

Feature assessment and selection using sparse clustering

Avgoustinos Vouros¹

¹PhD student,
Department of Computer Science,
University of Sheffield

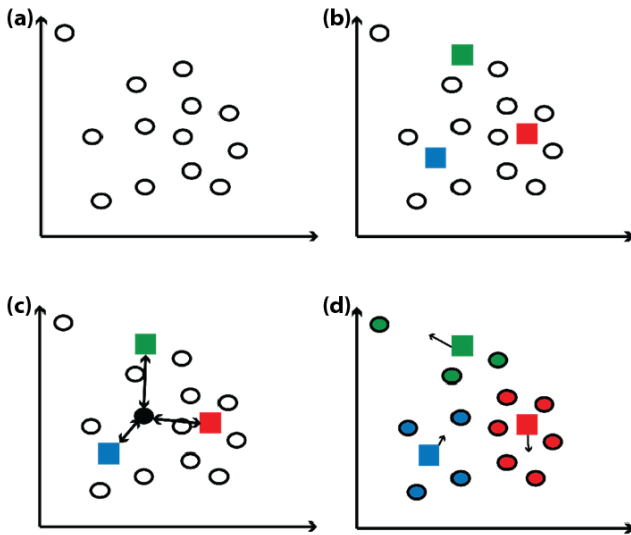
Supervised by Prof Eleni Vasilaki



- The K-Means Algorithm (Lloyd's)
- Sparse K-Means
 - Theory
 - ★ Regression (quick)
 - Algorithm
 - Tuning
 - ★ Gap Statistic
- Ongoing Research

The K-Means Algorithm (Lloyd's)

The K-Means Algorithm (Lloyd's)



[1] Lloyd, Stuart. "Least squares quantization in PCM." IEEE transactions on information theory 28.2 (1982): 129-137.

The K-Means Algorithm (Lloyd's)

Minimize:

$$WCSS = \sum_{k=1}^K \sum_{\substack{i=1 \\ x_i \in c_k}}^n \sum_{j=1}^p (x_{ij} - m_{kj})^2$$

Maximize:

$$BCSS = \sum_{j=1}^p \left(\sum_{i=1}^n (x_{ij} - M_{1j})^2 - \sum_{k=1}^K \sum_{\substack{i=1 \\ x_i \in c_k}}^n (x_{ij} - m_{kj})^2 \right)$$

The K-Means Algorithm (Lloyd's)

Advantages:

- Simple and easy to implement.
- Versatile.
- Guaranteed to converge.
- Invariant to data ordering.

Disadvantages:

- Detects only spherical and well-separated clusters.
- Sensitive to noise and outliers (Euclidean).
- Converges to a local minimum.

[1] Celebi, M. Emre, Hassan A. Kingravi, and Patricio A. Vela. "A comparative study of efficient initialization methods for the k-means clustering algorithm." Expert systems with applications 40.1 (2013): 200-210.

The K-Means Algorithm (Lloyd's)

Advantages:

- Simple and easy to implement.
- Versatile.
- Guaranteed to converge.
- Invariant to data ordering.

Disadvantages:

- Detects only spherical and well-separated clusters.
- Sensitive to noise and outliers (Euclidean).
- Converges to a local minimum.

In general:

- Sensitive to initial centroids location.

[1] Celebi, M. Emre, Hassan A. Kingravi, and Patricio A. Vela. "A comparative study of efficient initialization methods for the k-means clustering algorithm." Expert systems with applications 40.1 (2013): 200-210.

The K-Means Algorithm (Lloyd's)

Advantages:

- Simple and easy to implement.
- Versatile.
- Guaranteed to converge.
- Invariant to data ordering.

Disadvantages:

- Detects only spherical and well-separated clusters.
- Sensitive to noise and outliers (Euclidean).
- Converges to a local minimum.

In general:

- Sensitive to initial centroids location.
- Sensitive to features (variables/attributes).

[1] Celebi, M. Emre, Hassan A. Kingravi, and Patricio A. Vela. "A comparative study of efficient initialization methods for the k-means clustering algorithm." Expert systems with applications 40.1 (2013): 200-210.

Sparse K-Means

- Theory -

Sparse K-Means Theory

$$\begin{aligned} & \underset{c_1, \dots, c_k, W}{\text{maximize}} \left\{ \sum_{j=1}^p w_{jj} \left(\sum_{i=1}^n (x_{ij} - M_{1j})^2 - \sum_{k=1}^K \sum_{\substack{i=1 \\ (x_i \in c_k)}}^n (x_{ij} - m_{kj})^2 \right) \right\} \\ & \text{subject to} \quad \sum_{j=1}^p w_{jj}^2 \leq 1, \quad \sum_{j=1}^p |w_{jj}| \leq s, \quad w_{jj} \geq 0 \quad \forall j \end{aligned}$$

[1] Witten, Daniela M., and Robert Tibshirani. "A framework for feature selection in clustering." *Journal of the American Statistical Association* 105.490 (2010): 713-726.

Sparse K-Means Theory

$$\begin{aligned} & \underset{c_1, \dots, c_k, W}{\text{maximize}} \left\{ \sum_{j=1}^p w_{jj} \left(\sum_{i=1}^n (x_{ij} - M_{1j})^2 - \sum_{k=1}^K \sum_{\substack{i=1 \\ (x_i \in c_k)}}^n (x_{ij} - m_{kj})^2 \right) \right\} \\ & \text{subject to} \quad \sum_{j=1}^p w_{jj}^2 \leq 1, \quad \sum_{j=1}^p |w_{jj}| \leq s, \quad w_{jj} \geq 0 \quad \forall j \end{aligned}$$

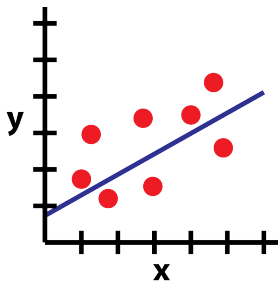
- w is a diagonal square p -by- p matrix.
- $\sum_{j=1}^p w_{jj}^2 \leq 1$ is the L_2 penalty or ridge regression ($\|w\|^2 \leq 1$) [2].
- $\sum_{j=1}^p |w_{jj}| \leq s$ is the L_1 penalty or lasso regression ($\|w\| \leq s$) [3].

[1] Witten, Daniela M., and Robert Tibshirani. "A framework for feature selection in clustering." *Journal of the American Statistical Association* 105.490 (2010): 713-726.

[2] Hoerl, Arthur E., and Robert W. Kennard. "Ridge regression: Biased estimation for nonorthogonal problems." *Technometrics* 12.1 (1970): 55-67.

[3] Tibshirani, Robert. "Regression shrinkage and selection via the lasso." *Journal of the Royal Statistical Society. Series B (Methodological)* (1996): 267-288.

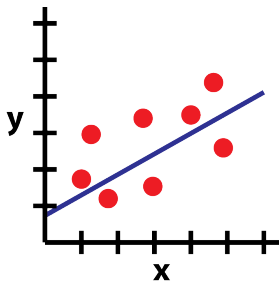
Regression



[1] StatQuest with Josh Starmer. "Regularization Part 1: Ridge Regression." YouTube. 2018. Online: <https://www.youtube.com/watch?v=Q81RR3yKn30>.

[2] StatQuest with Josh Starmer. "Regularization Part 2: Lasso Regression." YouTube. 2018. Online: <https://www.youtube.com/watch?v=NGf0voTMIcs>.

Regression

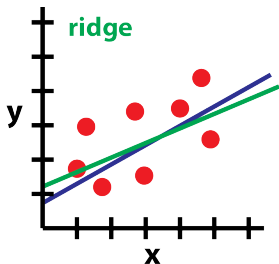


Line fitting using least squares:

$$y = y_{inter} + slope * x$$

Least squares minimizes the sum of the squared residuals.

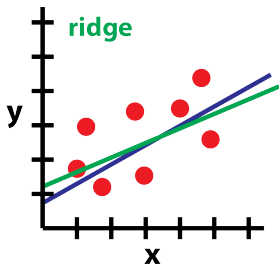
Regression



Line fitting using ridge regression:

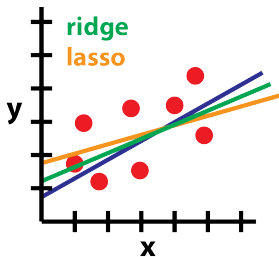
$$y = y_{inter} + slope * x + \lambda_2 * slope^2$$

Regression



Line fitting using ridge regression:

$$y = y_{inter} + slope * x + \lambda_2 * slope^2$$

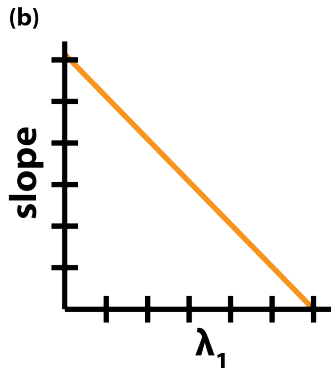
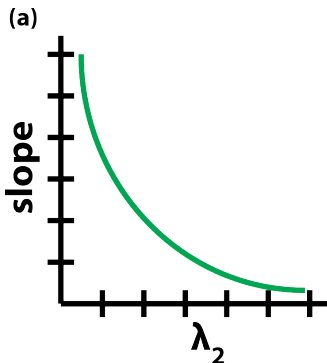


Line fitting using lasso regression:

$$y = y_{inter} + slope * x + \lambda_1 * |slope|$$

$$y = y_{inter} + slope * x + \lambda_2 * slope^2 \quad (a)$$

$$y = y_{inter} + slope * x + \lambda_1 * |slope| \quad (b)$$



Sparse K-Means Theory

$$\begin{aligned} & \underset{c_1, \dots, c_k, w}{\text{maximize}} \left\{ \sum_{j=1}^p w_{jj} \left(\sum_{i=1}^n (x_{ij} - M_{1j})^2 - \sum_{k=1}^K \sum_{\substack{i=1 \\ (x_{ij} \in c_k)}}^n (x_{ij} - m_{kj})^2 \right) \right\} \\ & \text{subject to} \quad \sum_{j=1}^p w_{jj}^2 \leq 1, \quad \sum_{j=1}^p |w_{jj}| \leq s, \quad w_{jj} \geq 0 \quad \forall j \end{aligned}$$

Sparse K-Means Theory

$$\underset{c_1, \dots, c_k, w}{\text{maximize}} \left\{ \sum_{j=1}^p w_{jj} \left(\sum_{i=1}^n (x_{ij} - M_{1j})^2 - \sum_{k=1}^K \sum_{\substack{i=1 \\ x_{ij} \in c_k}}^n (x_{ij} - m_{kj})^2 \right) \right\}$$

$$\text{subject to} \quad \sum_{j=1}^p w_{jj}^2 \leq 1, \quad \sum_{j=1}^p |w_{jj}| \leq s, \quad w_{jj} \geq 0 \quad \forall j$$

$$\underset{w}{\text{maximize}} \left\{ \sum_{j=1}^p w_{jj} a_j + \lambda \sum_{j=1}^p w_{jj}^2 + \delta \sum_{j=1}^p |w_{jj}| \right\}$$

Sparse K-Means Theory

$$\underset{w}{\text{maximize}} \left\{ \sum_{j=1}^p w_{jj} a_j + \lambda \sum_{j=1}^p w_{jj}^2 + \delta \sum_{j=1}^p |w_{jj}| \right\}$$

- $\lambda \sum_{j=1}^p w_{jj}^2$ and $\delta \sum_{j=1}^p |w_{jj}|$ are called Lagrange multipliers.
- When the constraints are having both equalities and inequalities we extend to Karush-Kuhn-Tucker (KKT) conditions.

Sparse K-Means Theory

If $w \neq 0$

$$\frac{\partial}{\partial w}|w| = \frac{\partial}{\partial w}\sqrt{w^2} = \frac{\partial}{\partial w}(w^2)^{\frac{1}{2}} = \frac{1}{2}(w^2)^{\frac{1}{2}} \cdot 2w = \frac{w}{\sqrt{w^2}} = \frac{w}{|w|}$$

else if $w = 0$

$$\frac{\partial}{\partial w}|w| = \lim_{w \rightarrow 0} \frac{|w| - 0}{w - 0} = \begin{cases} \lim_{w \rightarrow 0^+} \frac{w}{w} & , w > 0 \\ \lim_{w \rightarrow 0^-} \frac{-w}{w} & , w < 0 \end{cases} = \begin{cases} 1 & , w > 0 \\ -1 & , w < 0 \end{cases}$$

[1] proofwiki.org. "Derivative of Absolute Value Function."

wiki. 2018. Online:

https://proofwiki.org/wiki/Derivative_of_Absolute_Value_Function.

Sparse K-Means Theory

If $w \neq 0$

$$\frac{\partial}{\partial w}|w| = \frac{\partial}{\partial w}\sqrt{w^2} = \frac{\partial}{\partial w}(w^2)^{\frac{1}{2}} = \frac{1}{2}(w^2)^{\frac{1}{2}} \cdot 2w = \frac{w}{\sqrt{w^2}} = \frac{w}{|w|}$$

else if $w = 0$

$$\frac{\partial}{\partial w}|w| = \lim_{w \rightarrow 0} \frac{|w| - 0}{w - 0} = \begin{cases} \lim_{w \rightarrow 0^+} \frac{w}{w} & , w > 0 \\ \lim_{w \rightarrow 0^-} \frac{-w}{w} & , w < 0 \end{cases} = \begin{cases} 1 & , w > 0 \\ -1 & , w < 0 \end{cases}$$

[1] proofwiki.org. "Derivative of Absolute Value Function."

wiki. 2018. Online:

https://proofwiki.org/wiki/Derivative_of_Absolute_Value_Function.



Sparse K-Means Theory

Proposition: The solution to this convex problem is,

$$w_{jj} = \frac{(\text{sign}(a_j)(|a_j| - \delta))_+}{(\text{sign}(a_j)(|a_j| - \delta))_+^2}$$

where the $+$ subscript indicates the positive part of the function, $\delta = 0$ if that results in $\sum_{j=1}^p |w_{jj}| < s$ or $\delta > 0$ is chosen so that $\sum_{j=1}^p |w_{jj}| = s$ and it is assumed that $1 \leq s \leq \sqrt{p}$.

[1] Witten, Daniela M., Robert Tibshirani, and Trevor Hastie. "A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis." *Biostatistics* 10.3 (2009): 515-534.

Sparse K-Means Theory

Proposition: The solution to this convex problem is,

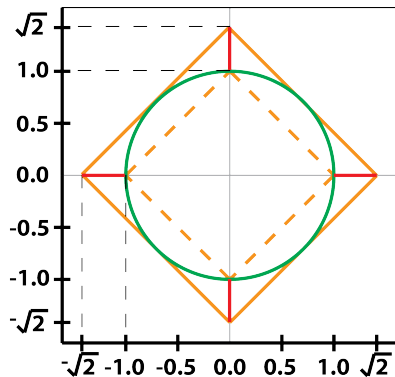
$$w_{jj} = \frac{(\text{sign}(a_j)(|a_j| - \delta))_+}{(\text{sign}(a_j)(|a_j| - \delta))_+^2}$$

where the $+$ subscript indicates the positive part of the function, $\delta = 0$ if that results in $\sum_{j=1}^p |w_{jj}| < s$ or $\delta > 0$ is chosen so that $\sum_{j=1}^p |w_{jj}| = s$ and it is assumed that $1 \leq s \leq \sqrt{p}$.

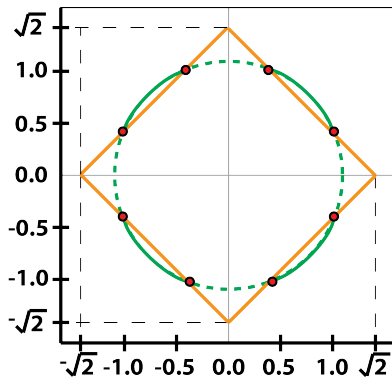
How do we find δ ?

[1] Witten, Daniela M., Robert Tibshirani, and Trevor Hastie. "A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis." *Biostatistics* 10.3 (2009): 515-534.

Sparse K-Means Theory



- $\|w\|^2 \leq 1$
- s : L1, L2 active
- $|w|_1 = 1.0$
- $|w|_1 = \sqrt{2}$



- for some s between 1.0 and $\sqrt{2}$*
- $\delta = 0$
 - L1, L2 active

[1] Witten, Daniela M., Robert Tibshirani, and Trevor Hastie. "A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis." *Biostatistics* 10.3 (2009): 515-534.

Sparse K-Means Theory

Using the Binary Search algorithm:

- Let δ be between lim_1 and lim_2
- $lim_1 = 0, lim_2 = \max(a)$
- Iterate...

$$u = \frac{\sum_{j=1}^p (\text{sign}(a_j)(|a_j| - \frac{lim_1 + lim_2}{2}))}{\sum_{j=1}^p (\text{sign}(a_j)(|a_j| - \frac{lim_1 + lim_2}{2}))^2}$$

$$\left\{ \begin{array}{ll} lim_2 = \frac{lim_1 + lim_2}{2} & , u < s \\ lim_1 = \frac{lim_1 + lim_2}{2} & , u \geq s \end{array} \right\}$$

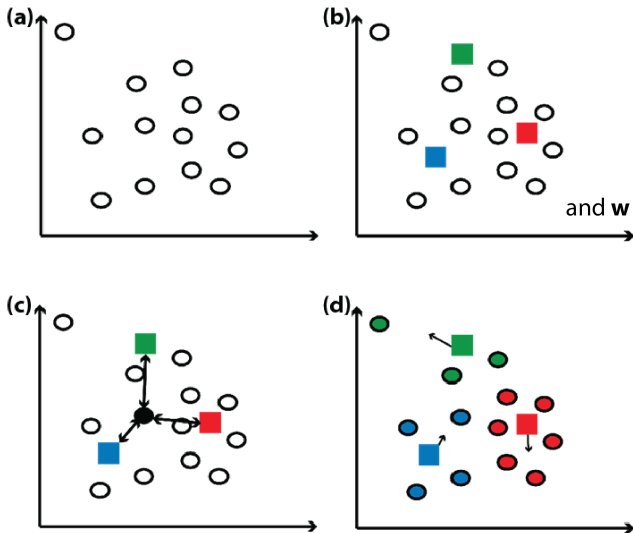
- $\delta = \frac{lim_1 + lim_2}{2}$

Sparse K-Means

- Algorithm -

Sparse K-Means Algorithm

Phase 1



Sparse K-Means Algorithm

Phase 2

- Execute Binary Search to find δ
i.e. $\delta = \text{BinarySearch}(wBCSS, s)$
- Compute all the new weights using

$$w_{jj} = \frac{(\text{sign}(a_j)(|a_j| - \delta))_+}{(\text{sign}(a_j)(|a_j| - \delta))_+^2}$$

Phase 3

- Update dataset based on w .
- Repeat from **Phase 1c** until convergence.

Sparse K-Means

- Tuning -

Sparse K-Means Tuning

How do we decide k and s ?

- Normally for k we use a performance index. But...

1	Internal clustering criteria	3
1.1	Algebraic background and notations	3
1.1.1	Total dispersion	3
1.1.2	Within-group scatter	4
1.1.3	Between-group scatter	6
1.1.4	Pairs of points	6
1.2	Internal indices	7
1.2.1	The Ball-Hall index	7
1.2.2	The Banfield-Raftery index	9
1.2.3	The C index	9
1.2.4	The Calinski-Harabasz index	9
1.2.5	The Davies-Bouldin index	10
1.2.6	The Det_Ratio index	10
1.2.7	The Dunn index	10
1.2.8	The Baker-Hubert Gamma index	11
1.2.9	The GDI index	12
1.2.10	The G-plus index	13
1.2.11	The Ksq_DetW index	13
1.2.12	The Log_Det_Ratio index	13
1.2.13	The Log_SS_Ratio index	13
1.2.14	The McClain-Rao index	13
1.2.15	The PBM index	14
1.2.16	The Point-Biserial index	14
1.2.17	The Ratkowsky-Lance index	15
1.2.18	The Ray-Turi index	16
1.2.19	The Scott-Symons index	16
1.2.20	The SD index	16
1.2.21	The S_Dbw index	17
1.2.22	The Silhouette index	18
1.2.23	The Tau index	19
1.2.24	The Trace_W index	19
1.2.25	The Trace_WiB index	20

[1] Desgraupes, Bernard. "Clustering indices." University of Paris Ouest-Lab Modal'X 1 (2013): 34.

Sparse K-Means Tuning

How do we decide k and s ?

- Normally for k we use a performance index. But...

1	Internal clustering criteria	3
1.1	Algebraic background and notations	3
1.1.1	Total dispersion	3
1.1.2	Within-group scatter	4
1.1.3	Between-group scatter	6
1.1.4	Pairs of points	6
1.2	Internal indices	7
1.2.1	The Ball-Hall index	7
1.2.2	The Banfield-Raftery index	9
1.2.3	The C index	9
1.2.4	The Calinski-Harabasz index	9
1.2.5	The Davies-Bouldin index	10
1.2.6	The Det_Ratio index	10
1.2.7	The Dunn index	10
1.2.8	The Baker-Hubert Gamma index	11
1.2.9	The GDI index	12
1.2.10	The G-plus index	13
1.2.11	The Ksq_DetW index	13
1.2.12	The Log_Det_Ratio index	13
1.2.13	The Log_SS_Ratio index	13
1.2.14	The McClain-Rao index	13
1.2.15	The PBM index	14
1.2.16	The Point-Biserial index	14
1.2.17	The Ratkowsky-Lance index	15
1.2.18	The Ray-Turi index	16
1.2.19	The Scott-Symons index	16
1.2.20	The SD index	16
1.2.21	The S_Dbw index	17
1.2.22	The Silhouette index	18
1.2.23	The Tau index	19
1.2.24	The Trace_W index	19
1.2.25	The Trace_WiB index	20

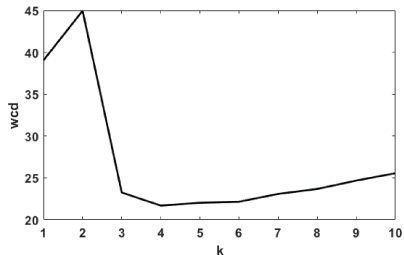
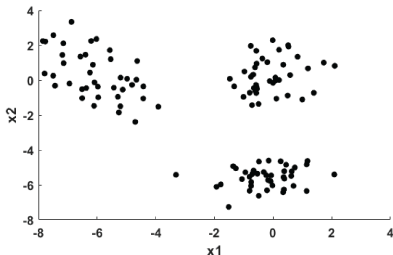
[1] Desgraupes, Bernard. "Clustering indices." University of Paris Ouest-Lab Modal'X 1 (2013): 34.

- In the studies of [2] and [3] the gap statistic is proposed. But...

[2] Witten, Daniela M., and Robert Tibshirani. "A framework for feature selection in clustering." Journal of the American Statistical Association 105.490 (2010): 713-726.

[3] Brodinova, Sarka, et al. "Robust and sparse k-means clustering for high-dimensional data." arXiv preprint arXiv:1709.10012 (2017).

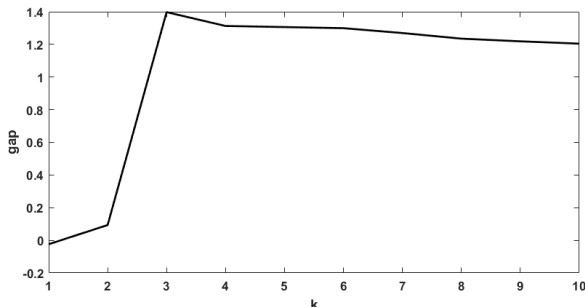
The Gap Statistic



[1] Tibshirani, Robert, Guenther Walther, and Trevor Hastie. "Estimating the number of clusters in a data set via the gap statistic." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63.2 (2001): 411-423.

The Gap Statistic

- Standardize the graph of $\log(WCSS_k)$ by comparing it with its expectation under an appropriate null reference distribution of the dataset.
- The optimal number of clusters is then the k for which $\log(WCSS_k)$ falls the farthest below the reference curve.
- $Gap_n(k) = E_n^*\{\log(WCSS_k)\} - \log(WCSS_k)$



The Gap Statistic: Algorithm

- Given dataset D , cluster it with different values for k and keep $\log(J_k)$, where J_k specifies the final value of the objective function of the clustering algorithm for a given k .
- Create B perturbations of D and for each repeat the above step, which results to: $\log(J_k^b) = [\log(J_k^1), \dots, \log(J_k^B)]$, for each k .
- Compute the estimated gap statistic,

$$Gap_k = \frac{1}{B} \sum_{b=1}^B \log(J_k^b) - \log(J_k)$$

- Given that $\overline{m_k^*} = \frac{1}{B} \sum_{b=1}^B \log(J_k^b)$, compute the simulation error,

$$SE_k = \sqrt{\frac{1}{B} \sum_{b=1}^B \left(\log(J_k^b) - \overline{m_k^*} \right)^2} + \sqrt{1 + \frac{1}{B}}$$

- Optimal k , \hat{k} , is the smallest k for which $Gap_k \geq Gap_{k+1} - SE_{k+1}$

The Gap Statistic: Algorithm

How many permutations, B ?

The Gap Statistic: Algorithm

How many permutations, B ?

- Witten, Daniela M., and Robert Tibshirani. "A framework for feature selection in clustering." Journal of the American Statistical Association 105.490 (2010): 713-726. **$B=25$** .
- [3] Brodinova, Sarka, et al. "Robust and sparse k-means clustering for high-dimensional data." arXiv preprint arXiv:1709.10012 (2017). **$B=10$** .

The Gap Statistic: Algorithm

How many permutations, B ?

- Witten, Daniela M., and Robert Tibshirani. "A framework for feature selection in clustering." Journal of the American Statistical Association 105.490 (2010): 713-726. **$B=25$** .
- [3] Brodinova, Sarka, et al. "Robust and sparse k-means clustering for high-dimensional data." arXiv preprint arXiv:1709.10012 (2017). **$B=10$** .
- MATLAB, **$B=100$** .
- R, **$B=100$** .

The Gap Statistic: Algorithm

How many permutations, B ?

- Witten, Daniela M., and Robert Tibshirani. “A framework for feature selection in clustering.” *Journal of the American Statistical Association* 105.490 (2010): 713-726. **$B=25$** .
- [3] Brodinova, Sarka, et al. “Robust and sparse k-means clustering for high-dimensional data.” *arXiv preprint arXiv:1709.10012* (2017). **$B=10$** .
- MATLAB, **$B=100$** .
- R, **$B=100$** .
- Let's say $B = 45$. We would like to test 10 different values for k and 5 for s .

The Gap Statistic: Algorithm

How many permutations, B ?

- Witten, Daniela M., and Robert Tibshirani. “A framework for feature selection in clustering.” *Journal of the American Statistical Association* 105.490 (2010): 713-726. **$B=25$** .
- [3] Brodinova, Sarka, et al. “Robust and sparse k-means clustering for high-dimensional data.” *arXiv preprint arXiv:1709.10012* (2017). **$B=10$** .
- MATLAB, **$B=100$** .
- R, **$B=100$** .
- Let's say $B = 45$. We would like to test 10 different values for k and 5 for s .
- We need to execute our algorithmic framework,
 $45 + 1(B) \times 10(k) \times 5(s) = 2300$ times!

Ongoing Research

Ongoing Research





- Find a computationally less expensive criterion than the gap statistic.
- Preliminary results indicated that:
 - ★ Silhouette index for s .
 - ★ Correlation for k .