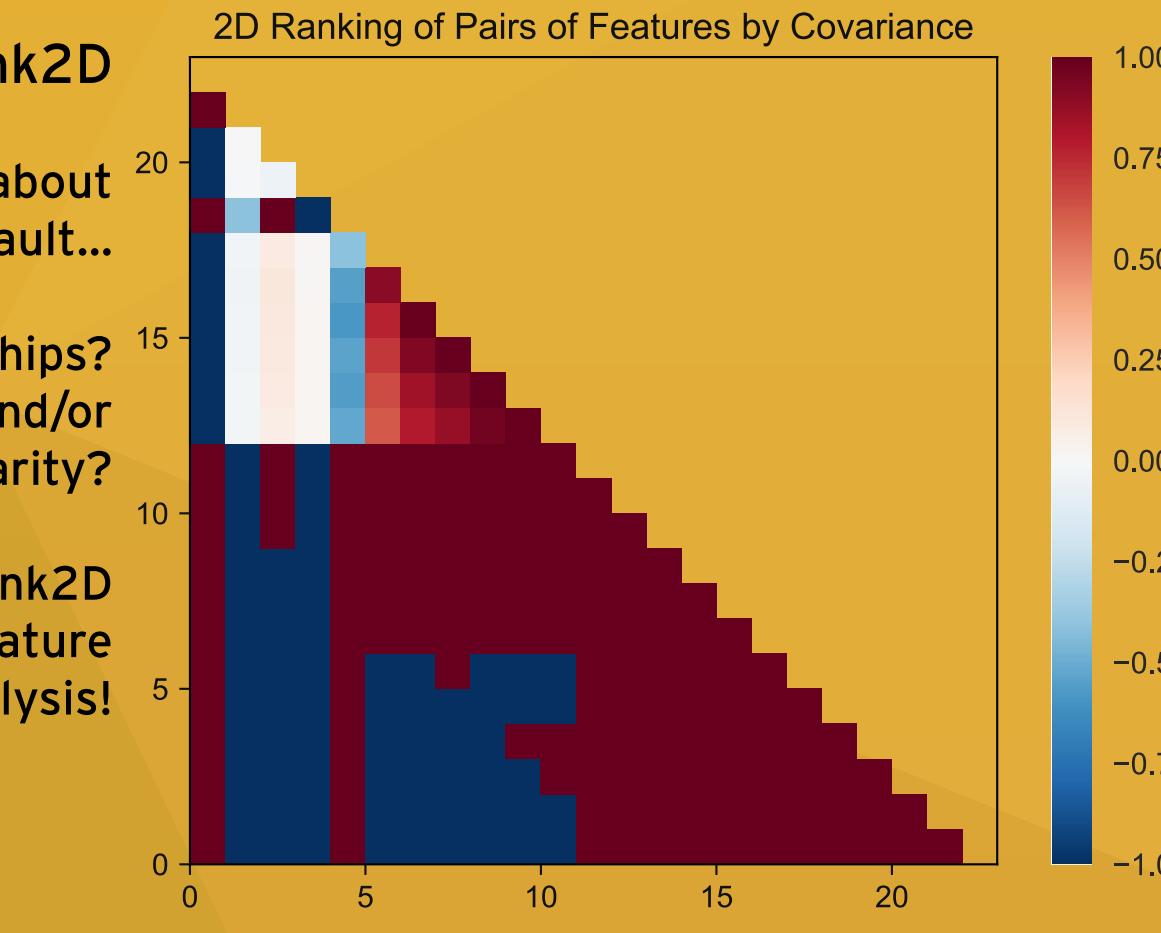
- Which features are most predictive?
- Empty or occupied?



Radviz and ParallelCoordinates
Use Yellowbrick Radial Visualizations or Parallel Coordinates to look for class separability!



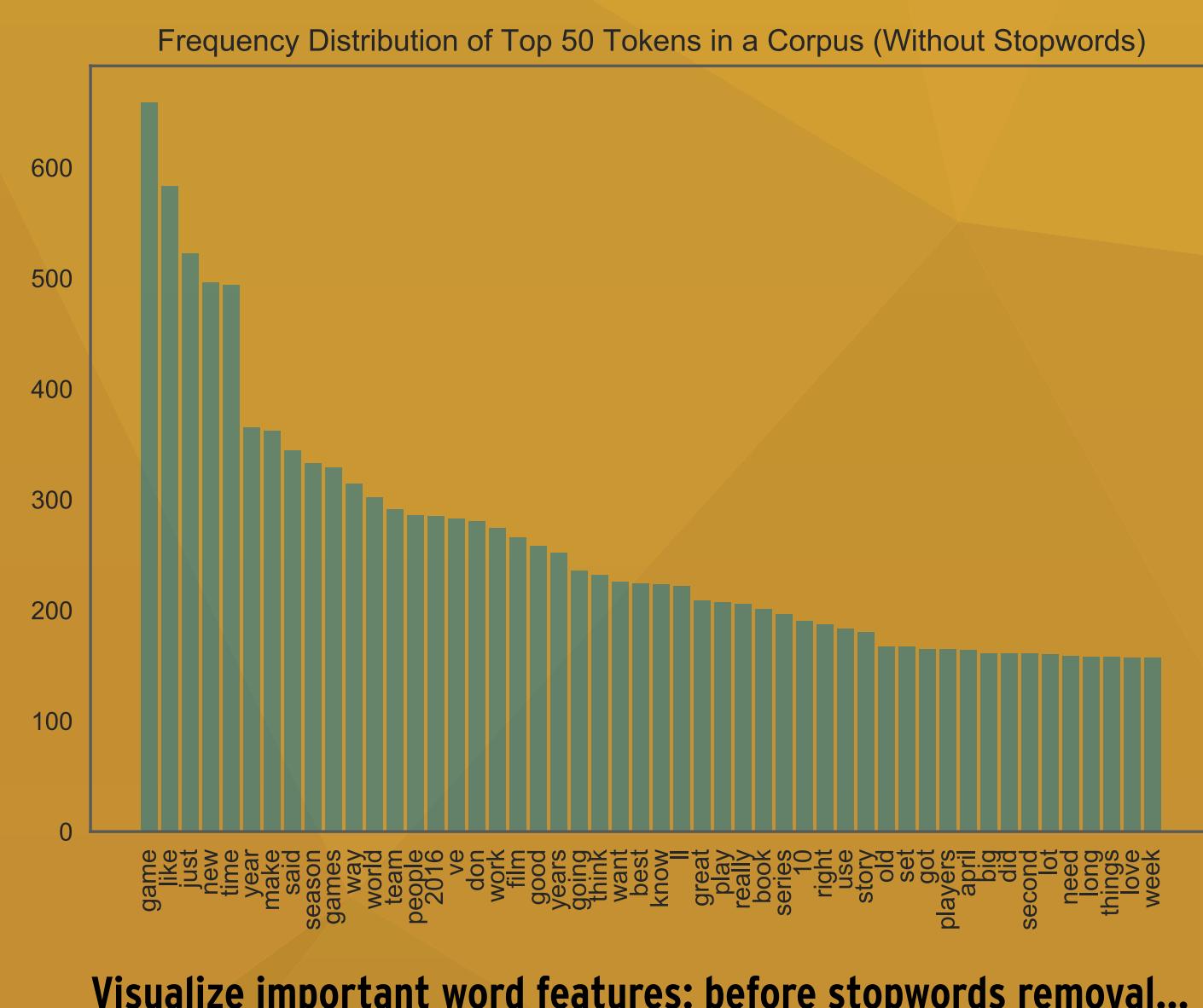
Rank2D
Given labeled data about credit card default...
• Feature relationships?
• Correlations and/or collinearity?
Use Yellowbrick Rank2D for pairwise feature analysis!



WORKING WITH TEXT DATA

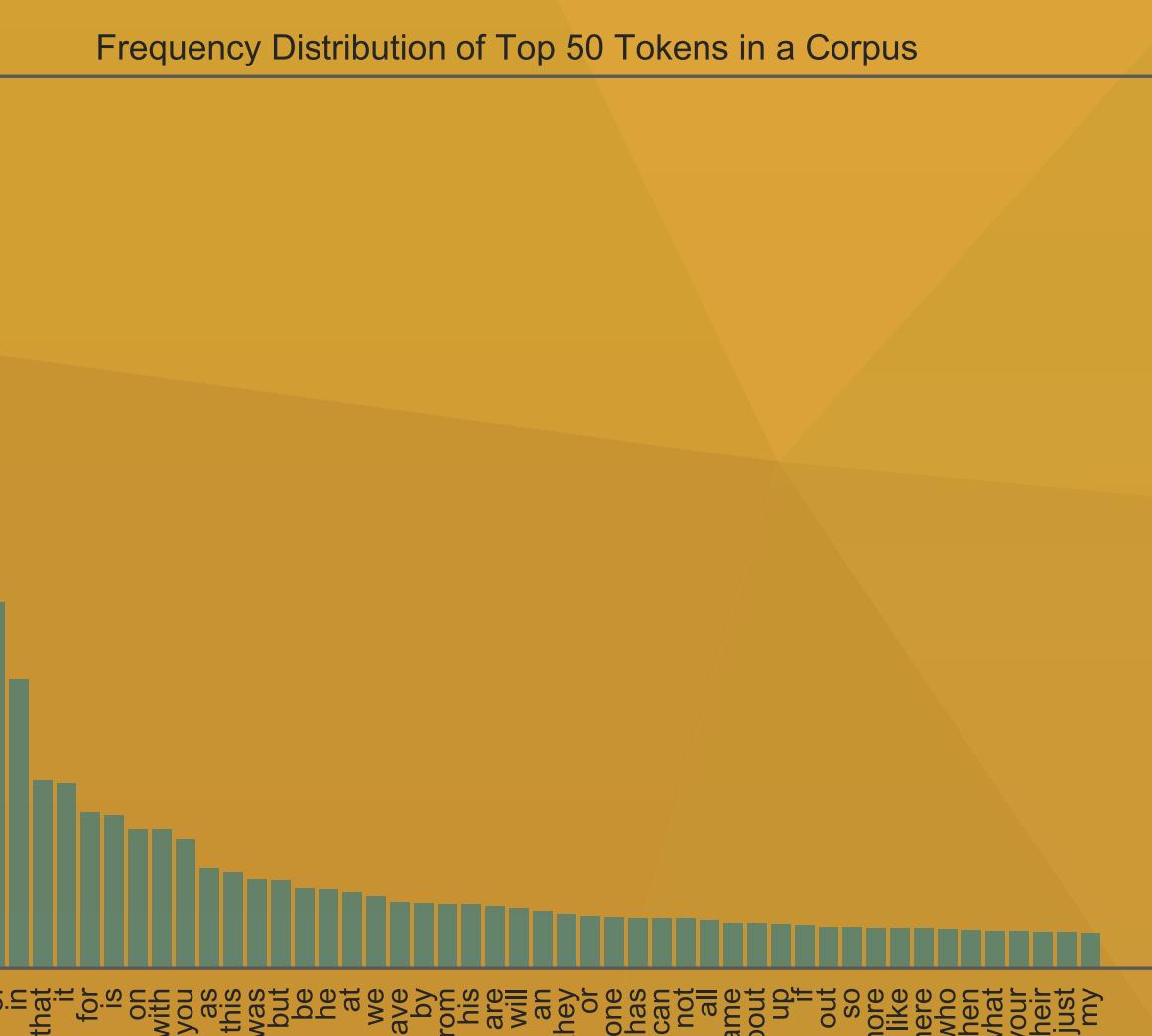
Text data is notoriously high-dimensional and hard to visualize.

Yellowbrick can help!



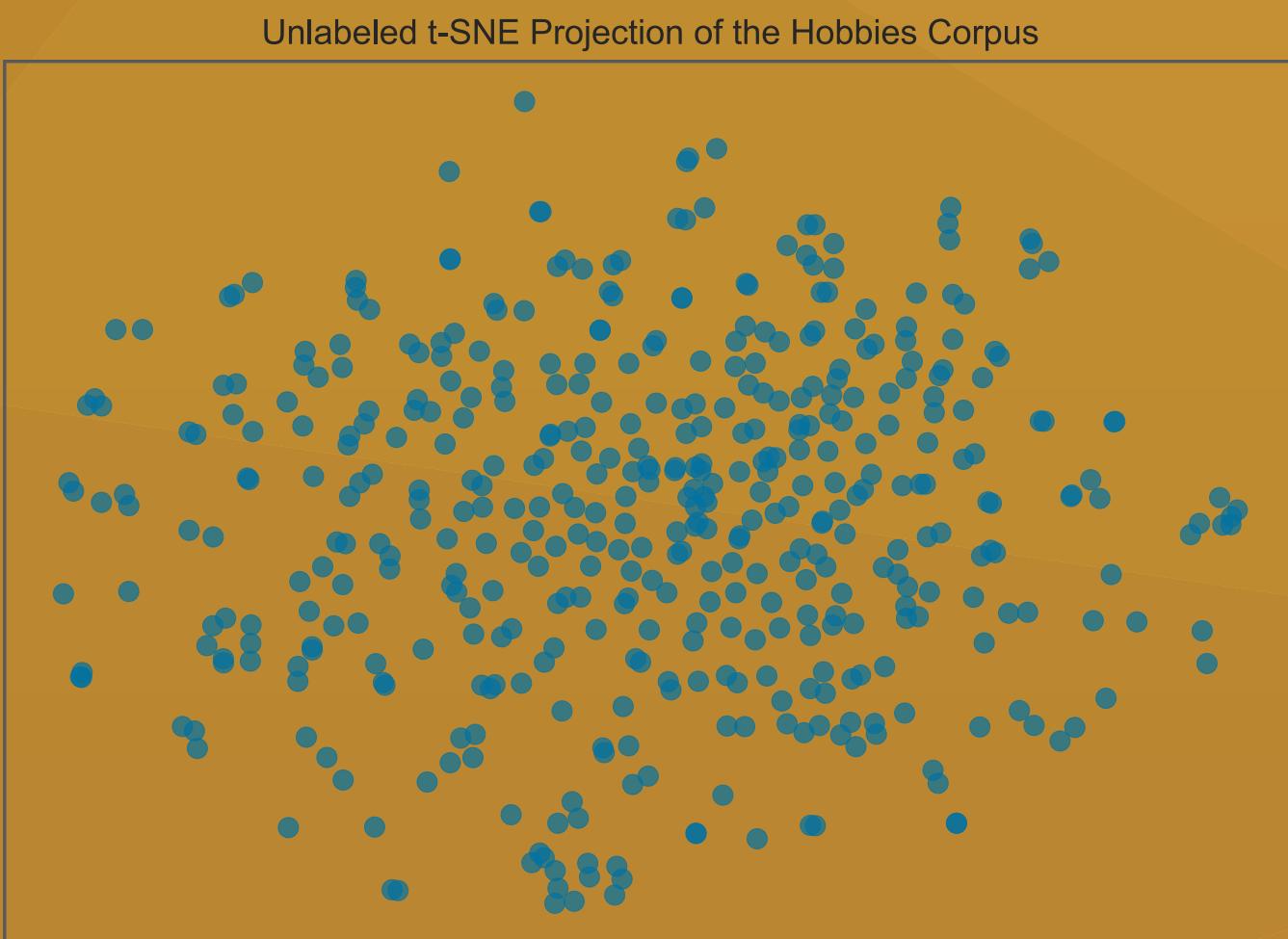
Visualize important word features: before stopwords removal...

Frequency Distributions

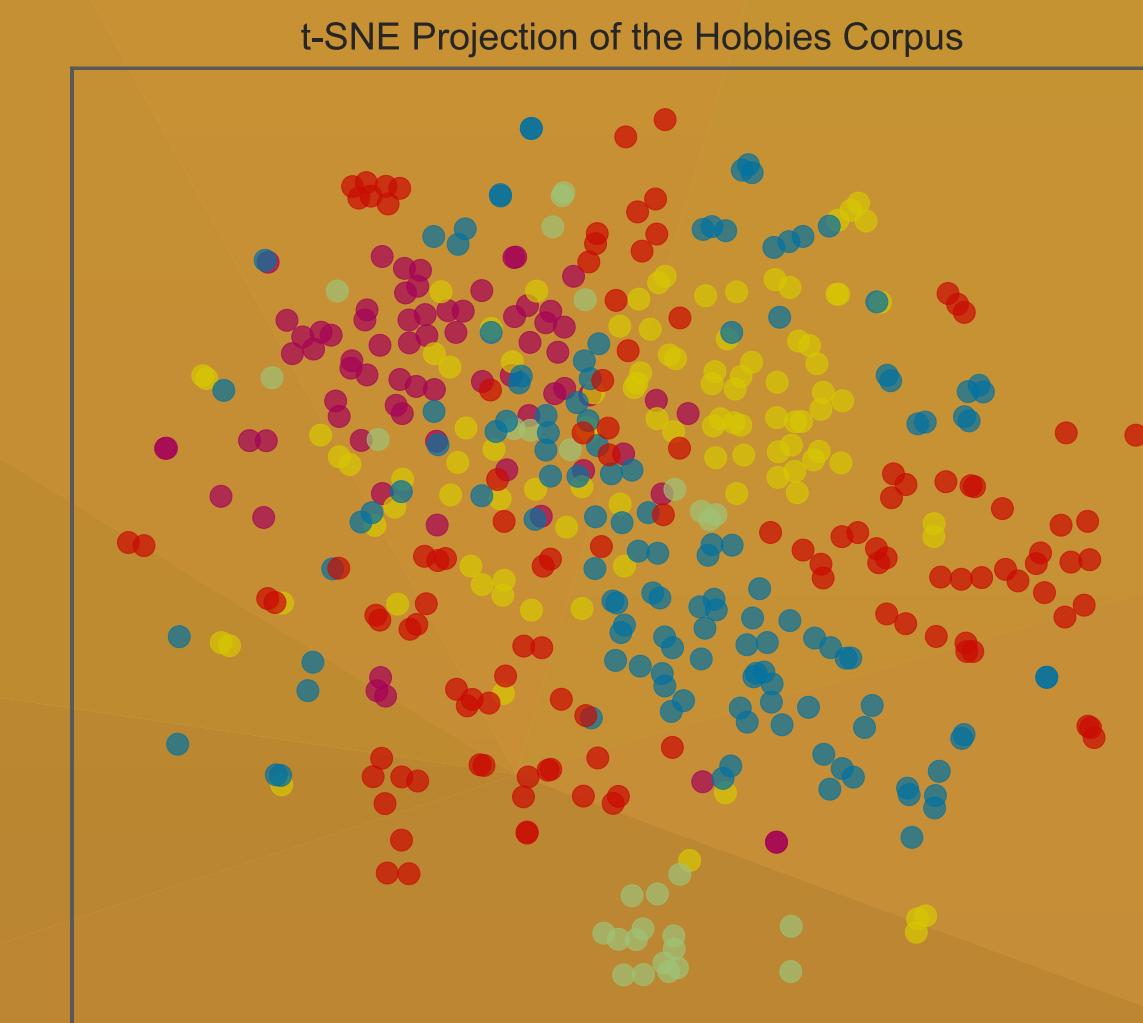


and after!

t-SNE



Visualize the distribution of corpus documents in 2 dimensions



ENTER YELLOWBRICK scikit-yb.org

Yellowbrick is a new Python library that:

- extends the Scikit-Learn API.
- enhances the model selection process.
- provides visual tools for feature analysis, diagnostics & steering.

In Scikit-Learn:

```
# Import the estimator
from sklearn.linear_model import Lasso

# Instantiate the estimator
model = Lasso()

# Fit the data to the estimator
model.fit(X_train, y_train)

# Generate a prediction
model.predict(X_test)
```

With Yellowbrick:

```
# Import the model and visualizer
from sklearn.linear_model import Lasso
from yellowbrick.regressor import PredictionError

# Instantiate the visualizer
visualizer = PredictionError(Lasso())

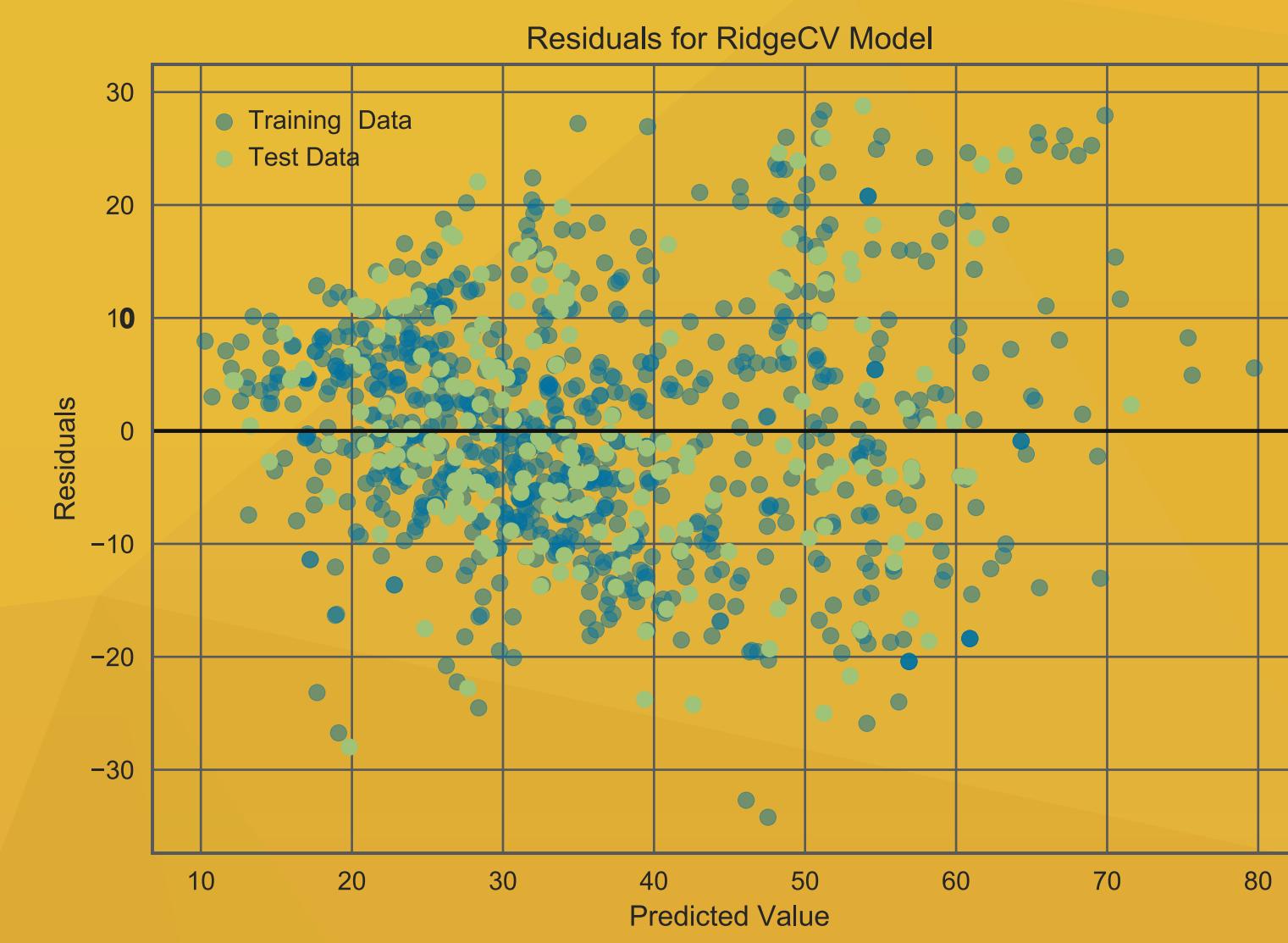
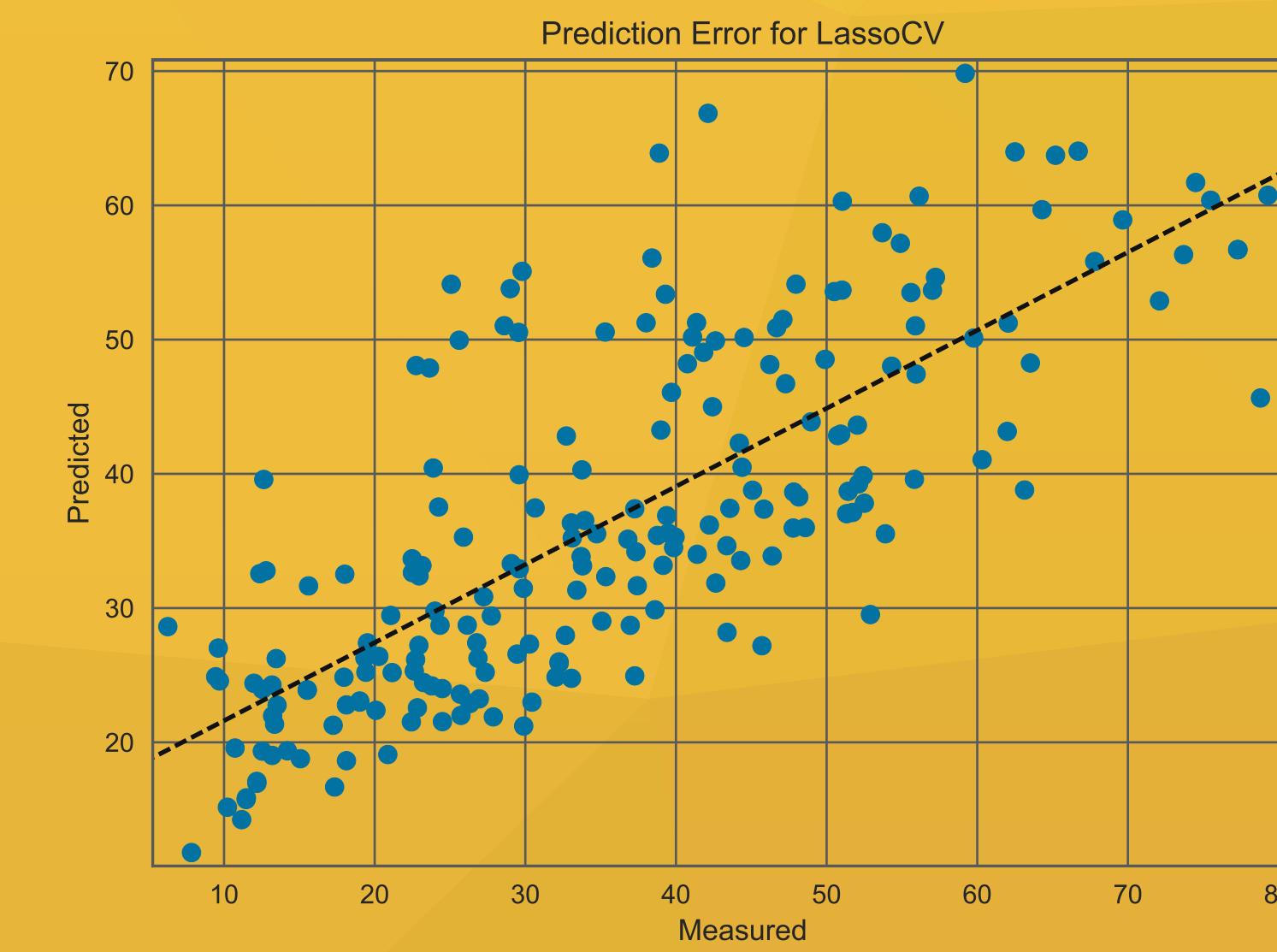
# Fit
visualizer.fit(X_train, y_train)

# Score and Visualize
visualizer.score(X_test, y_test)
visualizer.poof()
```

WHICH MODEL SHOULD I USE?

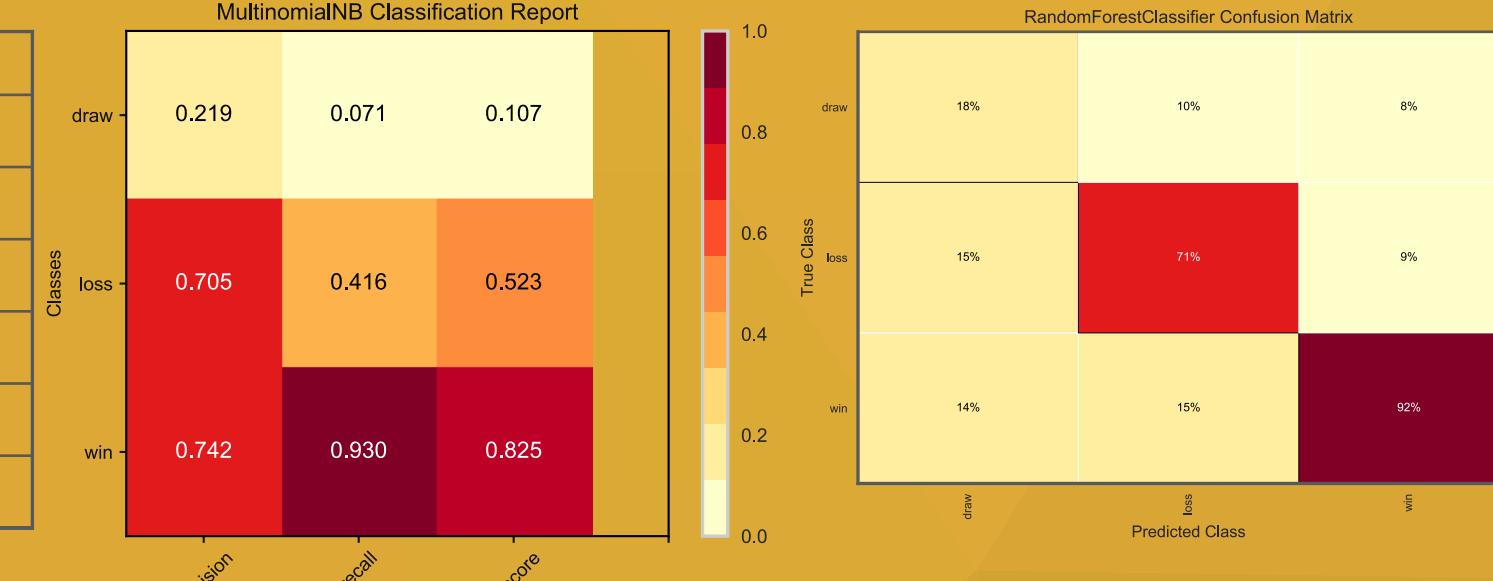
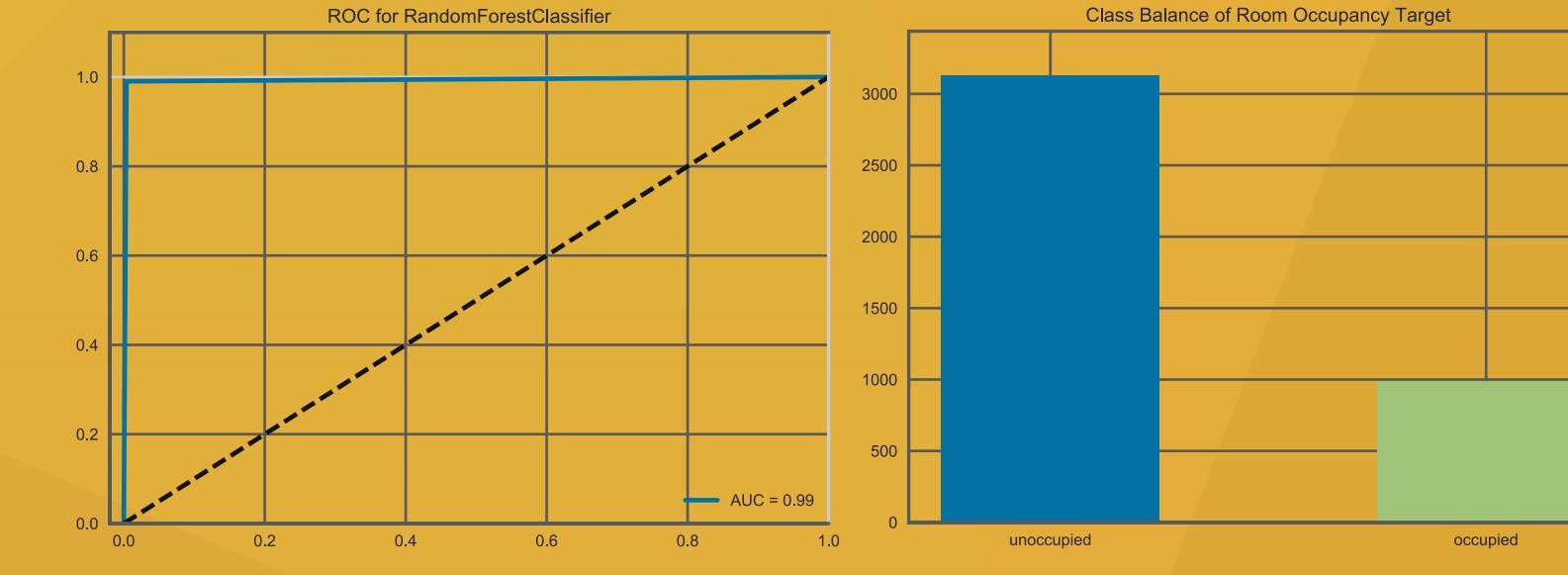
Prediction Error and Residuals Plot

Visualize the distribution of error to diagnose heteroscedasticity:



ROCAUC, Classification Report, Confusion Matrix, and Class Balance

ROCAUC helps us see overall accuracy; classification heatmap helps distinguish Type I & Type II error; and confusion matrix shows us error on a per-class basis. What to do with a low-accuracy classifier? Check for imbalance!

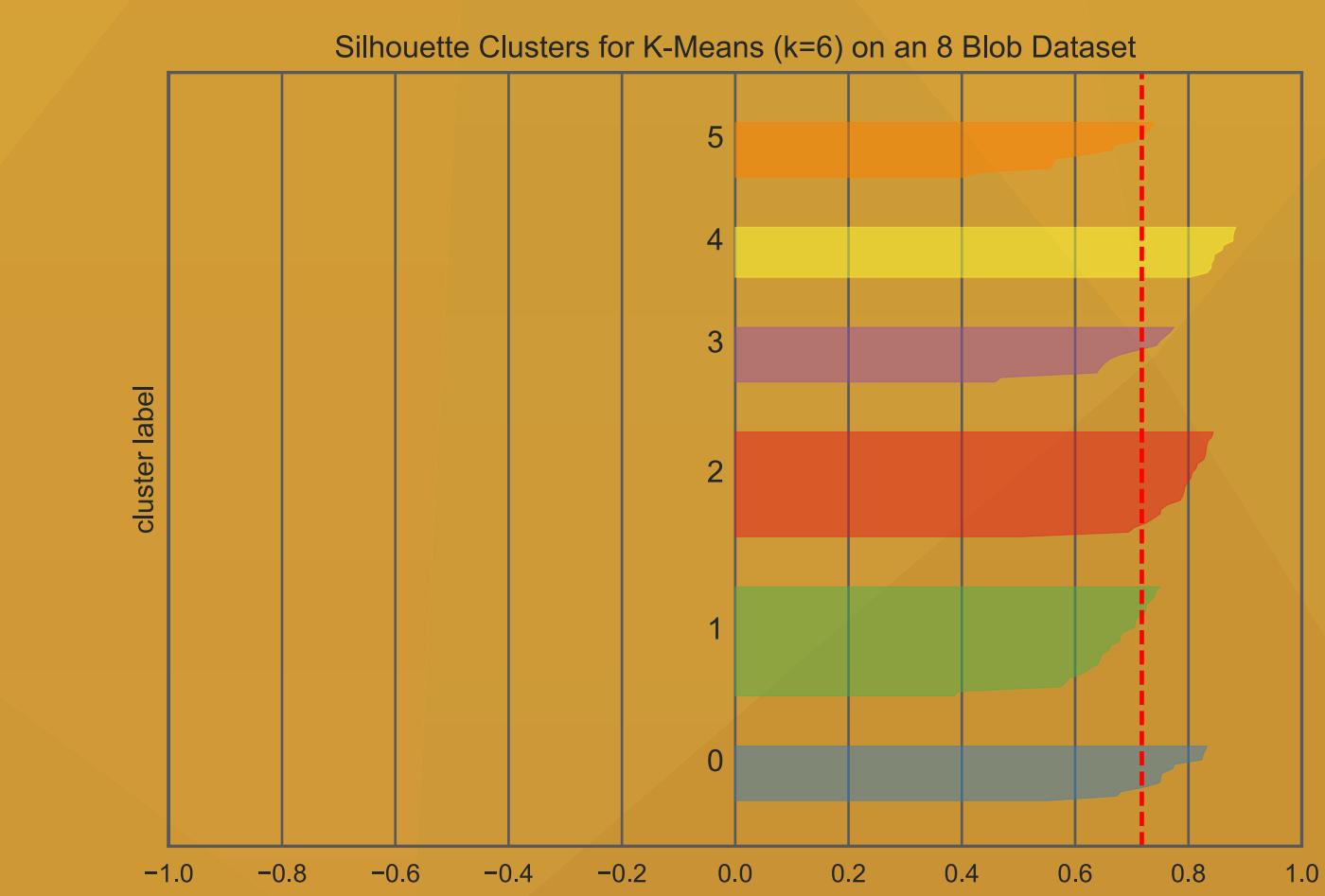
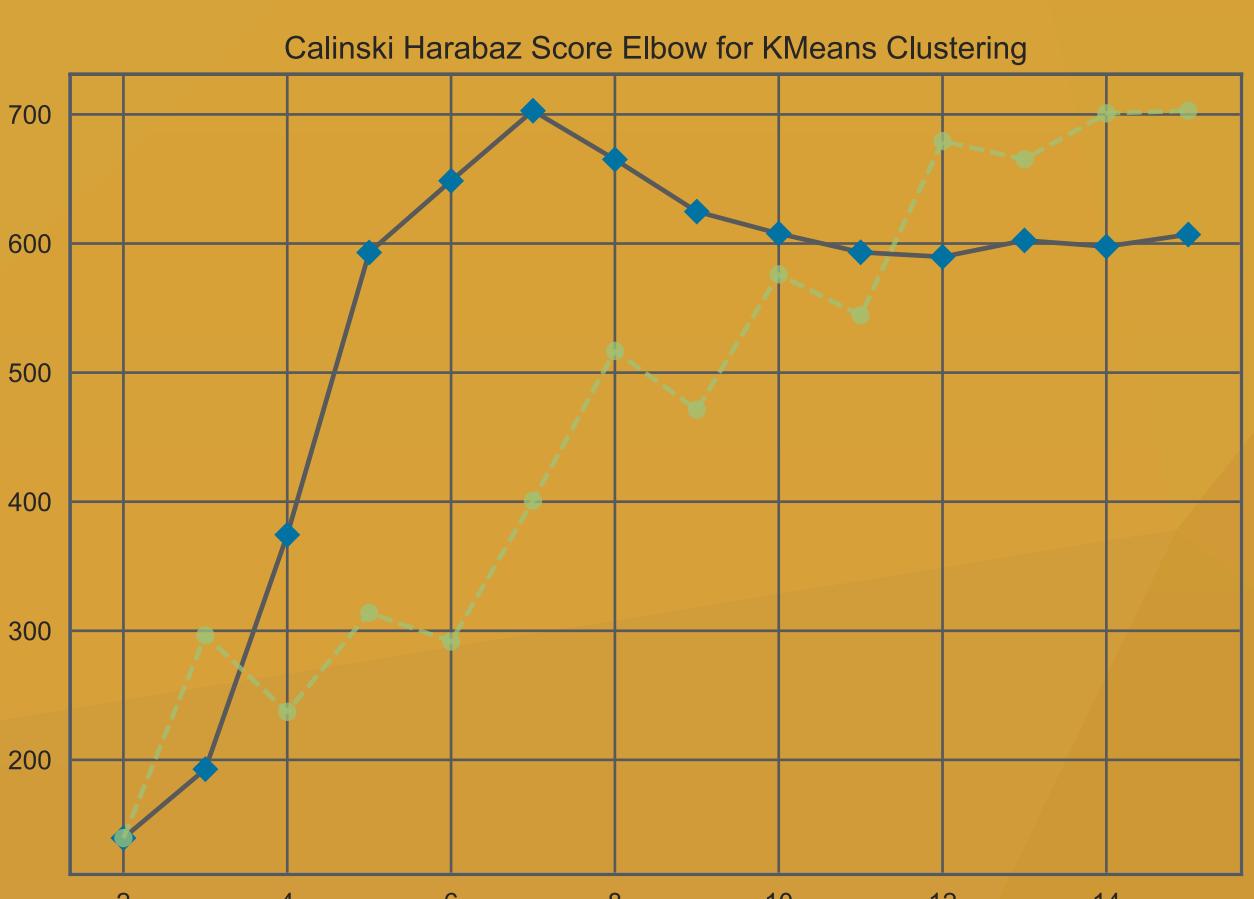


HOW DO I TUNE MY MODEL?

Elbow Curves and Silhouette Scores

- How do you pick an initial value for k in k-means clustering?
- How do you know whether to increase or decrease k?
- Is partitive clustering the right choice?

Higher silhouette scores mean denser, more separate clusters:



Alpha Selection

- Should I use Lasso, Ridge, or ElasticNet?
- Is regularization even working?

