# Salary distribution

# Salary per country



Countries Boxplot

# First approach: Regression

Using location as input and salary as target

| | $R^2$ |
|---|---|
| Decision Tree Regressor | 0.0887 |
| Random Forest Regressor | 0.1397 |

# First approach: Random Forest Classification

Using location as a binary input and salary as target

|  | Predicted Low | Predicted High |
|---|---|---|
| Low | 42 | 43 |
| High | 17 | 70 |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.63 | 0.50 | 0.56 | 78 |
| 1 | 0.65 | 0.76 | 0.70 | 94 |
| avg / total | 0.64 | 0.64 | 0.63 | 172 |

Accuracy:   0.0448

# Count Vectorizer

# Count Vectorizer: Random Forest Classification

Using title and location as a binary input and salary as target

|  | Predicted Low | Predicted High |
|---|---|---|
| Low | 89 | 6 |
| High | 23 | 54 |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.79 | 0.94 | 0.86 | 95 |
| 1 | 0.90 | 0.70 | 0.79 | 77 |
| avg / total | 0.84 | 0.83 | 0.83 | 172 |

Accuracy:   0.7729

# Count Vectorizer: Regression

|  | Accuracy |
|---|---|
| KNN (8 neighbors) | 0.1395 |
| Random Forest | 0.1881 |

# New Threshold: Random Forest Classification

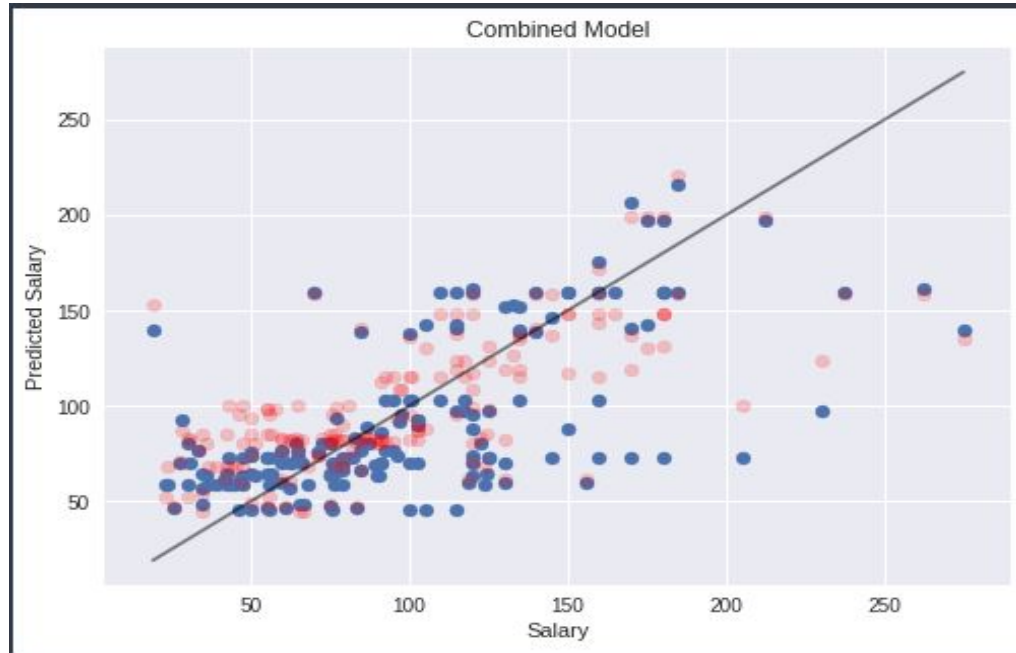Using 110 as threshold on salary to classify the target variable

| | Predicted Low | Predicted High |
|------|---------------|----------------|
| Low  | 124           | 11             |
| High | 11            | 26             |

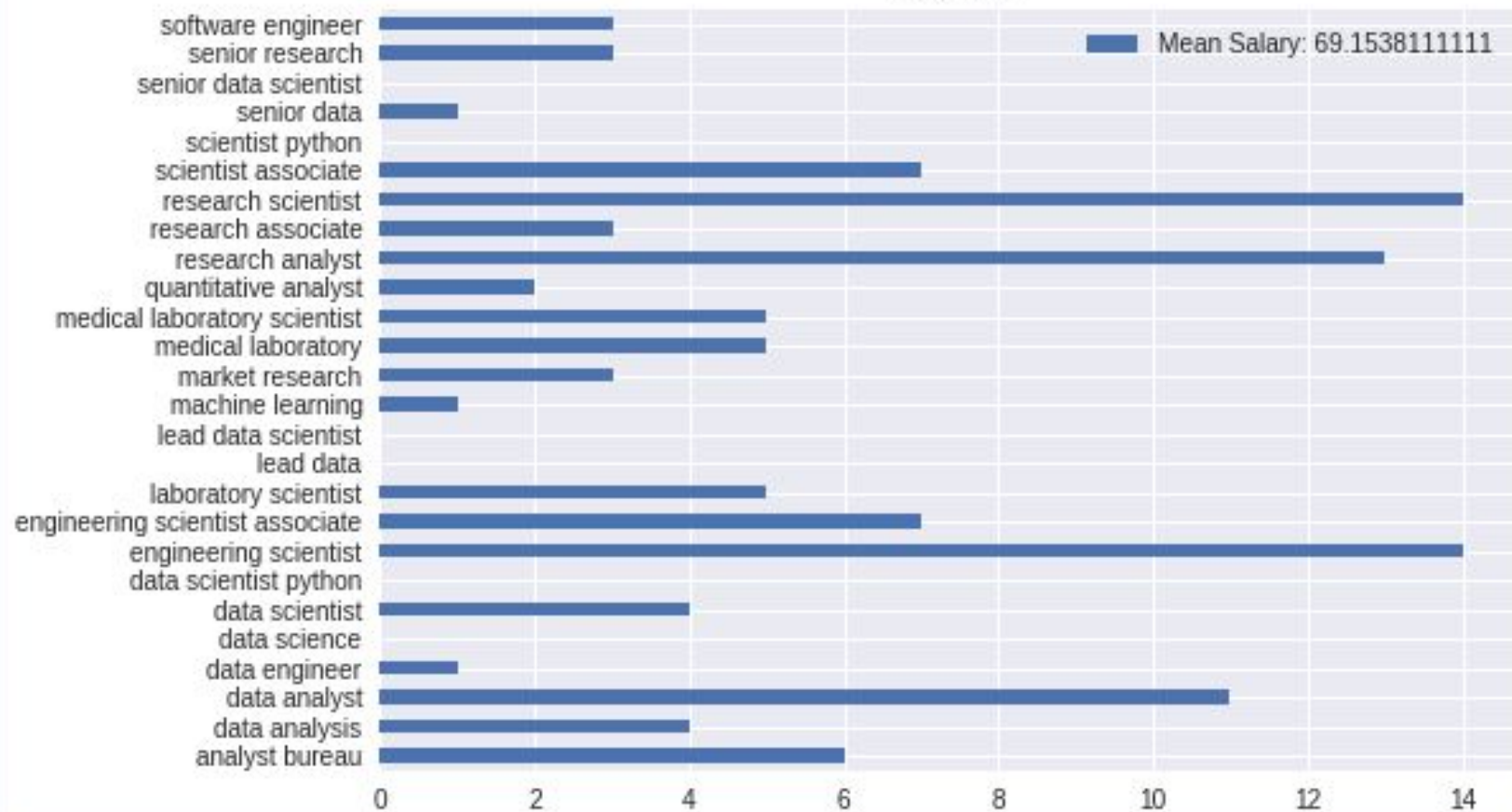|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0          | 0.92      | 0.92   | 0.92     | 135     |
| 1          | 0.70      | 0.70   | 0.70     | 37      |
| avg / total| 0.87      | 0.87   | 0.87     | 172     |

Accuracy:   0.8428

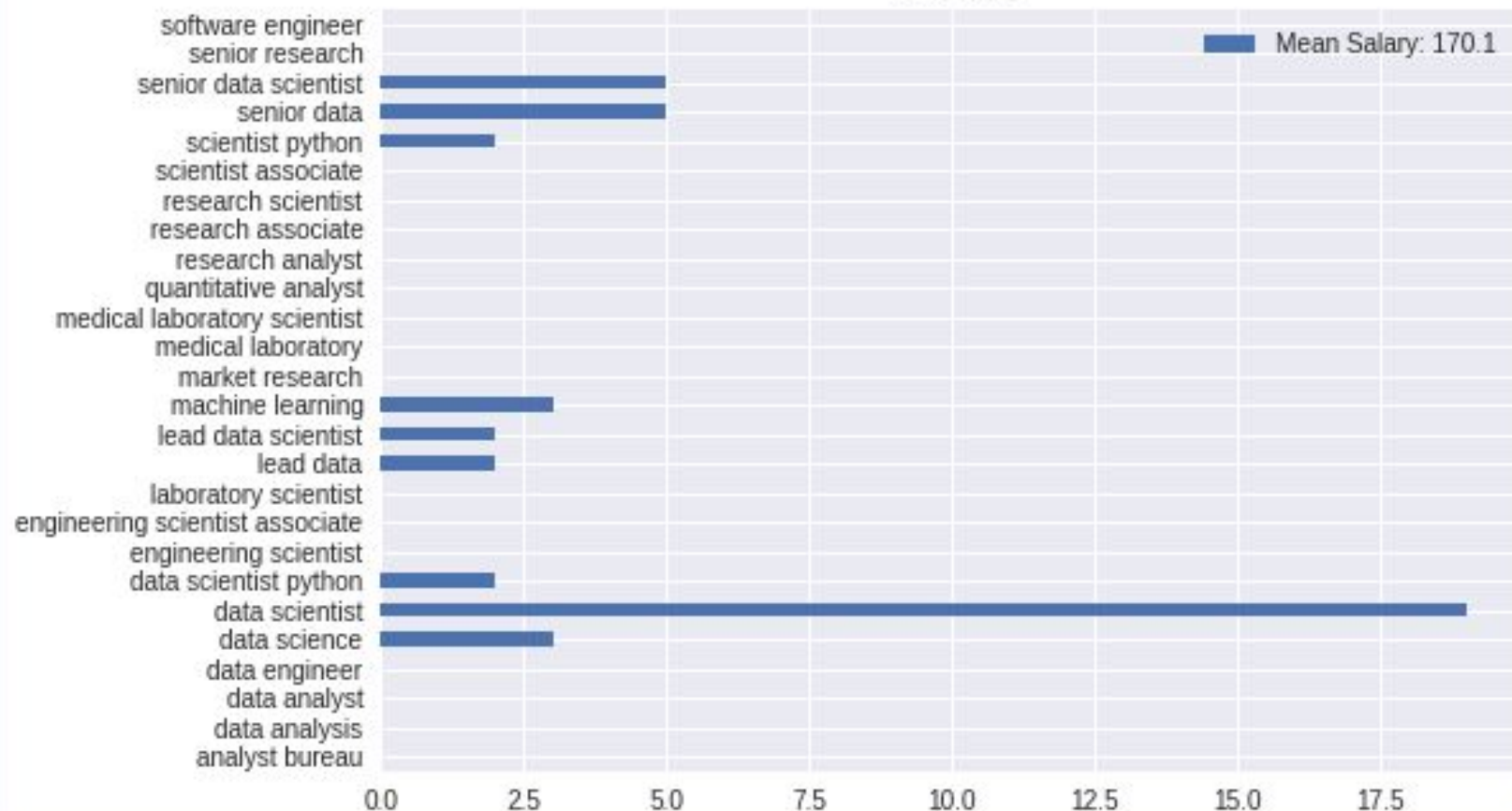# Combined Model

Classification + Random Forest Regressors R2 = 0.448

# Unsupervised: Clusters

Cluster 0

Cluster 1

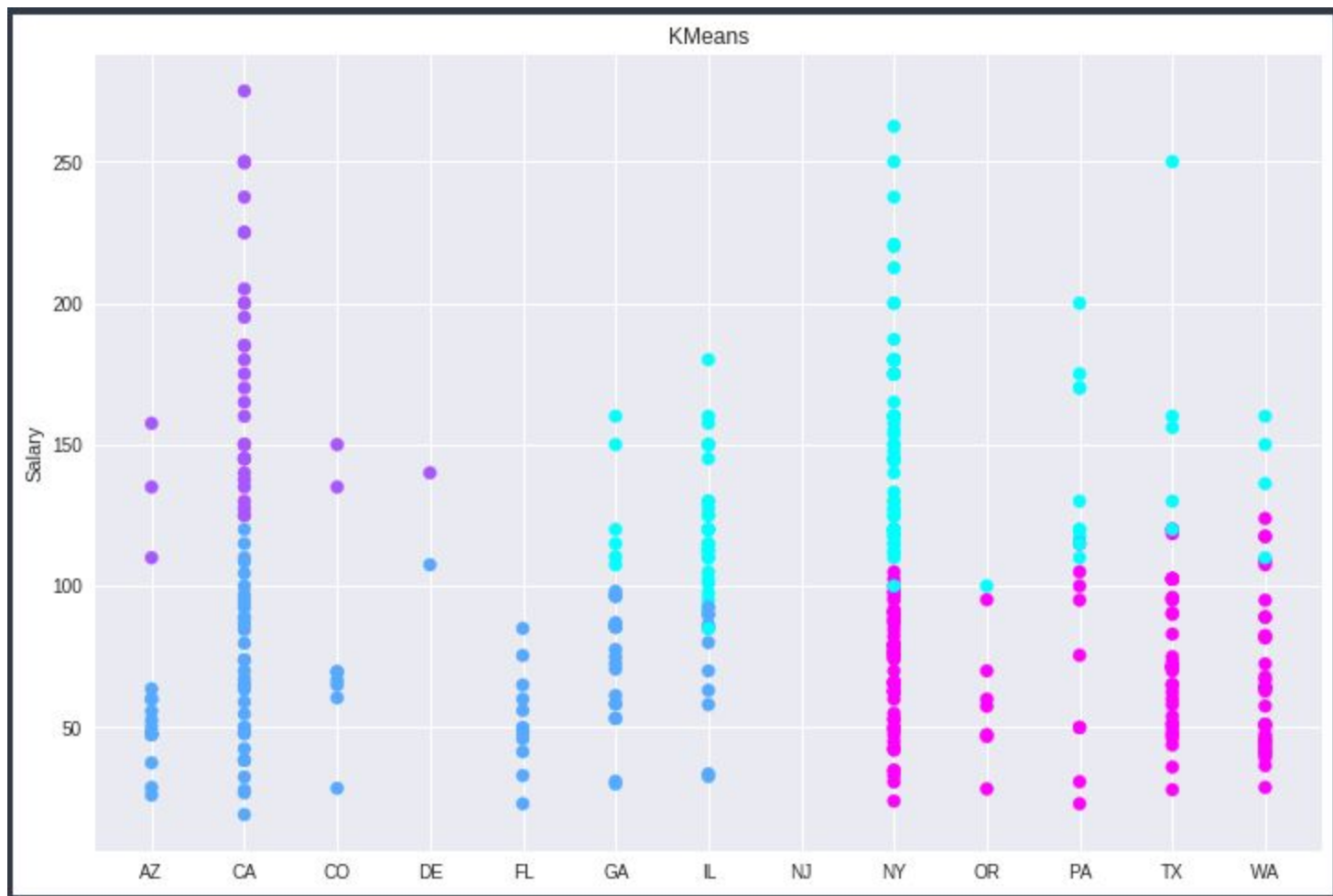Cluster 2
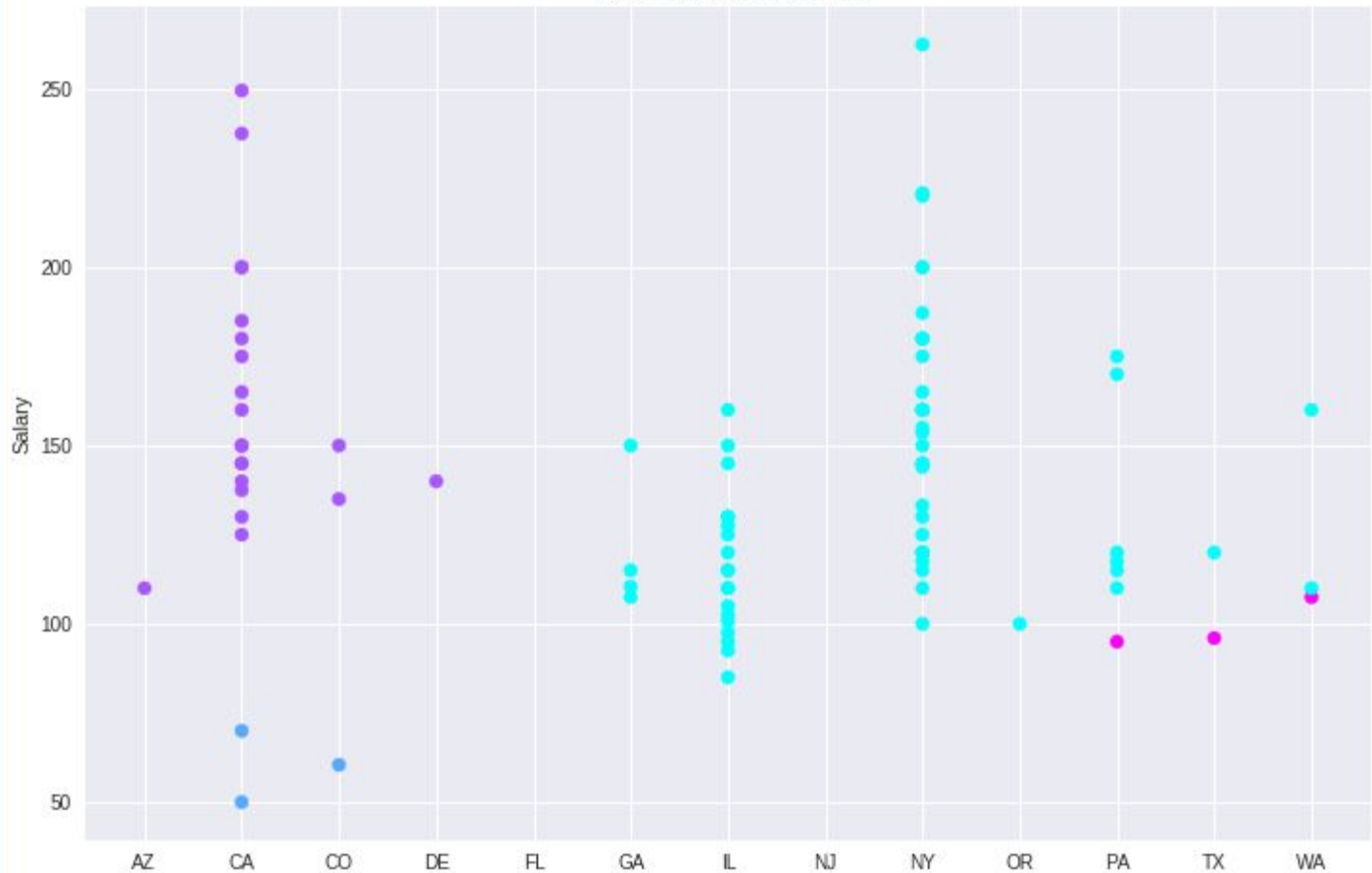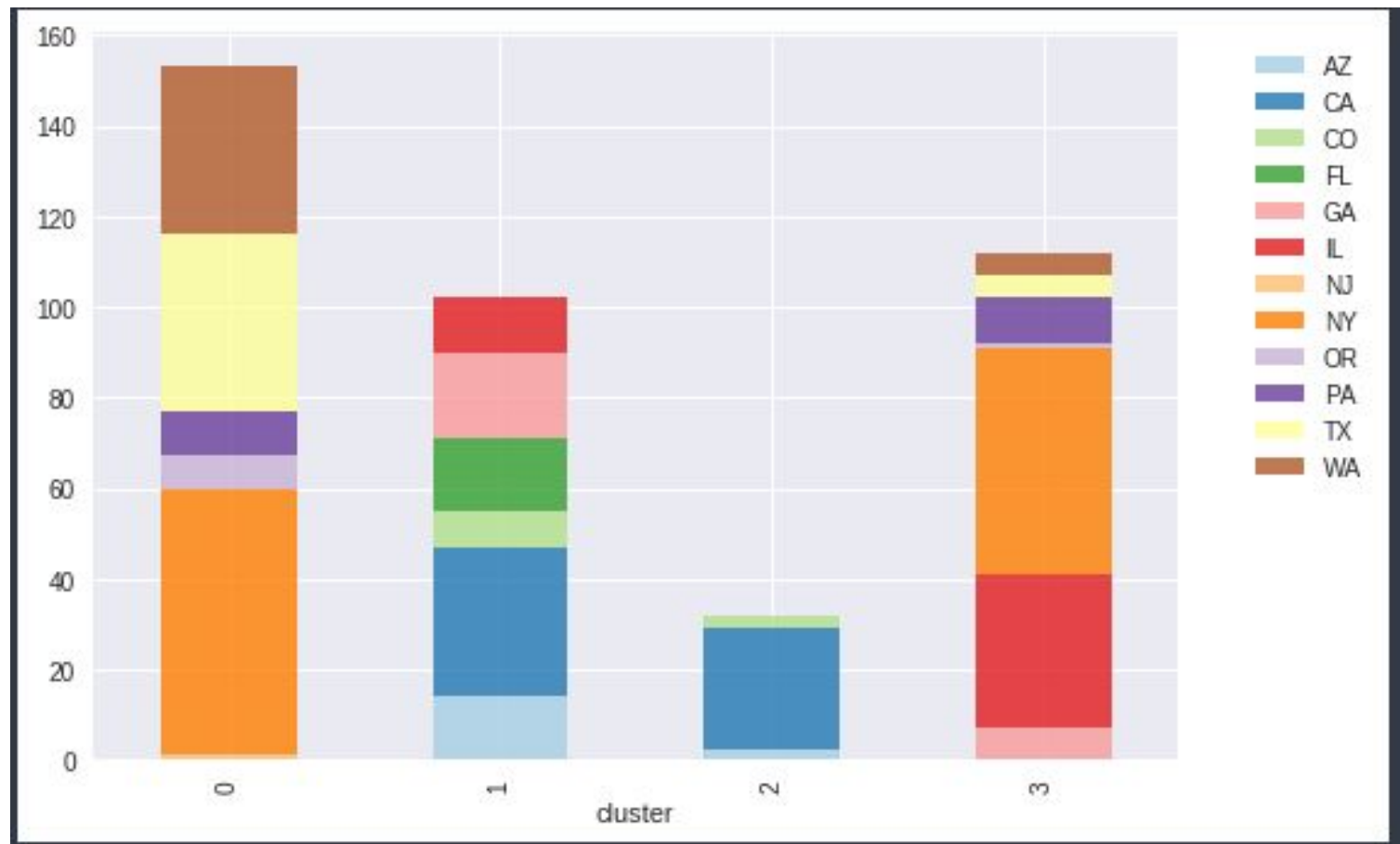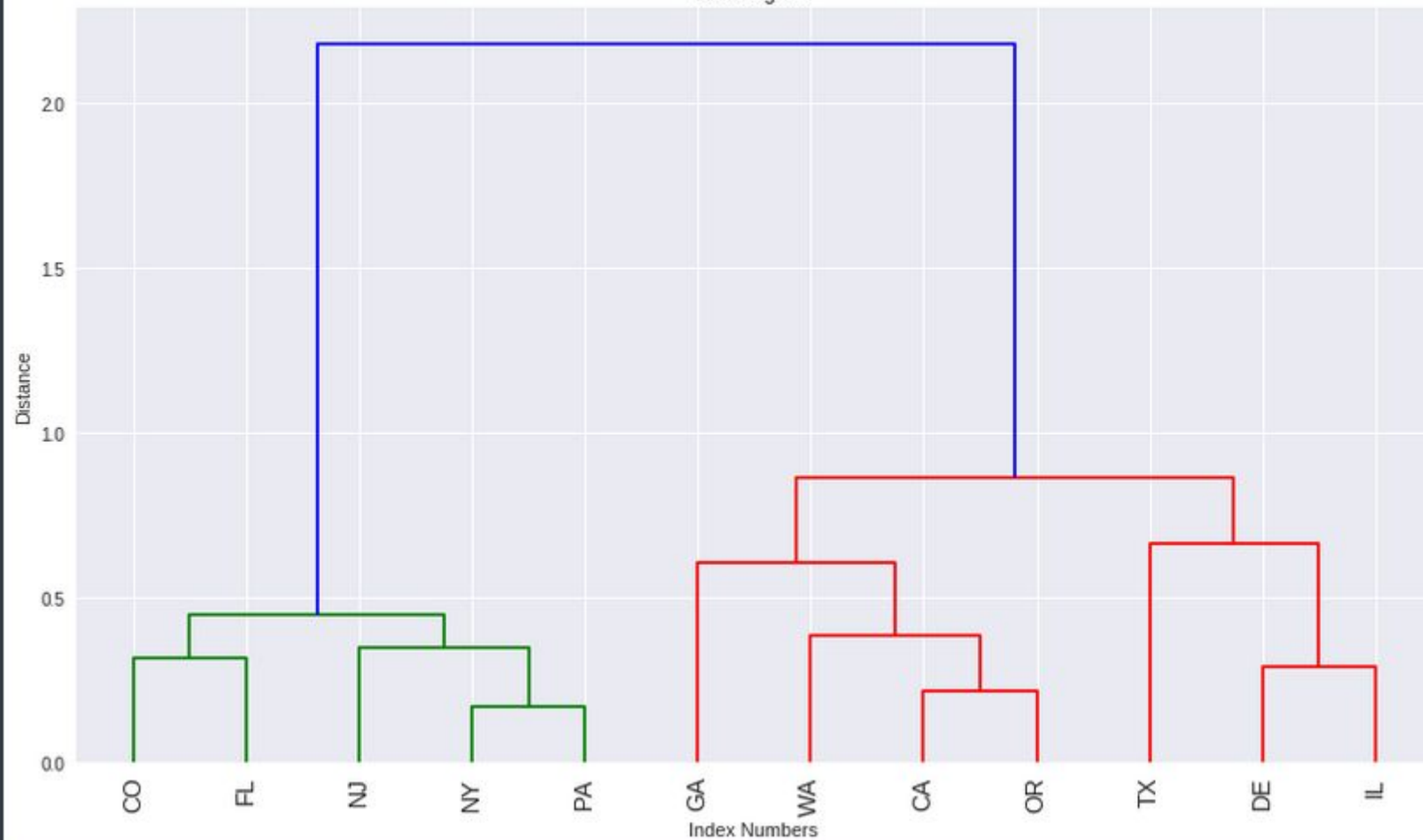
Cluster 3

Mean Salary: 140.474597321

KMeans for Data Science

Dendrogram

Questions
Questions
Questions
Questions
Questions
Questions
Questions
Questions
Goodbye