

K-MEANS CLUSTERING

APMA4903 TALK

Anji Zhao

9 November 2015

INTRODUCTION

K-MEANS AND LLOYD'S ALGORITHM

GAUSSIAN MIXTURE MODELS

APPLICATION: CLUSTERING WEATHER DATA

FUTURE

REFERENCES

- ▶ H. Daumé, *A Course in Machine Learning*, 2015.
- ▶ T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning*, Springer, 2013.
- ▶ D. Hsu, Lecture Slides, *COMS 4771 - Elementary Machine Learning*, Columbia University, 2015.
- ▶ D. Arthur, S. Vassilvitskii, “k-means++: The Advantages of Careful Seeding”, *SODA '07 Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 1027-1035, 2007.
- ▶ S. D. Roy, G. Lotan, “Detecting geo-spatial weather clusters using dynamic heuristic subspaces”, *Information Reuse and Integration (IRI), 2014 IEEE 15th International Conference on*, 811-818, 2014.
- ▶ Wikipedia, “k-means clustering”,
https://en.wikipedia.org/wiki/K-means_clustering.

UNSUPERVISED LEARNING

goal: find hidden structure behind unlabeled data

examples:

- ▶ partitioning data into clusters
- ▶ dimensionality reduction

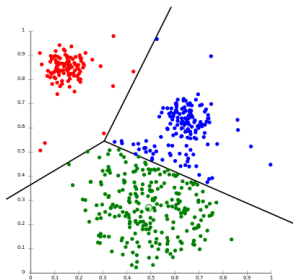
no direct measure of success

CLUSTERING

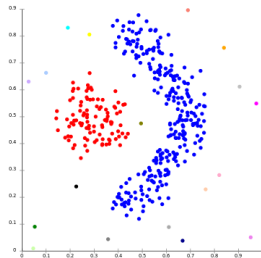
partition a dataset into groups of “similar” data points

some types of clustering:

- ▶ centroid-based (k-means)
- ▶ connectivity-based (hierarchical clustering)
- ▶ distribution-based



k-means



hierarchical

K-MEANS CLUSTERING

The Problem:

- ▶ input: n points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, $k \in \mathbb{N}$
- ▶ output: k centers $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ and (optionally) n cluster assignments $z_1, z_2, \dots, z_n \in \{1, 2, \dots, k\}$
- ▶ objective: choose $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ to minimize the within-cluster sum of squares:

$$SSE(\mathbf{x}, \mathbf{c}) = \sum_{i=1}^n \min_{j \in \{1, 2, \dots, k\}} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

- ▶ the k-means problem is NP-hard

LLOYD'S ALGORITHM

- ▶ also known as the k-means algorithm or Lloyd-Forgy algorithm
- ▶ iterative greedy algorithm which converges to a local optimum
- ▶ pseudocode:

```
for  $k$  in range( $K$ ):  
     $\mathbf{c}_k$  = random vector in  $\mathbb{R}^d$   
    repeat until convergence:  
        for  $i$  in range( $n$ ):  
             $z_i = \operatorname{argmin}_k \|\mathbf{c}_k - \mathbf{x}_i\|$   
        for  $k$  in range( $K$ ):  
             $S_k = \{\mathbf{x}_i : z_i = k\}$   
             $\mathbf{c}_k = \operatorname{mean}(S_k)$   
    return  $\mathbf{c}, \mathbf{z}$ 
```

CONVERGENCE

For any set of n points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and number of clusters $k \in \mathbb{N}$, Lloyd's algorithm converges in a finite number of iterations.

$$\begin{aligned} SSE(\mathbf{x}, \mathbf{c}) &= \sum_{i=1}^n \min_{j \in \{1, 2, \dots, k\}} \|\mathbf{x}_i - \mathbf{c}_j\|^2 \\ &= \sum_{j=1}^k \sum_{i: z_i = j} \|\mathbf{x}_i - \mathbf{c}_j\|^2 \\ &\text{where } z_i = \operatorname{argmin}_k \|\mathbf{c}_k - \mathbf{x}_i\| \end{aligned}$$

“Convergence” = SSE stops changing

quick proof on board

DRAWBACKS OF LLOYD'S ALGORITHM

- ▶ converges to a local, not global minimum
- ▶ arbitrarily bad clusters depending on initialization
- ▶ yields poor results when k is chosen incorrectly

K-MEANS++ ALGORITHM

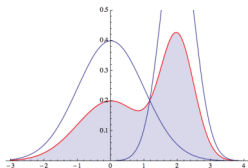
- ▶ proposed in 2007 by David Arthur and Sergei Vassilvitskii
- ▶ a method of choosing initial points to obtain more accurate clusterings than standard k-means
- ▶ algorithm:
 - ▶ choose one center c_1 uniformly at random from $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$
 - ▶ until we have k centers, choose center c_i from the dataset, picking x with probability $\frac{D(x)^2}{\sum_x D(x)^2}$, where $D(x)$ is the distance from x to the closest center we have already chosen.
 - ▶ using these initial centers, proceed with the standard k-means algorithm
- ▶ guarantees $E[\phi] \leq 8(\ln k + 2)\phi_{OPT}$

GAUSSIAN MIXTURE MODELS

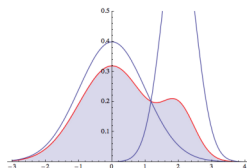
- ▶ Model dataset by a mixture of k probability distributions, where the j th component is a Gaussian distribution with μ_j and Σ_j , $j = 1, 2, \dots, k$.
- ▶ Formally: $(\mathbf{X}, Y) \sim P_\theta$, a distribution over $\mathbb{R}^d \times [k]$, where:
 - ▶ $Y \sim \pi$
 - ▶ $\mathbf{X}|Y = j \sim N(\mu_j, \Sigma_j)$
- ▶ P_θ has parameters $\theta = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_k, \mu_k, \Sigma_k)$
- ▶ Modeling assumption: our data $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times [k]$ is an iid sample from P , but we only know $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$.
- ▶ Gaussian mixture model:

$$\mathbf{X} \sim \sum_{j=1}^k \pi_j N(\mu_j, \Sigma_j)$$

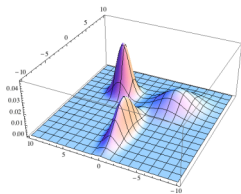
GAUSSIAN MIXTURES IN \mathbb{R}^1 AND \mathbb{R}^2



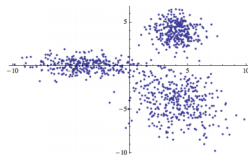
$$\frac{1}{2}N(0, 1) + \frac{1}{2}N(2, \frac{1}{4})$$



$$\frac{4}{5}N(0, 1) + \frac{1}{5}N(2, \frac{1}{4})$$



Mixture density



Observed sample points

GAUSSIAN PROBABILITY DENSITIES

- ▶ in one dimension ($x \in \mathbb{R}^1$):

$$p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- ▶ in d dimensions ($\mathbf{x} \in \mathbb{R}^d$):

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

- ▶ for a mixture of k Gaussians in d dimensions ($\mathbf{x} \in \mathbb{R}^d$):

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{j=1}^k \pi_j \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma}_j)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j)\right)$$

SOFT CLUSTERING

- ▶ instead of assigning each point to one component (as in k-means), we find the probability that the point belongs to each component
- ▶ suppose we are given the parameters of a Gaussian mixture model, and $(\mathbf{X}, Y) \sim P_\theta$
- ▶ let $\Phi \in \{0, 1\}^k$ be the vector of assignment variables, where $\Phi_j = \mathbb{1}\{Y = j\}$
- ▶ soft assignment of a data point \mathbf{x} to component j :

$$\begin{aligned} E_\theta[\Phi_j | \mathbf{X} = \mathbf{x}] &= Pr_\theta[Y = j | \mathbf{X} = \mathbf{x}] \\ &= \frac{Pr_\theta[Y = j] \cdot Pr_\theta[\mathbf{X} = \mathbf{x} | Y = j]}{Pr_\theta[\mathbf{X} = \mathbf{x}]} \\ &= \frac{\pi_j \cdot \sqrt{\det(\Sigma_j^{-1})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j)\right)}{\sum_{i=1}^k \left(\pi_i \cdot \sqrt{\det(\Sigma_i^{-1})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)\right)\right)} \end{aligned}$$

MAXIMUM LIKELIHOOD ESTIMATOR

Maximum Likelihood Estimation is a method used to estimate the parameters of a model, given a set of data points.

The MLE for a model P is

$$\theta_{ML} = \operatorname{argmax}_{\theta} \prod_{i=1}^n p(\mathbf{x}_i; \theta)$$

We can try to use the MLE to estimate the parameters of a Gaussian mixture, but we get...

$$\begin{aligned}\theta_{ML} &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln p(\mathbf{x}_i, \theta) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln (...)\end{aligned}$$

EXPECTATION-MAXIMIZATION IDEA

Iterative local optimization method for estimating parameters.

Given labeled data $(\mathbf{x}_1, \phi_1), (\mathbf{x}_2, \phi_2), \dots, (\mathbf{x}_n, \phi_n) \in \mathbb{R}^d \times \{0, 1\}^k$, the “complete log-likelihood” of $\theta = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_k, \mu_k, \Sigma_k)$ is

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^k \phi_{i,j} \ln \left(\pi_j \cdot \sqrt{\det(\Sigma_j^{-1})} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x} - \mu_j) \right) \right) \\ &= \sum_{i=1}^n \sum_{j=1}^k \phi_{i,j} \left(\ln \pi_j + \frac{1}{2} \ln \det(\Sigma_j^{-1}) - \frac{1}{2} (\mathbf{x} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x} - \mu_j) \right) \end{aligned}$$

(can also use soft assignments $w_{i,j} = E_{\theta}[\phi_{i,j} | \mathbf{X} = \mathbf{x}_i]$ instead of $\phi_{i,j}$)

E-M ALGORITHM

Algorithm:

- ▶ initialize $\theta = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_k, \mu_k, \Sigma_k)$
- ▶ E step: calculate the expectation of the unknown labels / soft assignments given the parameters θ :
 - ▶ $w_{i,j} = E_{\theta}[\phi_{i,j} | \mathbf{X} = \mathbf{x}_i], \forall i \in \{1, 2, \dots, n\}, j \in \{1, 2, \dots, k\}$
- ▶ M step: maximize the expected complete log-likelihood w.r.t. each parameter
 - ▶ $\pi_j = \frac{1}{n} \sum_{i=1}^n w_{i,j}$
 - ▶ $\mu_j = \frac{1}{n\pi_j} \sum_{i=1}^n w_{i,j} \mathbf{x}_i$
 - ▶ $\Sigma_j = \frac{1}{n\pi_j} \sum_{i=1}^n w_{i,j} (\mathbf{x}_i - \mu_j)(\mathbf{x}_i - \mu_j)^T$

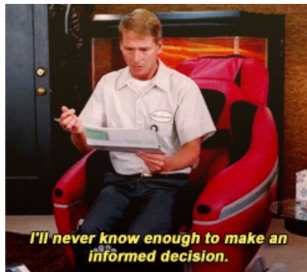
GAUSSIAN MIXTURE MODELS VS. K-MEANS

- ▶ k-means is a special case of GMM where we restrict $\Sigma_i = I \ \forall i \in [k]$ and $\pi_i = \pi_j \ \forall i, j \in [k]$
- ▶ in k-means we use hard assignment in the E-step
- ▶ k-means converges faster, but GMM is more flexible

APPLICATION: CLUSTERING WEATHER DATA

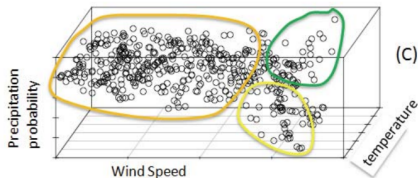
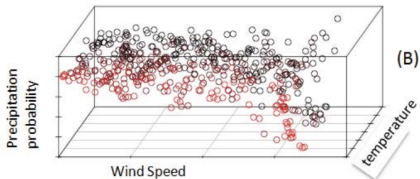
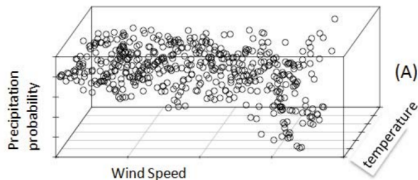


› **REMEMBER TO
VOTE!**



Clear skies with temps dipping into 15C. Do you like this weather? If yes, retweet! If not, favorite -- sorry, I mean "like."

APPLICATION: CLUSTERING WEATHER DATA



Weather Data distribution
at 12PM, Feb 16, 2011, in
Stamford, CT.

- A. The original data
- B. Two clusters using
typical k-means
- C. Three clusters after
heuristic splits and
rounding

Image by Roy, Lotan

FUTURE & QUESTIONS

Current Areas of Research:

- ▶ improving performance of existing clustering algorithms
- ▶ handling high-dimensional data
- ▶ specific applications to different fields