# Estimating K with the Gap Statistic

Ryan Hicks

Department of Statistics, Colorado State University

December 3, 2015

# Motivation

- Minimizing within-cluster variation does not work

- $W(K) = \sum\limits_{k=1}^{K} \sum\limits_{i \in C_k} \| X_i - \overline{X_k} \|_2^2$

- Maximizing between-cluster variation does not work

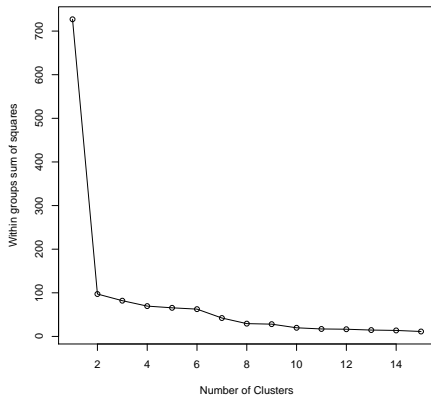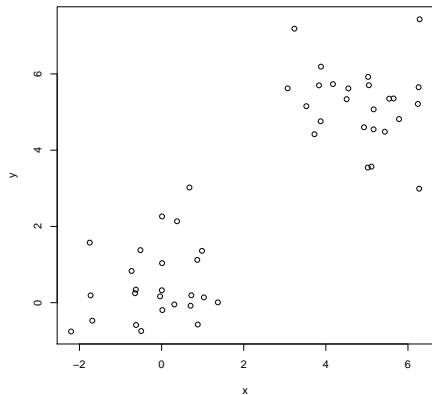- $B(K) = \sum\limits_{k=1}^{K} |C_k| \| \overline{X}_k - \overline{X} \|_2^2$

# Calinski-Harabasz Index

- $CH(K) = \frac{B(K)/(K-1)}{W(K)/(n-K)}$

- Choose a maximum number of clusters then find
  $\hat{K} = \arg\max_{K \in \{2,\ldots,K_{max}\}} CH(K)$

- However, CH(K) is undefined for K = 1; a big disadvantage.

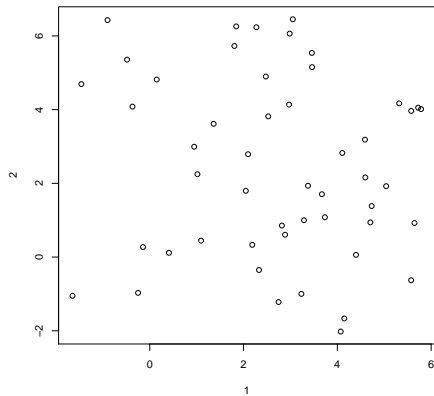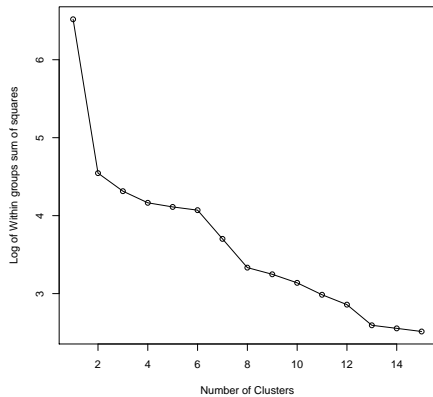- There may not be any underlying clusters in the data.

# Gap Statistic

- $\forall i, 1 \leq i \leq K_{max}$, run a clustering method on the dataset to find $i$ clusters, and sum the distance of all points from their cluster mean.

- Generate B reference datasets, easily found by uniformly sampling from a bounding rectangle of the original dataset, though there are more complex approaches.

- Define the gap statistic by $Gap_n(k) = E_n^*\{ln(W_k)\} - ln(W_k)$

- Choose the number of clusters that maximizes the gap.

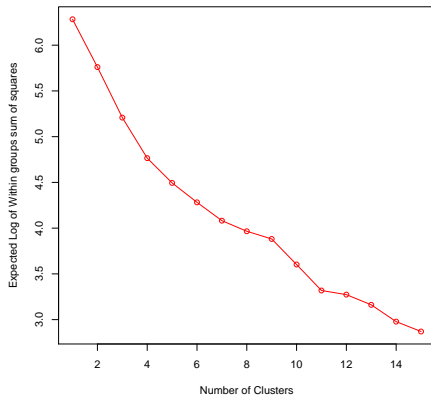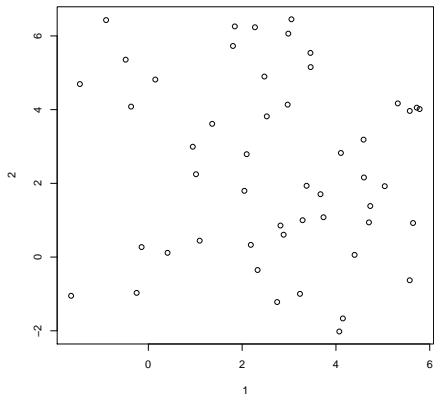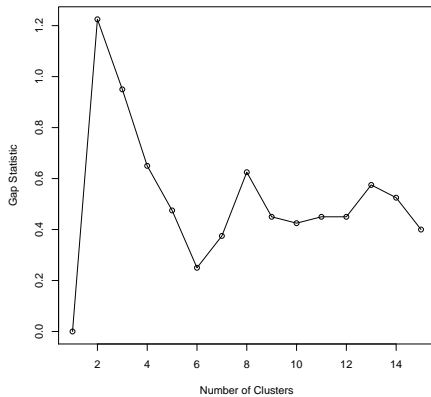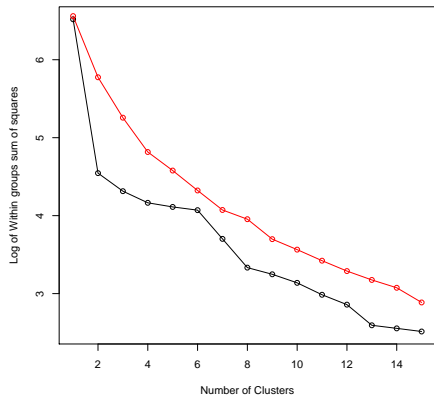# Example (Using K-means)

# Example

# Example

# Example

# The End

Thank You!

Questions?