

Text as Data

Justin Grimmer

Associate Professor
Department of Political Science
University of Chicago

January 22nd, 2018

Discovery

Search for new ways to organize text

- Complement, Not Replace, Organizations of Text
- There is No Ground Truth Conceptualization
- Once you have a conceptualization it is yours

Clustering: partition of documents

- Discover categories
- Assign documents to categories

Fully Automated Clustering

- 1) Notion of distance
- 2) Definition of “good” clustering
- 3) Optimization method

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

$$\boldsymbol{\theta}_k = (\theta_{1k}, \theta_{2k}, \dots, \theta_{Jk})$$

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

$$\boldsymbol{\theta}_k = (\theta_{1k}, \theta_{2k}, \dots, \theta_{Jk})$$

$\boldsymbol{\theta}_k = \text{exemplar}$ for cluster k

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

$$\boldsymbol{\theta}_k = (\theta_{1k}, \theta_{2k}, \dots, \theta_{Jk})$$

$\boldsymbol{\theta}_k = \text{exemplar}$ for cluster k

- 2) \mathbf{T} is an $N \times K$ matrix. Each row is an indicator vector.

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

$$\boldsymbol{\theta}_k = (\theta_{1k}, \theta_{2k}, \dots, \theta_{Jk})$$

$\boldsymbol{\theta}_k = \text{exemplar}$ for cluster k

- 2) \mathbf{T} is an $N \times K$ matrix. Each row is an indicator vector.

If observation i is from cluster k , then

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

$$\boldsymbol{\theta}_k = (\theta_{1k}, \theta_{2k}, \dots, \theta_{Jk})$$

$\boldsymbol{\theta}_k = \text{exemplar}$ for cluster k

- 2) \mathbf{T} is an $N \times K$ matrix. Each row is an indicator vector.

If observation i is from cluster k , then

$$\boldsymbol{\tau}_i = (0, 0, \dots, 0, \underbrace{1}_{k^{th}}, 0, \dots, 0)$$

K-Means \rightsquigarrow Objective Function

N documents $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ (normalized)

Goal \rightsquigarrow Partition documents into K clusters.

Two parameters to estimate

- 1) $K \times J$ matrix of cluster centers Θ .

Cluster k has center

$$\boldsymbol{\theta}_k = (\theta_{1k}, \theta_{2k}, \dots, \theta_{Jk})$$

$\boldsymbol{\theta}_k = \text{exemplar}$ for cluster k

- 2) \mathbf{T} is an $N \times K$ matrix. Each row is an indicator vector.

If observation i is from cluster k , then

$$\boldsymbol{\tau}_i = (0, 0, \dots, 0, \underbrace{1}_{k^{th}}, 0, \dots, 0)$$

Hard Assignment

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions
 - If $K = N$ then $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = 0$ (Minimum)

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions
 - If $K = N$ then $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = 0$ (Minimum)
 - Each observation in its own cluster

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions
 - If $K = N$ then $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = 0$ (Minimum)
 - Each observation in its own cluster
 - $\theta_i = \mathbf{x}_i$

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions
 - If $K = N$ then $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = 0$ (Minimum)
 - Each observation in its own cluster
 - $\theta_i = \mathbf{x}_i$
 - If $K = 1$, $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = N \times \sum_{j=1}^J \sigma_j^2$

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions
 - If $K = N$ then $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = 0$ (Minimum)
 - Each observation in its own cluster
 - $\theta_i = x_i$
 - If $K = 1$, $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = N \times \sum_{j=1}^J \sigma_j^2$
 - Each observation in same cluster

K-Means \rightsquigarrow Objective Function

Assume squared euclidean distance

$$f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{\tau_{ik}}_{\text{cluster indicator}} \underbrace{\left(\sum_{j=1}^J (x_{ij} - \theta_{kj})^2 \right)}_{\text{Squared Euclidean Distance}}$$

- Calculate squared euclidean distance from center
- **Only** for the assigned cluster
- Two trivial solutions
 - If $K = N$ then $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = 0$ (Minimum)
 - Each observation in its own cluster
 - $\theta_i = x_i$
 - If $K = 1$, $f(\mathbf{X}, \mathbf{T}, \mathbf{\Theta}) = N \times \sum_{j=1}^J \sigma_j^2$
 - Each observation in same cluster
 - $\theta_1 = \text{Average across documents}$

K-Means \rightsquigarrow Optimization

Coordinate descent

K-Means \rightsquigarrow Optimization

Coordinate descent \rightsquigarrow iterate between labels and centers.

K-Means \rightsquigarrow Optimization

Coordinate descent \rightsquigarrow iterate between labels and centers.

Iterative algorithm: each iteration t

K-Means \rightsquigarrow Optimization

Coordinate descent \rightsquigarrow iterate between labels and centers.

Iterative algorithm: each iteration t

- Conditional on Θ^{t-1} (from previous iteration), choose \mathbf{T}^t

K-Means \rightsquigarrow Optimization

Coordinate descent \rightsquigarrow iterate between labels and centers.

Iterative algorithm: each iteration t

- Conditional on Θ^{t-1} (from previous iteration), choose T^t
- Conditional on T^t , choose Θ^t

K-Means \rightsquigarrow Optimization

Coordinate descent \rightsquigarrow iterate between labels and centers.

Iterative algorithm: each iteration t

- Conditional on Θ^{t-1} (from previous iteration), choose \mathcal{T}^t
- Conditional on \mathcal{T}^t , choose Θ^t

Repeat until convergence \rightsquigarrow as measured as change in f dropping below threshold ϵ

K-Means \rightsquigarrow Optimization

Coordinate descent \rightsquigarrow iterate between labels and centers.

Iterative algorithm: each iteration t

- Conditional on Θ^{t-1} (from previous iteration), choose \mathbf{T}^t
- Conditional on \mathbf{T}^t , choose Θ^t

Repeat until convergence \rightsquigarrow as measured as change in f dropping below threshold ϵ

$$\text{Change} = f(\mathbf{X}, \mathbf{T}^t, \Theta^t) - f(\mathbf{X}, \mathbf{T}^{t-1}, \Theta^{t-1})$$

K-Means \rightsquigarrow Optimization

K-Means \rightsquigarrow Optimization

1) initialize K cluster centers $\theta_1^t, \theta_2^t, \dots, \theta_K^t$.

K-Means \rightsquigarrow Optimization

- 1) initialize K cluster centers $\theta_1^t, \theta_2^t, \dots, \theta_K^t$.
- 2) Choose T^t

K-Means \rightsquigarrow Optimization

- 1) initialize K cluster centers $\theta_1^t, \theta_2^t, \dots, \theta_K^t$.
- 2) Choose T^t

$$\tau_{im}^t = \begin{cases} 1 & \text{if } m = \arg \min_k \sum_{j=1}^J (x_{ij} - \theta_{kj}^t)^2 \\ 0 & \text{otherwise,} \end{cases}.$$

K-Means \rightsquigarrow Optimization

1) initialize K cluster centers $\theta_1^t, \theta_2^t, \dots, \theta_K^t$.

2) Choose \mathbf{T}^t

$$\tau_{im}^t = \begin{cases} 1 & \text{if } m = \arg \min_k \sum_{j=1}^J (x_{ij} - \theta_{kj}^t)^2 \\ 0 & \text{otherwise,} \end{cases}.$$

In words: Assign each document \mathbf{x}_i to the closest center θ_m^t

K-Means \rightsquigarrow Optimization

K-Means \rightsquigarrow Optimization

3) Choose $\Theta^t \rightsquigarrow$ Focus on the center for cluster k

K-Means \rightsquigarrow Optimization

3) Choose $\Theta^t \rightsquigarrow$ Focus on the center for cluster k

$$f(\mathbf{X}, \mathbf{T}^t, \mathbf{\Theta})_k = \sum_{i=1}^N \tau_{ik}^t \left(\sum_{j=1}^J (x_{ij} - \theta_{jk})^2 \right)$$

K-Means \rightsquigarrow Optimization

3) Choose $\Theta^t \rightsquigarrow$ Focus on the center for cluster k

$$f(\mathbf{X}, \mathbf{T}^t, \mathbf{\Theta})_k = \sum_{i=1}^N \tau_{ik}^t \left(\sum_{j=1}^J (x_{ij} - \theta_{jk})^2 \right)$$

$$\frac{\partial f(\mathbf{X}, \mathbf{T}^t, \mathbf{\Theta})_k}{\partial \theta_{kj}} = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk})$$

K-Means \rightsquigarrow Optimization

3) Choose $\Theta^t \rightsquigarrow$ Focus on the center for cluster k

$$f(\mathbf{X}, \mathbf{T}^t, \Theta)_k = \sum_{i=1}^N \tau_{ik}^t \left(\sum_{j=1}^J (x_{ij} - \theta_{jk})^2 \right)$$

$$\frac{\partial f(\mathbf{X}, \mathbf{T}^t, \Theta)_k}{\partial \theta_{kj}} = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk})$$

$$0 = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk}^*)$$

K-Means \rightsquigarrow Optimization

3) Choose $\Theta^t \rightsquigarrow$ Focus on the center for cluster k

$$f(\mathbf{X}, \mathbf{T}^t, \Theta)_k = \sum_{i=1}^N \tau_{ik}^t \left(\sum_{j=1}^J (x_{ij} - \theta_{jk})^2 \right)$$

$$\frac{\partial f(\mathbf{X}, \mathbf{T}^t, \Theta)_k}{\partial \theta_{kj}} = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk})$$

$$0 = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk}^*)$$

$$= \sum_{i=1}^N \tau_{ij}^t x_{ij} - \theta_{jk}^* \sum_{i=1}^N \tau_{ij}^t$$

K-Means \rightsquigarrow Optimization

3) Choose $\Theta^t \rightsquigarrow$ Focus on the center for cluster k

$$f(\mathbf{X}, \mathbf{T}^t, \Theta)_k = \sum_{i=1}^N \tau_{ik}^t \left(\sum_{j=1}^J (x_{ij} - \theta_{jk})^2 \right)$$

$$\frac{\partial f(\mathbf{X}, \mathbf{T}^t, \Theta)_k}{\partial \theta_{kj}} = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk})$$

$$0 = -2 \sum_{i=1}^N \tau_{ij}^t (x_{ij} - \theta_{jk}^*)$$

$$= \sum_{i=1}^N \tau_{ij}^t x_{ij} - \theta_{jk}^* \sum_{i=1}^N \tau_{ij}^t$$

$$\frac{\sum_{i=1}^N \tau_{ik}^t x_{ij}}{\sum_{i=1}^N \tau_{ik}^t} = \theta_{jk}^*$$

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}}$$

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .
Optimization algorithm:

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .

Optimization algorithm:

- Initialize centers

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .

Optimization algorithm:

- Initialize centers
- Do until converged:

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .

Optimization algorithm:

- Initialize centers
- Do until converged:
 - For each document, find closest center $\rightsquigarrow \tau_i^t$

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .
Optimization algorithm:

- Initialize centers
- Do until converged:
 - For each document, find closest center $\rightsquigarrow \tau_i^t$
 - For each center, take average of assigned documents $\rightsquigarrow \boldsymbol{\theta}_k^t$

K-Means \rightsquigarrow Optimization

$$\boldsymbol{\theta}^{t+1} = \frac{\sum_{i=1}^N \tau_{ik} \mathbf{x}_i}{\sum_{i=1}^N \tau_{ik}} \propto \sum_{i=1}^N \tau_{ik} \mathbf{x}_i$$

In words: $\boldsymbol{\theta}^{t+1}$ is the average of the documents assigned to k .

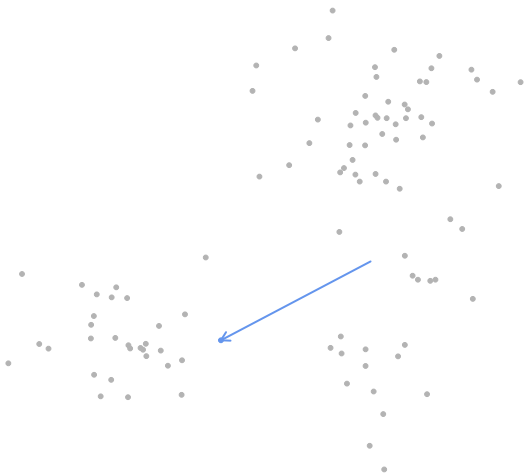
Optimization algorithm:

- Initialize centers
- Do until converged:
 - For each document, find closest center $\rightsquigarrow \tau_i^t$
 - For each center, take average of assigned documents $\rightsquigarrow \boldsymbol{\theta}_k^t$
 - Update change $f(\mathbf{X}, \mathbf{T}^t, \boldsymbol{\Theta}^t) - f(\mathbf{X}, \mathbf{T}^{t-1}, \boldsymbol{\Theta}^{t-1})$

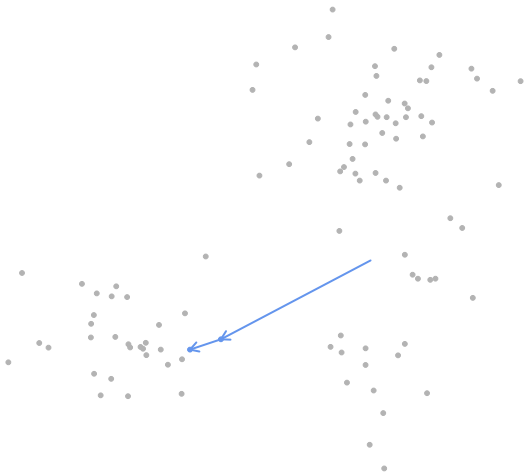
Visual Example



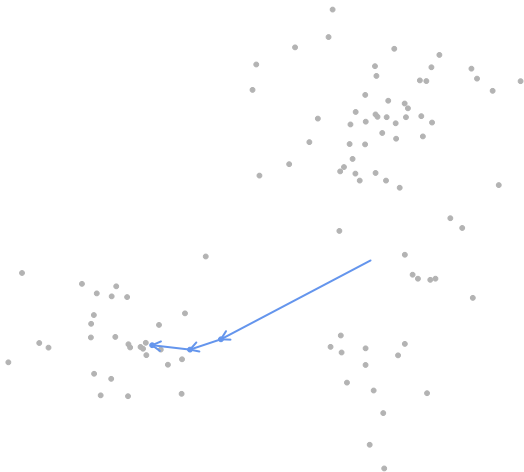
Visual Example



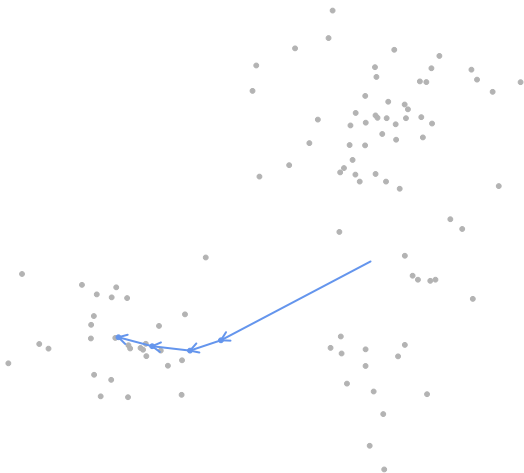
Visual Example



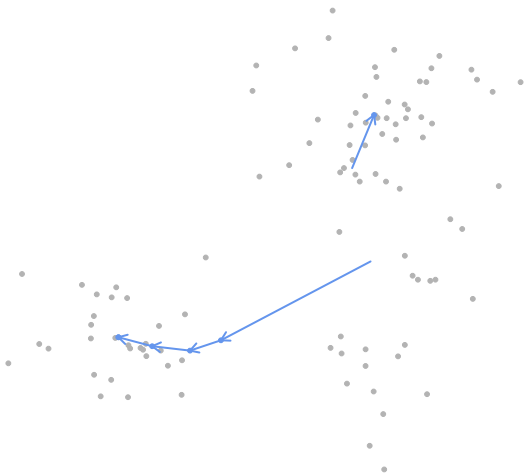
Visual Example



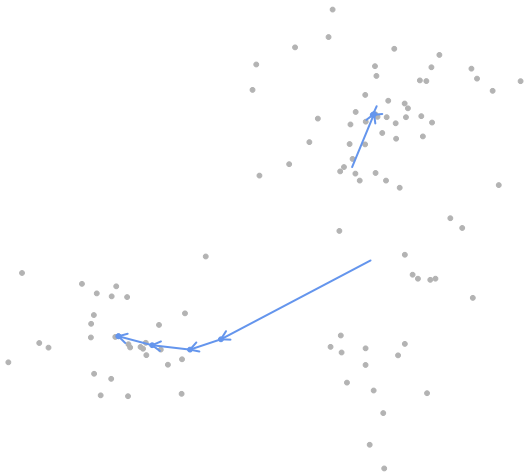
Visual Example



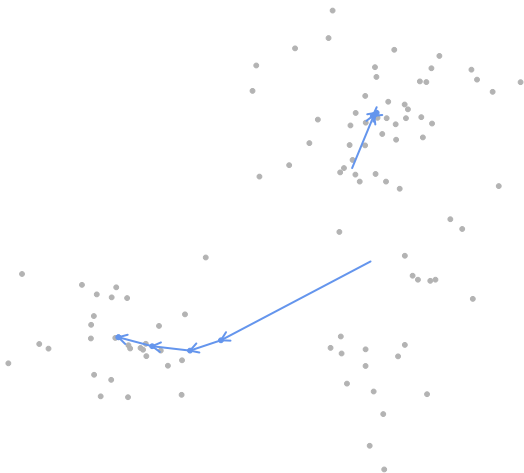
Visual Example



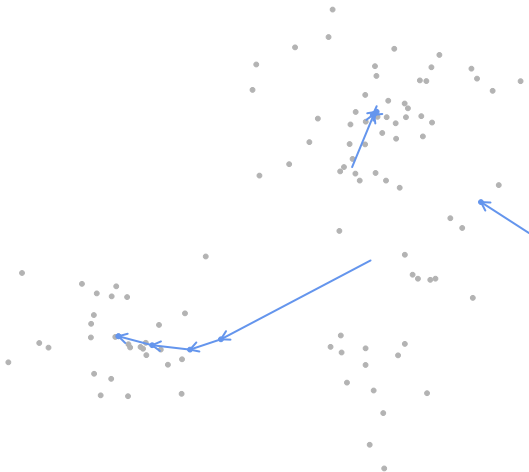
Visual Example



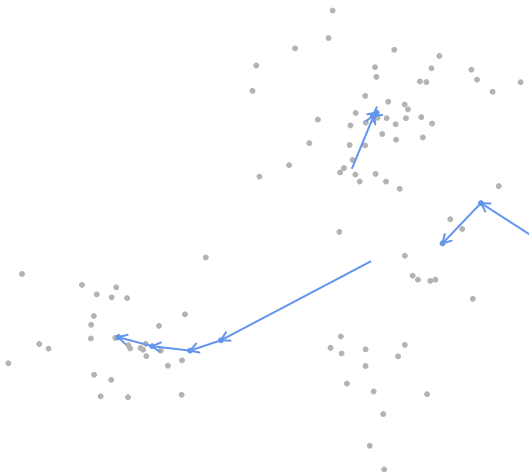
Visual Example



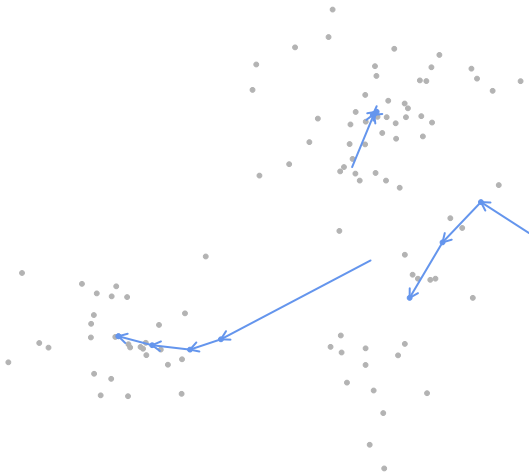
Visual Example



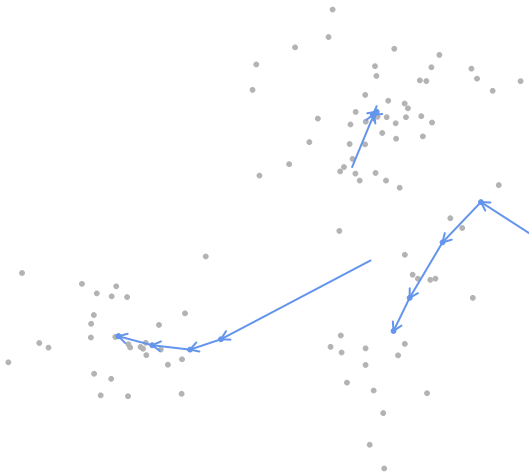
Visual Example



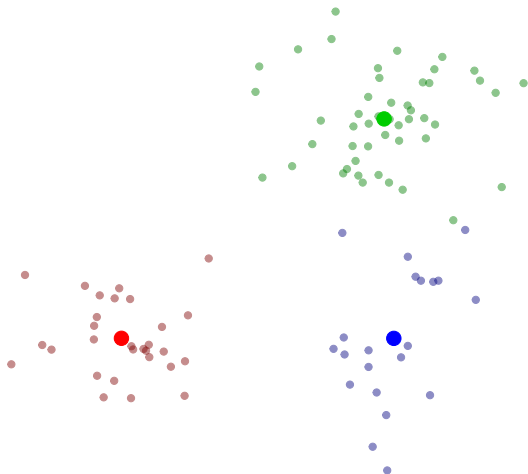
Visual Example



Visual Example



Visual Example



An Example: Jeff Flake

To the R Code!

Interpreting Cluster Components

Unsupervised methods

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words
 - Use these to help infer differences across clusters

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words
 - Use these to help infer differences across clusters
- Transparency

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words
 - Use these to help infer differences across clusters
- **Transparency**
 - Debate what clusters are

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words
 - Use these to help infer differences across clusters
- **Transparency**
 - Debate what clusters are
 - Debate what they mean

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words
 - Use these to help infer differences across clusters
- **Transparency**
 - Debate what clusters are
 - Debate what they mean
 - Provide documents + organizations

Interpreting Cluster Components

Unsupervised methods \rightsquigarrow low startup costs, high post-model costs

- Apply clustering methods, we have groups of documents
- How to interpret the groups?
- Two (broad) methods:
 - Manual identification (Quinn et al 2010)
 - Sample set of documents from same cluster
 - Read documents
 - Assign cluster label
 - Automatic identification
 - Know label classes
 - Use methods to identify separating words
 - Use these to help infer differences across clusters
- Transparency
 - Debate what clusters are
 - Debate what they mean
 - Provide documents + organizations

back to the R code!

How Do We Choose K ?

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

Think!

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

Think!

- No one statistic captures how you want to use your data

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

Think!

- No one statistic captures how you want to use your data
- But, can help guide your selection

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

Think!

- No one statistic captures how you want to use your data
- But, can help guide your selection
- Combination statistic + manual search

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

Think!

- No one statistic captures how you want to use your data
- But, can help guide your selection
- Combination statistic + manual search \leadsto discuss statistical methods/experimental methods on Thursday
- Humans should be the final judge

How Do We Choose K ?

- Previous Analysis Assumed We Know Number of Clusters
- How Do We Choose Cluster Number?
- Cannot Compare f across clusters
 - Sum squared errors decreases as K increases
 - Trivial answer: each document in own cluster (useless)
 - Modeling problem: Fit often increases with features
- How do we choose number of clusters?

Think!

- No one statistic captures how you want to use your data
- But, can help guide your selection
- Combination statistic + manual search \leadsto discuss statistical methods/experimental methods on Thursday
- **Humans should be the final judge**
 - Compare insights across clusterings

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of distribution data generating process:

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of distribution data generating process:

$$\tau_i | \boldsymbol{\pi} \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of distribution data generating process:

$$\tau_i | \boldsymbol{\pi} \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1 \sim \text{Distribution}(\text{parameters}_k)$$

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of distribution data generating process:

$$\tau_i | \boldsymbol{\pi} \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1 \sim \text{Distribution}(\text{parameters}_k)$$

In words:

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of distribution data generating process:

$$\tau_i | \boldsymbol{\pi} \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1 \sim \text{Distribution}(\text{parameters}_k)$$

In words:

- Draw a cluster label

Mixture of Unigram Models (Mixture of Multinomials)

Mixture models \rightsquigarrow wide range of applications

Single distribution data generating process:

$$\mathbf{x}_i \sim \text{Distribution}(\text{parameters})$$

Mixture of distribution data generating process:

$$\tau_i | \boldsymbol{\pi} \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1 \sim \text{Distribution}(\text{parameters}_k)$$

In words:

- Draw a cluster label
- Given distribution, draw realization

Mixture of Unigram Models (Mixture of Multinomials)

A mixture of unigram-language models

$$\boldsymbol{\pi} \sim \text{Dirichlet}(\mathbf{1})$$

$$\boldsymbol{\theta} \sim \text{Dirichlet}(\mathbf{1})$$

$$\tau_i | \boldsymbol{\pi} \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1, \boldsymbol{\theta}_k \sim \text{Multinomial}(N_i, \boldsymbol{\theta}_k)$$

Mixture of Unigram Models (Mixture of Multinomials)

This implies the following posterior distribution:

$$p(\mathbf{T}, \mathbf{\Theta}, \pi | \mathbf{X})$$

Mixture of Unigram Models (Mixture of Multinomials)

This implies the following posterior distribution:

$$p(\mathbf{T}, \mathbf{\Theta}, \pi | \mathbf{X}) \propto \overbrace{p(\pi)p(\boldsymbol{\theta})}^1 \underbrace{p(\mathbf{X}, \mathbf{T} | \pi, \boldsymbol{\theta})}_{\text{Complete data likelihood}}$$

Mixture of Unigram Models (Mixture of Multinomials)

This implies the following posterior distribution:

$$\begin{aligned} p(\mathbf{T}, \boldsymbol{\Theta}, \boldsymbol{\pi} | \mathbf{X}) &\propto \overbrace{p(\boldsymbol{\pi})p(\boldsymbol{\theta})}^1 \underbrace{p(\mathbf{X}, \mathbf{T} | \boldsymbol{\pi}, \boldsymbol{\theta})}_{\text{Complete data likelihood}} \\ &\propto \underbrace{\prod_{i=1}^N p(\boldsymbol{\tau}_i, \mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\pi})}_{\text{Complete data likelihood}} \end{aligned}$$

Mixture of Unigram Models (Mixture of Multinomials)

This implies the following posterior distribution:

$$\begin{aligned} p(\mathbf{T}, \boldsymbol{\Theta}, \boldsymbol{\pi} | \mathbf{X}) &\propto \overbrace{p(\boldsymbol{\pi})p(\boldsymbol{\theta})}^1 \underbrace{p(\mathbf{X}, \mathbf{T} | \boldsymbol{\pi}, \boldsymbol{\theta})}_{\text{Complete data likelihood}} \\ &\propto \underbrace{\prod_{i=1}^N p(\boldsymbol{\tau}_i, \mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\pi})}_{\text{Complete data likelihood}} \\ &\propto \prod_{i=1}^N p(\boldsymbol{\tau}_i | \boldsymbol{\pi}) p(\mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\tau}_i) \end{aligned}$$

Mixture of Unigram Models (Mixture of Multinomials)

This implies the following posterior distribution:

$$p(\mathbf{T}, \mathbf{\Theta}, \pi | \mathbf{X}) \propto \overbrace{p(\pi)p(\boldsymbol{\theta})}^1 \underbrace{p(\mathbf{X}, \mathbf{T} | \pi, \boldsymbol{\theta})}_{\text{Complete data likelihood}}$$

$$\propto \underbrace{\prod_{i=1}^N p(\boldsymbol{\tau}_i, \mathbf{x}_i | \boldsymbol{\theta}, \pi)}_{\text{Complete data likelihood}}$$

$$\propto \prod_{i=1}^N p(\boldsymbol{\tau}_i | \pi) p(\mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\tau}_i)$$

$$\propto \prod_{i=1}^N \prod_{k=1}^K \left[\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ik}} \right]^{\tau_{ik}}$$

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

1) Initialize parameters Θ^t, π^t

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

$$E[\log \text{Complete data} | \theta_k, \pi] = \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, \tau_{ik}^t | \theta_k, \pi_k) p(\tau_{ik}^t | \Theta, \pi_k)$$

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

$$E[\log \text{Complete data} | \theta_k, \pi] = \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, \tau_{ik}^t | \theta_k, \pi_k) p(\tau_{ik}^t | \Theta, \pi_k)$$

Obtain Θ^{t+1}, π^{t+1}

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

$$E[\log \text{Complete data} | \theta_k, \pi] = \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, \tau_{ik}^t | \theta_k, \pi_k) p(\tau_{ik}^t | \Theta, \pi_k)$$

Obtain Θ^{t+1}, π^{t+1}

- 4) Assess change

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

$$E[\log \text{Complete data} | \theta_k, \pi] = \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, \tau_{ik}^t | \theta_k, \pi_k) p(\tau_{ik}^t | \Theta, \pi_k)$$

Obtain Θ^{t+1}, π^{t+1}

- 4) Assess change

$$\begin{aligned} \text{Change} &= E[\log \text{Complete data} | \Theta^{t+1}, \pi^{t+1}] \\ &\quad - E[\log \text{Complete data} | \Theta^t, \pi^t] \end{aligned}$$

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

$$E[\log \text{Complete data} | \theta_k, \pi] = \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, \tau_{ik}^t | \theta_k, \pi_k) p(\tau_{ik}^t | \Theta, \pi_k)$$

Obtain Θ^{t+1}, π^{t+1}

- 4) Assess change

$$\begin{aligned} \text{Change} &= E[\log \text{Complete data} | \Theta^{t+1}, \pi^{t+1}] \\ &\quad - E[\log \text{Complete data} | \Theta^t, \pi^t] \end{aligned}$$

Our update steps will be strikingly similar to the K-Means algorithm

Mixture of Unigram Models (Mixture of Multinomials)

Obtain MAP estimates \rightsquigarrow EM Algorithm

- 1) Initialize parameters Θ^t, π^t
- 2) **Expectation step**: compute $p(\tau_i | \Theta^t, \pi^t, \mathbf{X}) \rightsquigarrow \mathbf{r}_i^t$
- 3) **Maximization step**: maximize with respect to Θ and π :

$$E[\log \text{Complete data} | \theta_k, \pi] = \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, \tau_{ik}^t | \theta_k, \pi_k) p(\tau_{ik}^t | \Theta, \pi_k)$$

Obtain Θ^{t+1}, π^{t+1}

- 4) Assess change

$$\begin{aligned} \text{Change} &= E[\log \text{Complete data} | \Theta^{t+1}, \pi^{t+1}] \\ &\quad - E[\log \text{Complete data} | \Theta^t, \pi^t] \end{aligned}$$

Our update steps will be strikingly similar to the K-Means algorithm

Mixture of Unigram Models (Mixture of Multinomials)

1) Initialize parameters Θ^t and π^t

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) E-Step

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) E-Step

$$p(\tau_{ik} | \Theta^t, \pi^t, \mathbf{X})$$

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) **E-Step**

$$p(\tau_{ik} | \Theta^t, \pi^t, \mathbf{X}) = \frac{\overbrace{p(\tau_{ik} | \pi^t) p(\mathbf{x}_i | \theta_k^t)}^{\text{general form}}}{\sum_{m=1}^K (p(\tau_{im} | \pi^t) p(\mathbf{x}_i | \theta_m^t))}$$

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) **E-Step**

$$\begin{aligned} p(\tau_{ik} | \Theta^t, \pi^t, \mathbf{X}) &= \frac{\overbrace{p(\tau_{ik} | \pi^t) p(\mathbf{x}_i | \theta_k^t)}^{\text{general form}}}{\sum_{m=1}^K (p(\tau_{im} | \pi^t) p(\mathbf{x}_i | \theta_m^t))} \\ &= \frac{\pi_k^t \prod_{j=1}^J (\theta_{jk}^t)^{x_{ij}}}{\sum_{m=1}^K (\pi_m^t \prod_{j=1}^J (\theta_{jm}^t)^{x_{ij}})} \end{aligned}$$

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) **E-Step**

$$\begin{aligned} p(\tau_{ik} | \Theta^t, \pi^t, \mathbf{X}) &= \frac{\overbrace{p(\tau_{ik} | \pi^t) p(\mathbf{x}_i | \theta_k^t)}^{\text{general form}}}{\sum_{m=1}^K (p(\tau_{im} | \pi^t) p(\mathbf{x}_i | \theta_m^t))} \\ &= \frac{\pi_k^t \prod_{j=1}^J (\theta_{jk}^t)^{x_{ij}}}{\sum_{m=1}^K (\pi_m^t \prod_{j=1}^J (\theta_{jm}^t)^{x_{ij}})} \end{aligned}$$

Define:

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) **E-Step**

$$\begin{aligned} p(\tau_{ik} | \Theta^t, \pi^t, \mathbf{X}) &= \frac{\overbrace{p(\tau_{ik} | \pi^t) p(\mathbf{x}_i | \theta_k^t)}^{\text{general form}}}{\sum_{m=1}^K (p(\tau_{im} | \pi^t) p(\mathbf{x}_i | \theta_m^t))} \\ &= \frac{\pi_k^t \prod_{j=1}^J (\theta_{jk}^t)^{x_{ij}}}{\sum_{m=1}^K (\pi_m^t \prod_{j=1}^J (\theta_{jm}^t)^{x_{ij}})} \end{aligned}$$

Define:

$$r_{ik}^t \equiv \frac{\pi_k^t \prod_{j=1}^J (\theta_{jk}^t)^{x_{ij}}}{\sum_{m=1}^K (\pi_m^t \prod_{j=1}^J (\theta_{jm}^t)^{x_{ij}})}$$

Mixture of Unigram Models (Mixture of Multinomials)

- 1) Initialize parameters Θ^t and π^t
- 2) **E-Step**

$$\begin{aligned} p(\tau_{ik} | \Theta^t, \pi^t, \mathbf{X}) &= \overbrace{\frac{p(\tau_{ik} | \pi^t) p(\mathbf{x}_i | \theta_k^t)}{\sum_{m=1}^K (p(\tau_{im} | \pi^t) p(\mathbf{x}_i | \theta_m^t))}}^{\text{general form}} \\ &= \frac{\pi_k^t \prod_{j=1}^J (\theta_{jk}^t)^{x_{ij}}}{\sum_{m=1}^K (\pi_m^t \prod_{j=1}^J (\theta_{jm}^t)^{x_{ij}})} \end{aligned}$$

Define: Avoid underflow

$$r_{ik}^t = \left[1 + \sum_{k' \neq k} \frac{\pi_{k'} \prod_{j=1}^J (\theta_{jk'}^t)^{x_{ij}}}{\pi_k \prod_{j=1}^J (\theta_{jk}^t)^{x_{ij}}} \right]^{-1}$$

Mixture of Unigram Models (Mixture of Multinomials)

3) **M-Step:**

Mixture of Unigram Models (Mixture of Multinomials)

3) M-Step:

$$E[\log \text{Complete data} | \boldsymbol{\theta}, \boldsymbol{\pi}] = \sum_{i=1}^N \sum_{k=1}^K E[\tau_{ik}] \log \left(\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ik}} \right)$$

Mixture of Unigram Models (Mixture of Multinomials)

3) M-Step:

$$\begin{aligned} E[\log \text{Complete data} | \boldsymbol{\theta}, \boldsymbol{\pi}] &= \sum_{i=1}^N \sum_{k=1}^K E[\tau_{ik}] \log \left(\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ik}} \right) \\ &= \sum_{i=1}^N \sum_{k=1}^K r_{ik}^t \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K \sum_{j=1}^J r_{ik}^t x_{ij} \log \theta_{jk} \end{aligned}$$

Mixture of Unigram Models (Mixture of Multinomials)

3) M-Step:

$$\begin{aligned} E[\log \text{Complete data} | \boldsymbol{\theta}, \boldsymbol{\pi}] &= \sum_{i=1}^N \sum_{k=1}^K E[\tau_{ik}] \log \left(\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ik}} \right) \\ &= \sum_{i=1}^N \sum_{k=1}^K r_{ik}^t \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K \sum_{j=1}^J r_{ik}^t x_{ij} \log \theta_{jk} \end{aligned}$$

Introducing constraints, differentiating, setting equal to zero and algebra yields:

Mixture of Unigram Models (Mixture of Multinomials)

3) M-Step:

$$\begin{aligned} E[\log \text{Complete data} | \boldsymbol{\theta}, \boldsymbol{\pi}] &= \sum_{i=1}^N \sum_{k=1}^K E[\tau_{ik}] \log \left(\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ik}} \right) \\ &= \sum_{i=1}^N \sum_{k=1}^K r_{ik}^t \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K \sum_{j=1}^J r_{ik}^t x_{ij} \log \theta_{jk} \end{aligned}$$

Introducing constraints, differentiating, setting equal to zero and algebra yields:

$$\pi_k^{t+1} = \frac{\sum_{i=1}^N r_{ik}^t}{N}$$

Mixture of Unigram Models (Mixture of Multinomials)

3) M-Step:

$$\begin{aligned} E[\log \text{Complete data} | \boldsymbol{\theta}, \boldsymbol{\pi}] &= \sum_{i=1}^N \sum_{k=1}^K E[\tau_{ik}] \log \left(\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ij}} \right) \\ &= \sum_{i=1}^N \sum_{k=1}^K r_{ik}^t \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K \sum_{j=1}^J r_{ik}^t x_{ij} \log \theta_{jk} \end{aligned}$$

Introducing constraints, differentiating, setting equal to zero and algebra yields:

$$\begin{aligned} \pi_k^{t+1} &= \frac{\sum_{i=1}^N r_{ik}^t}{N} \\ \theta_{jk}^{t+1} &= \frac{\sum_{i=1}^N r_{ik}^t x_{ij}}{\sum_{m=1}^J \sum_{i=1}^N r_{ik}^t x_{im}} \end{aligned}$$

Mixture of Unigram Models (Mixture of Multinomials)

3) M-Step:

$$\begin{aligned} E[\log \text{Complete data} | \boldsymbol{\theta}, \boldsymbol{\pi}] &= \sum_{i=1}^N \sum_{k=1}^K E[\tau_{ik}] \log \left(\pi_k \prod_{j=1}^J \theta_{jk}^{x_{ij}} \right) \\ &= \sum_{i=1}^N \sum_{k=1}^K r_{ik}^t \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K \sum_{j=1}^J r_{ik}^t x_{ij} \log \theta_{jk} \end{aligned}$$

Introducing constraints, differentiating, setting equal to zero and algebra yields:

$$\begin{aligned} \pi_k^{t+1} &= \frac{\sum_{i=1}^N r_{ik}^t}{N} \\ \theta_{jk}^{t+1} &= \frac{\sum_{i=1}^N r_{ik}^t x_{ij}}{\sum_{m=1}^J \sum_{i=1}^N r_{ik}^t x_{im}} \propto \sum_{i=1}^N r_{ik}^t \mathbf{x}_i \end{aligned}$$

Example: Jeff Flake Again!

To the R Code!

Fully Automated Clustering

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday
 - Validation: model based fit statistics

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday
 - Validation: model based fit statistics
- How do we know we have the “right” model?

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday
 - Validation: model based fit statistics
- How do we know we have the “right” model?

YOU DON'T!

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday
 - Validation: model based fit statistics
- How do we know we have the “right” model?

YOU DON'T! \rightsquigarrow And never will

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday
 - Validation: model based fit statistics
- How do we know we have the “right” model?

YOU DON'T! \rightsquigarrow And never will \rightsquigarrow but
still useful(!!!!)

Fully Automated Clustering

- Notion of similarity and “good” partition \rightsquigarrow clustering
- Many clustering methods:
 - Spectral clustering
 - Affinity Propagation
 - Non-parametric statistical models
 - Hierarchical clustering
 - Biclustering
 - ...
- How do we know we have something useful?
 - Validation: read the documents
 - Validation: experiments to assess cluster quality \rightsquigarrow Thursday
 - Validation: model based fit statistics
- How do we know we have the “right” model?

YOU DON'T! \rightsquigarrow And never will \rightsquigarrow but
still useful(!!!!)

Appendix: Why EM Works

Goal:

$$\operatorname{argmax}_{\theta} p(\mathbf{X}|\theta) = \sum_{\mathbf{T}} p(\mathbf{X}, \mathbf{T}|\theta)$$

Define:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{T}} q(\mathbf{T}) \log \left[\frac{p(\mathbf{X}, \mathbf{T}|\theta)}{q(\mathbf{T})} \right] \\ K(q||p) &= - \sum_{\mathbf{T}} q(\mathbf{T}) \log \left[\frac{p(\mathbf{T}|\mathbf{X}, \theta)}{q(\mathbf{T})} \right]\end{aligned}$$

Then:

$$\log p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + K(q||p)$$

Appendix: Why EM Works

$$\begin{aligned}\log p(\mathbf{X}|\boldsymbol{\theta}) &= \mathcal{L}(q, \boldsymbol{\theta}) + K(q||p) \\&= \sum_{\mathbf{T}} q(\mathbf{T}) \log \left[\frac{p(\mathbf{X}, \mathbf{T}|\boldsymbol{\theta})}{q(\mathbf{T})} \right] - \sum_{\mathbf{T}} q(\mathbf{T}) \log \left[\frac{p(\mathbf{T}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{T})} \right] \\&= \sum_{\mathbf{T}} q(\mathbf{T}) \log(p(\mathbf{X}|\boldsymbol{\theta})) + \sum_{\mathbf{T}} q(\mathbf{T}) \log(p(\mathbf{T}|\mathbf{X}, \boldsymbol{\theta})) \\&\quad - \sum_{\mathbf{T}} q(\mathbf{T}) \log q(\mathbf{T}) - \sum_{\mathbf{T}} q(\mathbf{T}) \log p(\mathbf{T}|\mathbf{X}, \boldsymbol{\theta}) + \sum_{\mathbf{T}} q(\mathbf{T}) \log q(\mathbf{T})\end{aligned}$$

Collect terms that cancel and recognize $\sum_{\mathbf{T}} q(\mathbf{T}) = 1$ and we see equivalence

Appendix: Why EM Works

$K(q||p) \geq 0$ with $K(q||p) = 0$ only if $q = p$. So, $\mathcal{L}(q, \theta)$ is a lower-bound on the log-likelihood.

E-step

$$\log p(\mathbf{X}|\theta) - K(q||p) = \mathcal{L}(q, \theta)$$

$\mathcal{L}(q, \theta) \rightsquigarrow$ biggest when $K(q||p) = 0$, so set

$$q(\mathbf{T}) = p(\mathbf{T}|\mathbf{X}, \theta)$$

M-step:

Given the new value of q , maximize parameters (expectation of the log complete data likelihood)

Change in log-likelihood will be greater \rightsquigarrow because new maximum induces non-zero KL-divergence. Changes in log-likelihood are greater than changes in lower bound.