



CLUSTER ANALYSIS IN R

Occupational Wage Data

Dmitriy (Dima) Gorenshteyn

Sr. Data Scientist,

Memorial Sloan Kettering Cancer Center



Occupational Wage Data

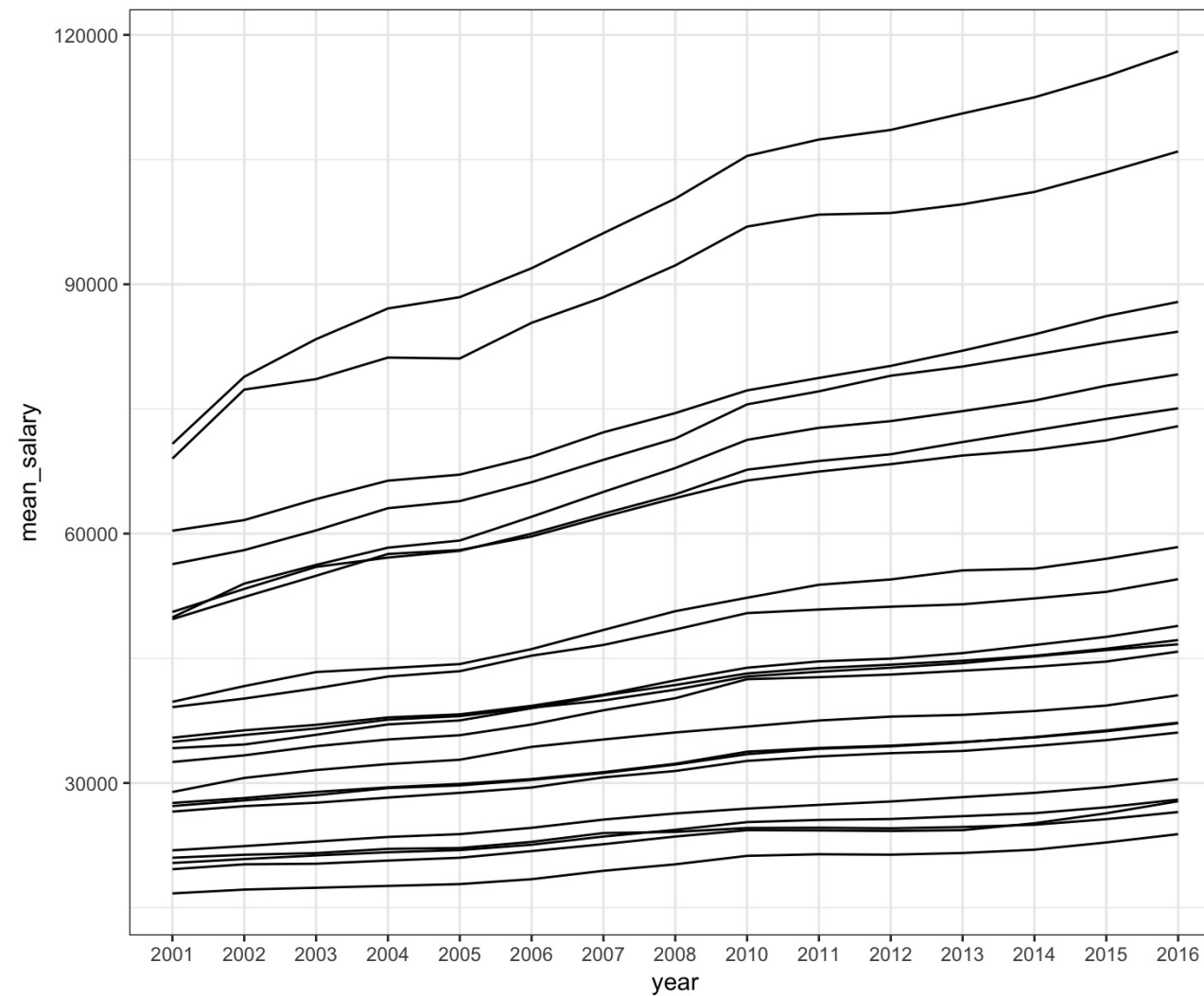
- 22 Occupation Observations
- 15 Measurements of Average Income from 2001-2016



Occupational Wage Data

```
print(oes)
      2001  2002  2003  2004  2005  ...
Management 70800 78870 83400 87090 88450 ...
Business Operations 50580 53350 56000 57120 57930 ...
Computer Science 60350 61630 64150 66370 67100 ...
Architecture/Engineering 56330 58020 60390 63060 63910 ...
Life/Physical/Social Sci. 49710 52380 54930 57550 58030 ...
Community Services 34190 34630 35800 37050 37530 ...
...      ...      ...      ...      ...      ...
```

Occupational Wage Data





Next Steps: Hierarchical Clustering

- Evaluate whether pre-processing is necessary
- Create a distance matrix
- Build a dendrogram
- Extract clusters from dendrogram
- Explore resulting clusters



CLUSTER ANALYSIS IN R

Let's practice!



CLUSTER ANALYSIS IN R

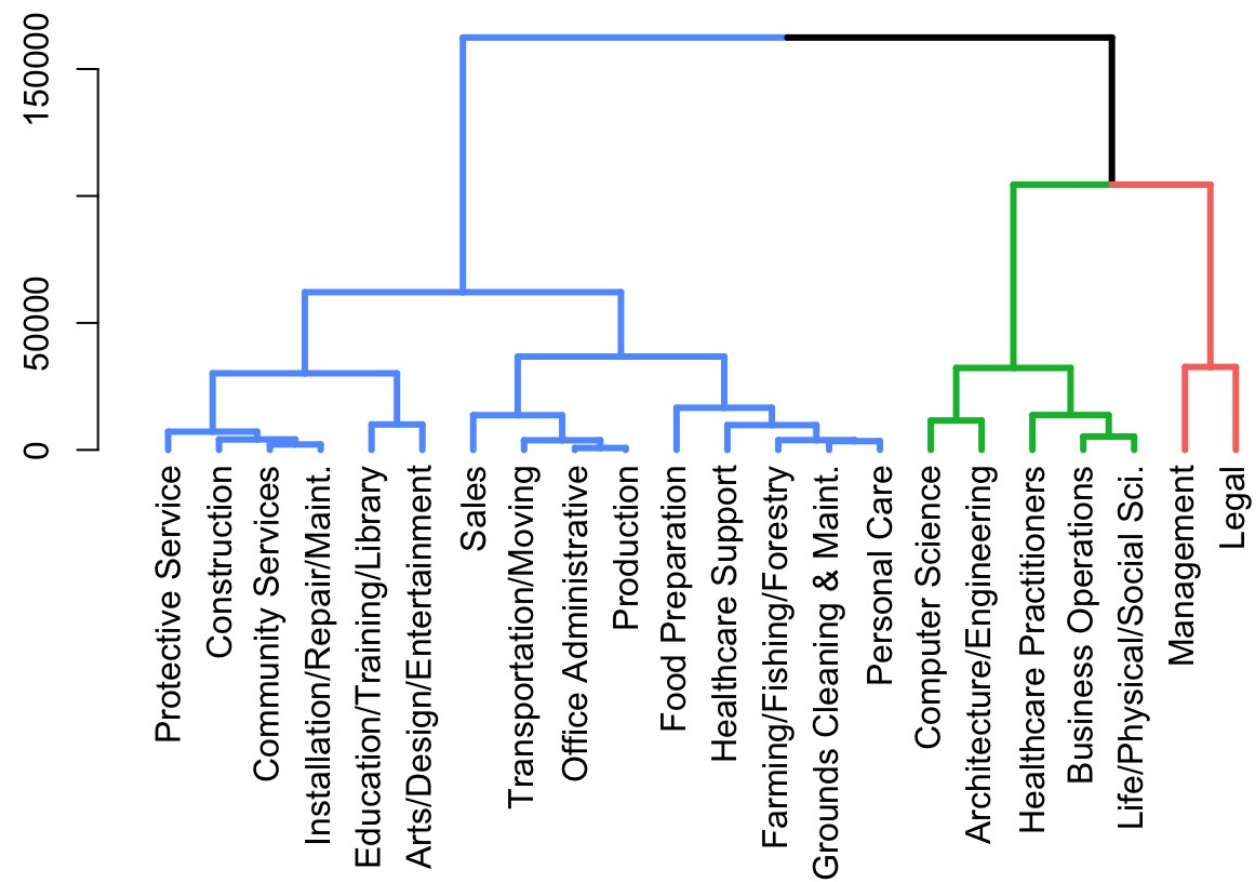
Reviewing the Hierarchical Clustering Results

Dmitriy (Dima) Gorenshteyn

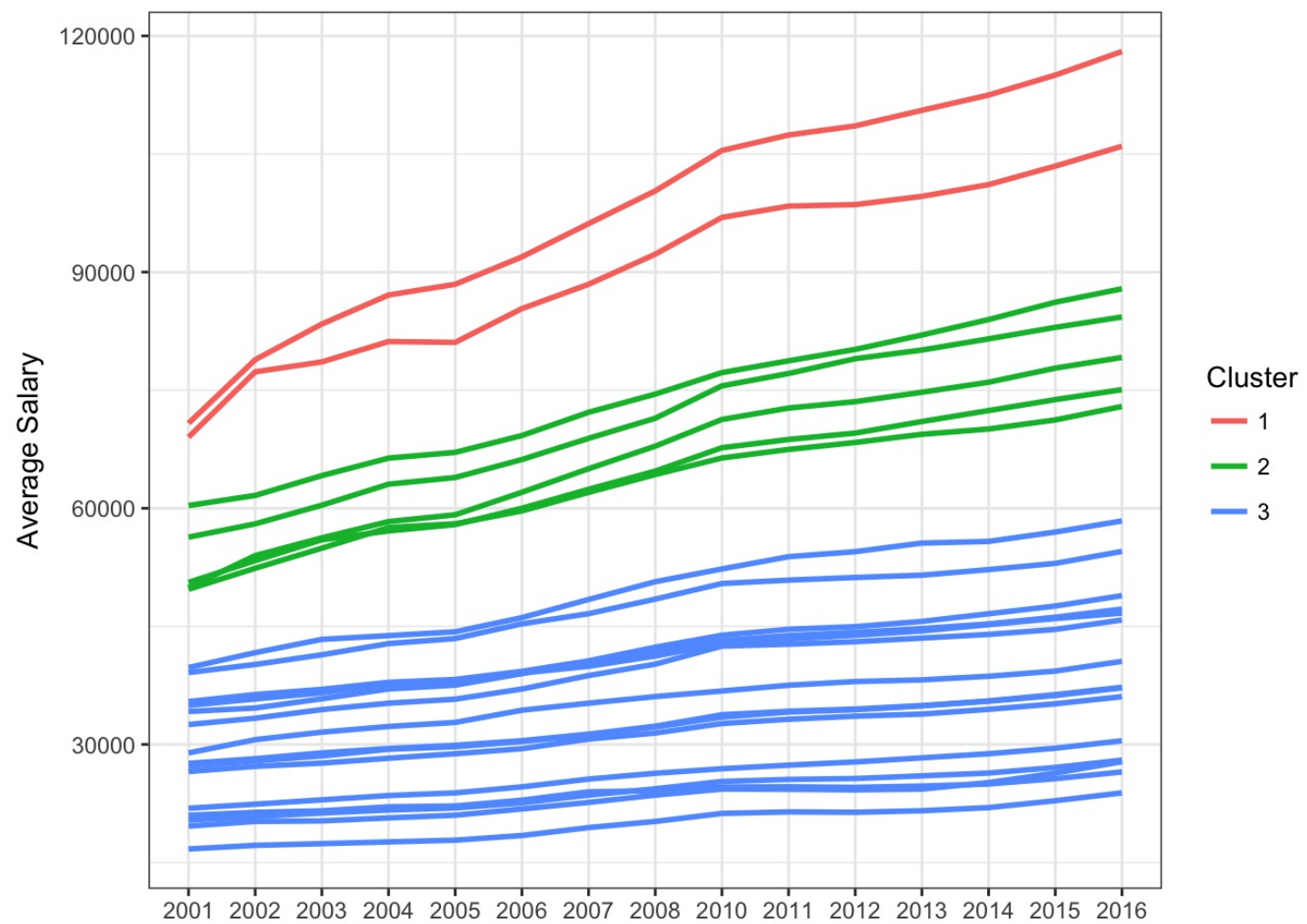
Sr. Data Scientist,

Memorial Sloan Kettering Cancer Center

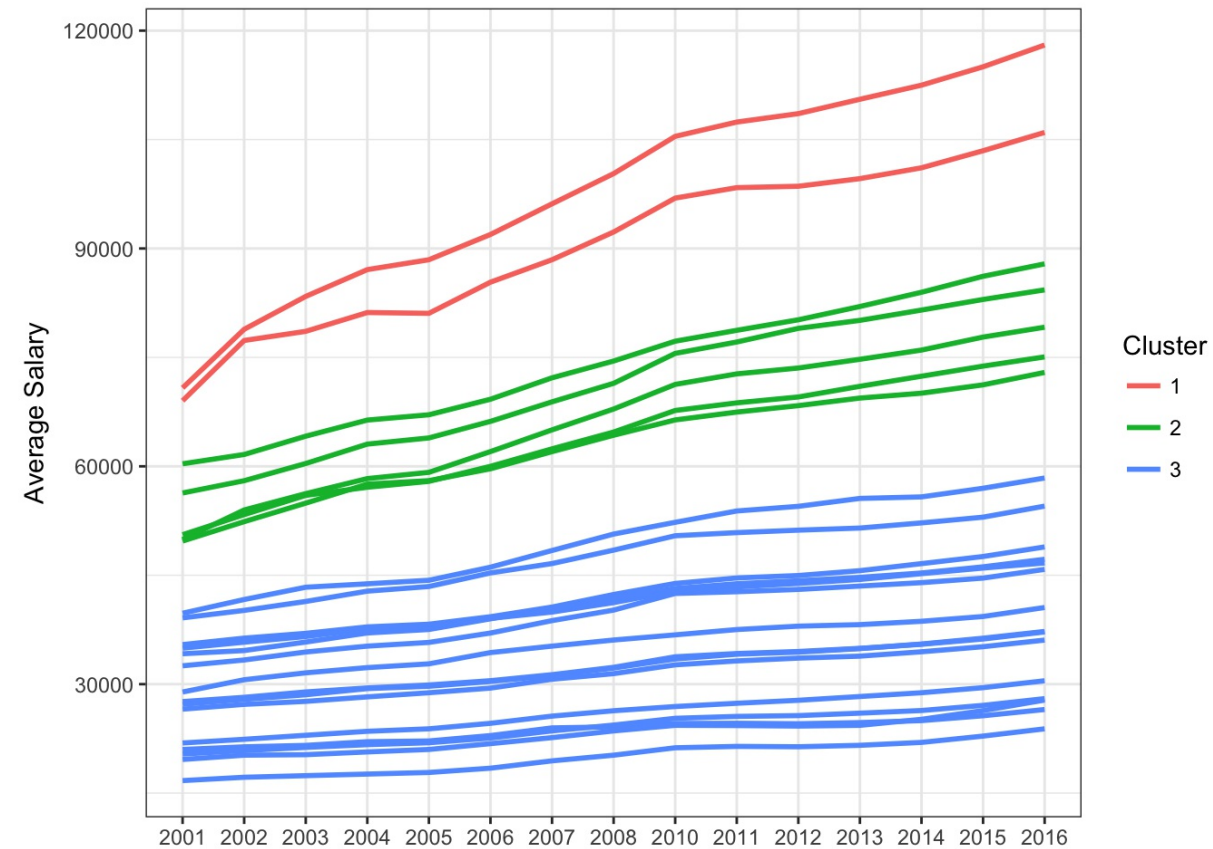
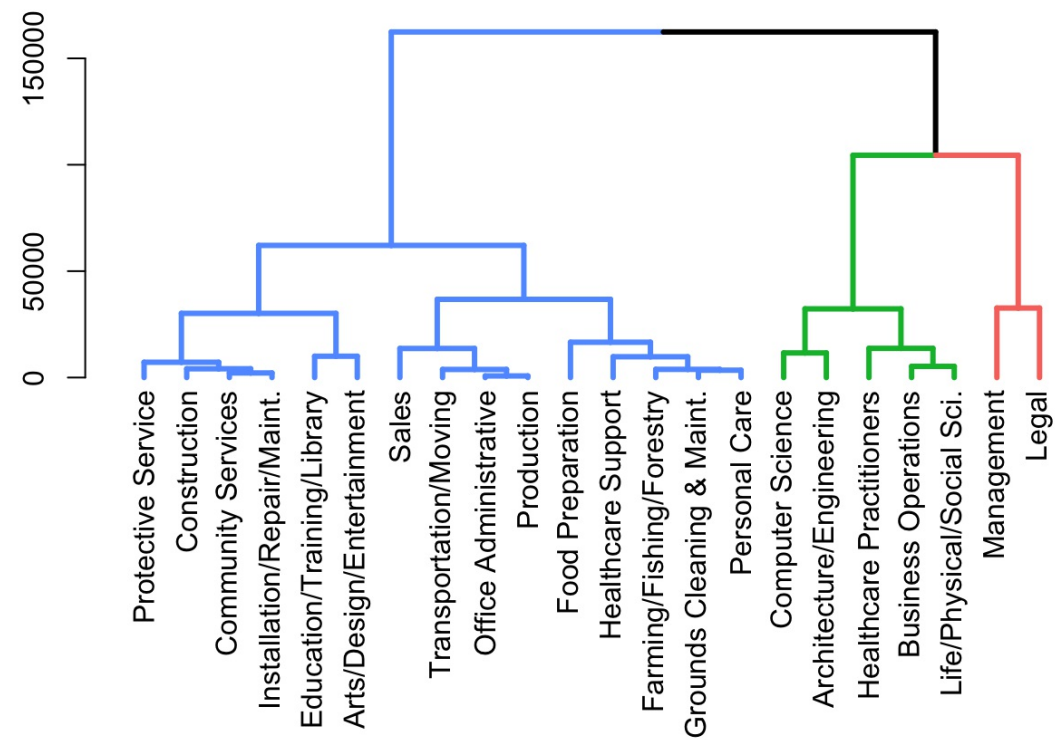
The Dendrogram



The Trends



Connecting The Two





Next Steps: k-means Clustering

- Evaluate whether pre-processing is necessary
- Estimate the "best" k using the elbow plot
- Estimate the "best" k using the maximum average silhouette width
- Explore resulting clusters



CLUSTER ANALYSIS IN R

Let's cluster!



CLUSTER ANALYSIS IN R

Review K-means Results

Dmitriy (Dima) Gorenshteyn

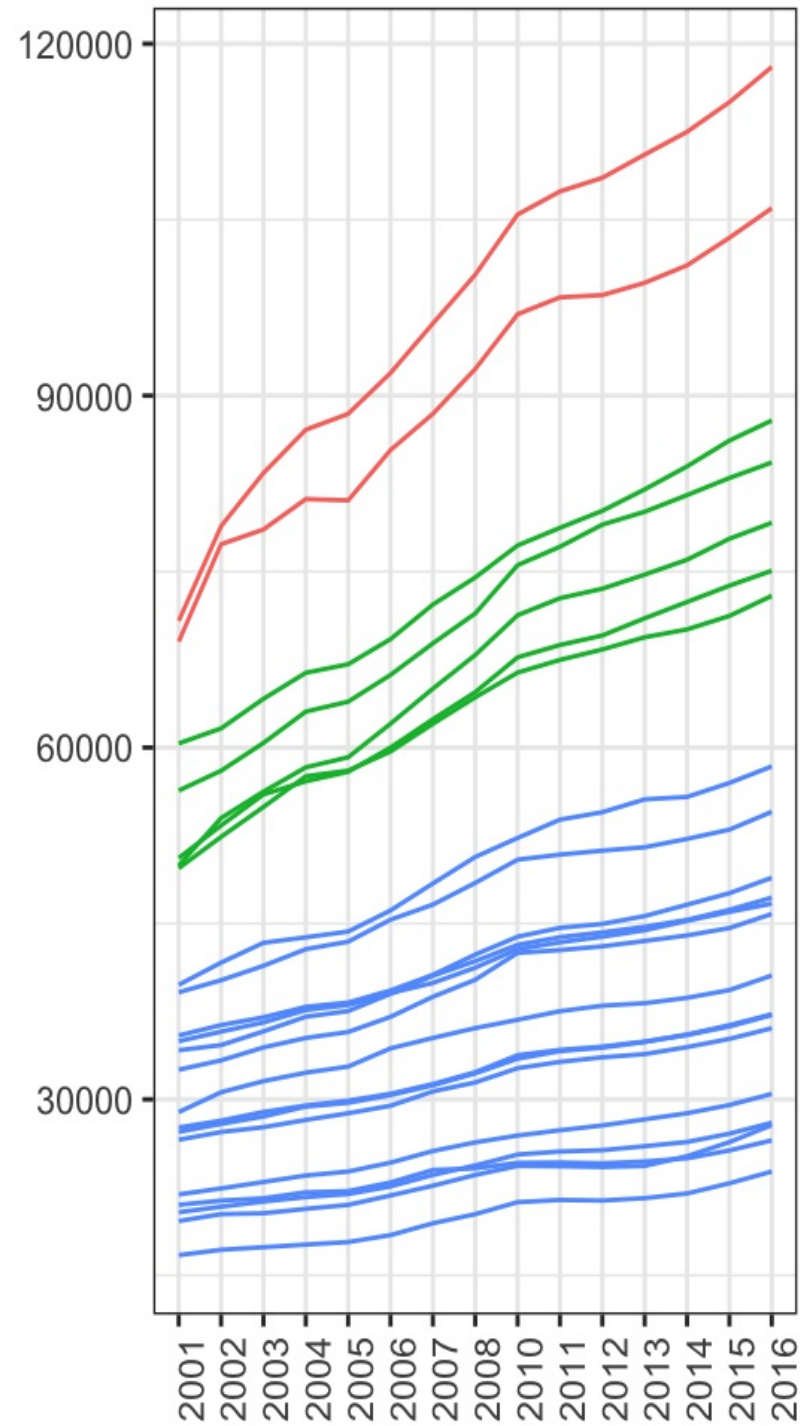
Sr. Data Scientist,

Memorial Sloan Kettering Cancer Center



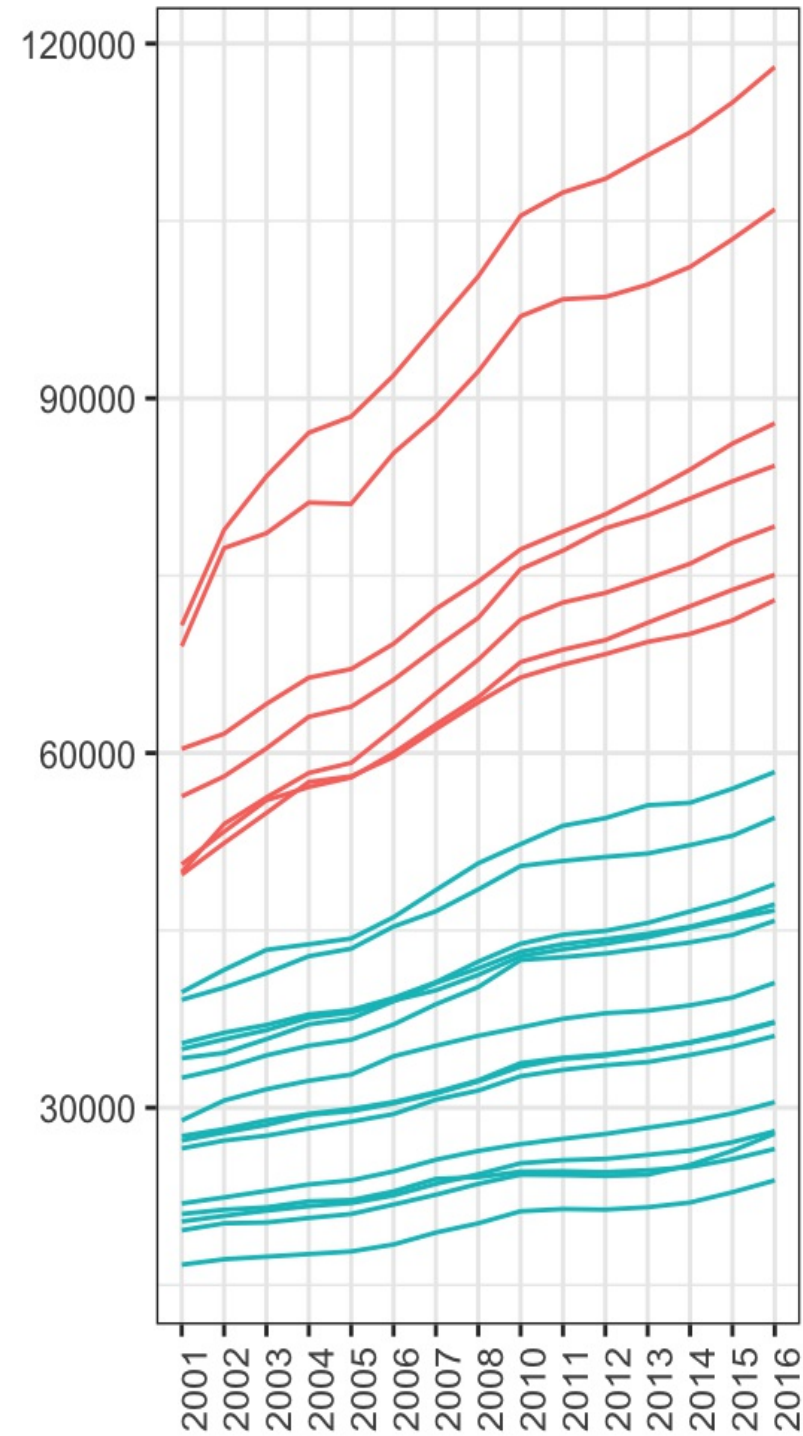
Hierarchical Clustering

Based on Dendrogram with Euclidean



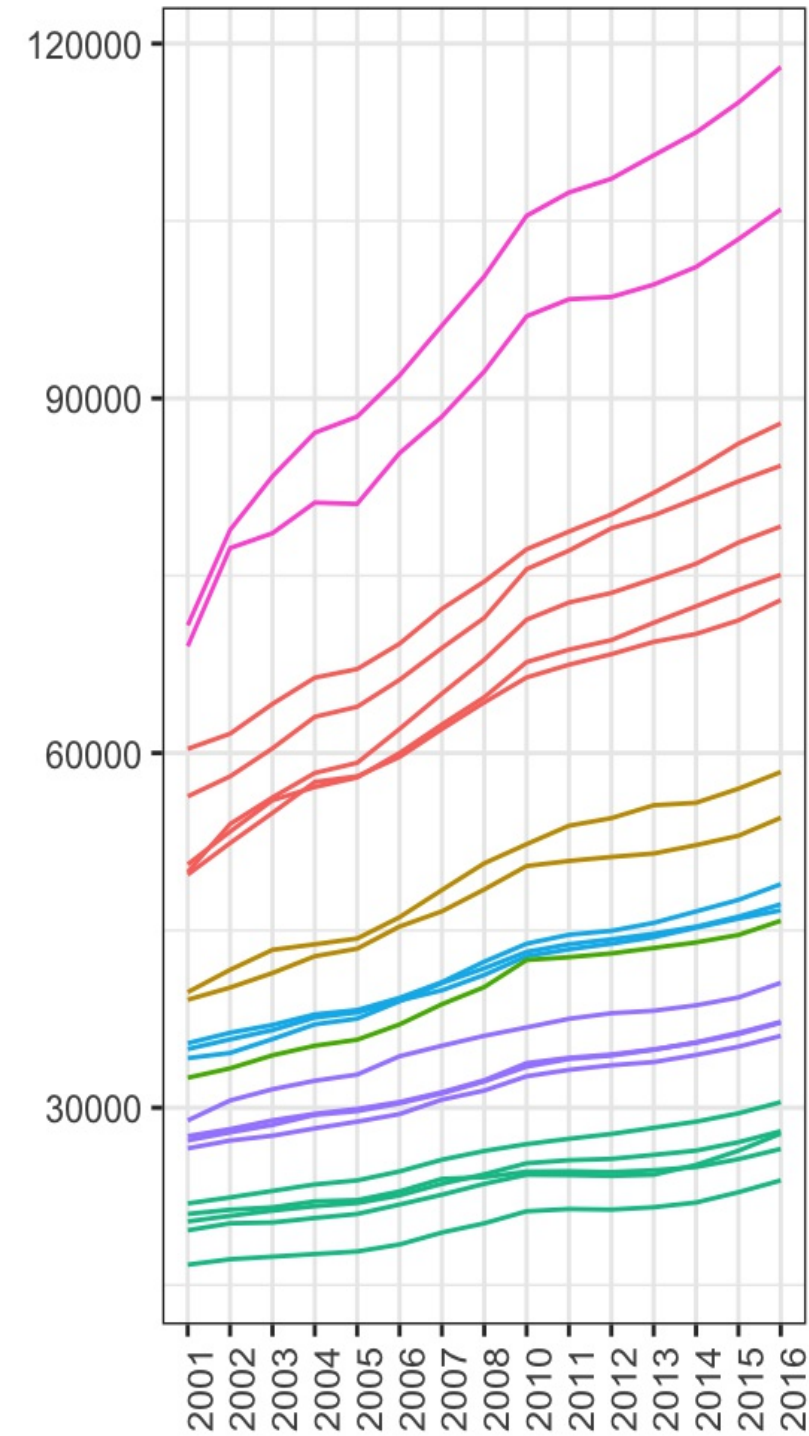
K-Means Clustering: k = 2

Based on Elbow Plot



K-Means Clustering: k = 7

Based on Silhouette Plot



Comparing The Two Clustering Methods

	Hierarchical Clustering	k-means
Distance Used:	virtually any	euclidean only
Results Stable:	Yes	No
Evaluating # of Clusters:	dendrogram, silhouette, elbow	silhouette, elbow
Computation Complexity:	Relatively Higher	Relatively Lower



What you have learned?

- **Chapter 1:**

- What is distance
- Why is scale important

- **Chapter 2:**

- How linkage works
- How the dendrogram is formed
- How to analyze your clusters

- **Chapter 3:**

- How k-means works
- How to estimate k
- How to analyze how well an observation fits in a cluster



Lot's More to Learn

- k-mediods
- DBSCAN
- Optics



CLUSTER ANALYSIS IN R

Congratulations!