

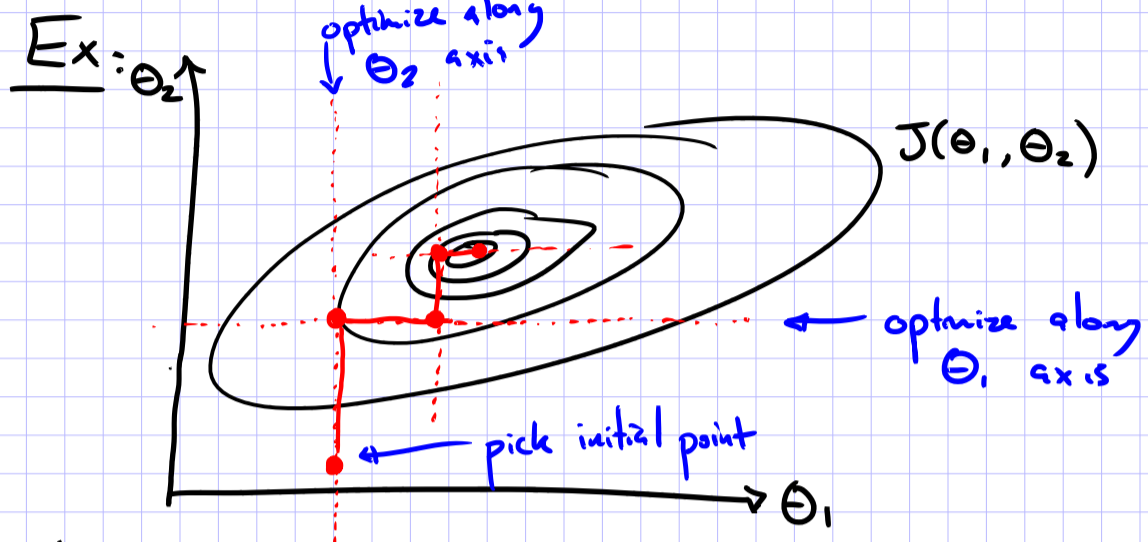
Lecture 15: 3/8/17

Optimization Background: Coordinate Descent

Goal: Minimize a function $J(\vec{\theta})$

e.g. $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$

Idea: Pick one dimension, and minimize along that dimension.



Algorithm:

- Choose initial point $\vec{\theta}$
- Repeat until stopping criterion is reached.

$$\theta_1 = \underset{\theta_1}{\operatorname{argmin}} J(\theta_1, \theta_2, \dots, \theta_M)$$

$$\theta_2 = \underset{\theta_2}{\operatorname{argmin}} J(\theta_1, \theta_2, \dots, \theta_M)$$

$$\vdots$$

$$\theta_M = \underset{\theta_M}{\operatorname{argmin}} J(\theta_1, \theta_2, \dots, \theta_M)$$

All the blue variables are "fixed" while we minimize over green vars.

each step is an exact line search along some axis

③ Return $\vec{\theta}$

Block Coordinate Descent

Here: An example w/ two blocks $\vec{\alpha}, \vec{\beta}$, where $\vec{\theta} = \begin{bmatrix} \vec{\alpha} \\ \vec{\beta} \end{bmatrix}$

Goal: $\vec{\alpha}, \vec{\beta} = \underset{\vec{\alpha}, \vec{\beta}}{\operatorname{argmin}} J(\vec{\alpha}, \vec{\beta}), \vec{\alpha} \in \mathbb{R}^A, \vec{\beta} \in \mathbb{R}^B$

Idea: Minimize over an entire group of variables at a time.

Algorithm:

- Choose initial point $\vec{\alpha}, \vec{\beta}$
- Repeat until stopping criterion

$$\vec{\alpha} = \underset{\vec{\alpha}}{\operatorname{argmin}} J(\vec{\alpha}, \vec{\beta})$$

$$\vec{\beta} = \underset{\vec{\beta}}{\operatorname{argmin}} J(\vec{\alpha}, \vec{\beta})$$

blue is "fixed" aka. held constant, while minimizing over green.

Clustering

our first example of unsupervised learning

GOAL: partition unlabeled instances into groups of "similar" points

Input: Unlabeled data: $D = \{\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(N)}\}$, $\vec{x}^{(i)} \in \mathbb{R}^M$

*** We do not know the labels of the training examples.**

Output:

View #1:

Labeled Instances

$$\{(\vec{x}^{(1)}, z^{(1)}), (\vec{x}^{(2)}, z^{(2)}), \dots, (\vec{x}^{(N)}, z^{(N)})\}$$

where $z^{(i)} \in \{1, \dots, K\}$

these cluster assignments are predictions

of clusters

View #2:

Clusterings

$$C_1, C_2, \dots, C_k$$

$$C_j = \{\vec{x}^{(i)} : z^{(i)} = j\}$$

the points in the j -th partition.

Important Questions:

- How many clusters are there?

- How do we define "similarity" between points?

Objective-Based Clustering

Ex: K-Means Objective.

Input: $D = \{\vec{x}^{(i)}\}_{i=1}^N$

lower case c .

Cluster Centers: $\vec{c} = \{\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k\}$

Decision Rule: Assign $\vec{x}^{(i)}$ to its nearest cluster center \vec{c}_j

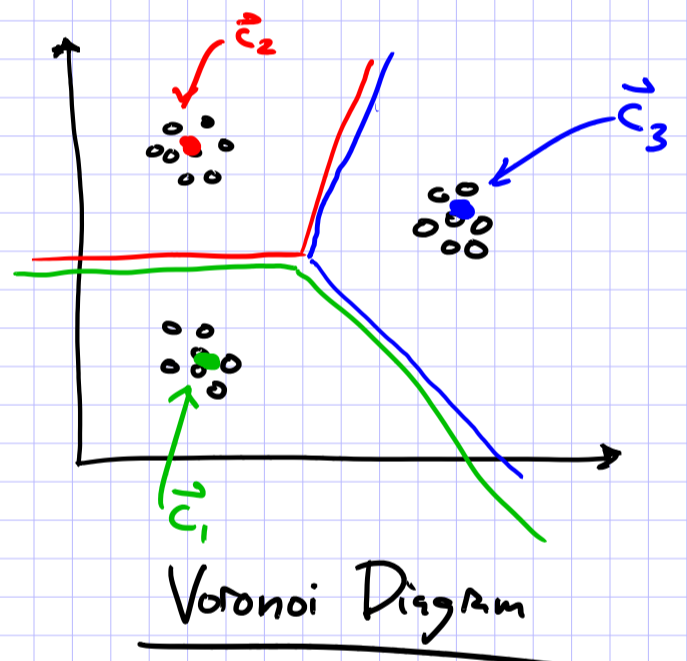
Objective:

$$\vec{c} = \underset{\vec{c}}{\operatorname{argmin}} \sum_{i=1}^N \min_{j \in \{1, \dots, K\}} \|\vec{x}^{(i)} - \vec{c}_j\|^2$$

\vec{c}_j is the cluster center

Each point should as close as possible to its cluster center.

All points should be close to their cluster centers.



Equivalent Objective:

Let \vec{z} be all the cluster assignments $\vec{z} = [z^{(1)}, \dots, z^{(N)}]$
 $z^{(i)} \in \{1, \dots, K\}$

$$\vec{c} = \underset{\vec{c}}{\operatorname{argmin}} \sum_{i=1}^N \min_{z^{(i)}} \|x^{(i)} - c_{z^{(i)}}\|^2$$

$$\equiv \vec{c}, \vec{z} = \underset{\vec{c}, \vec{z}}{\operatorname{argmin}} \sum_{i=1}^N \|x^{(i)} - c_{z^{(i)}}\|^2$$

push min out of sum.

call this $J_{\text{KMeans}}(\vec{c}, \vec{z})$

$$= \underset{\vec{c}, \vec{z}}{\operatorname{argmin}} J_{\text{KMeans}}(\vec{c}, \vec{z})$$

Question: How should we optimize $J_{\text{KMeans}}(\vec{c}, \vec{z})$?

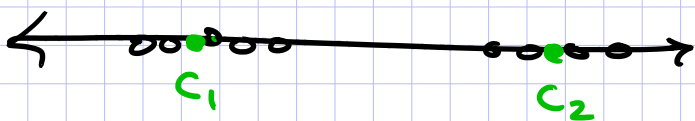
Computational Complexity of K-Means Obj. Minimization Problem

- ① Objective is nonconvex
- ② NP-Hard, even for $k=2$ ← # clusters
 even for $M=2$ ← # features

Easy Case #1: $K=1$

$$\begin{aligned} \vec{c}_1 &= \underset{\vec{c}_1}{\operatorname{argmin}} \sum_{i=1}^N \|\vec{x}^{(i)} - \vec{c}_1\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N \vec{x}^{(i)} \quad \leftarrow \text{mean} \end{aligned}$$

Easy Case #2: $M=1$



Dynamic programming in time $O(N^2K)$

K-Means in Practice

- Solve minimization problem heuristically w/ Block Coord. Descent.

K-Means Algorithm:

- ① Given $\vec{x}^{(1)}, \dots, \vec{x}^{(N)}$
- ② Initialize cluster centers $\vec{c} = \{\vec{c}_1, \dots, \vec{c}_k\}$
 Initialize cluster assignments $\vec{z} = \{z^{(1)}, \dots, z^{(N)}\}$
- ③ Repeat until objective stops changing.

$$\begin{aligned} \vec{a} &= \underset{a}{\operatorname{argmin}} f(a_1) + g(a_2) \\ &\equiv \\ a_1 &= \underset{a_1}{\operatorname{argmin}} f(a_1) \\ a_2 &= \underset{a_2}{\operatorname{argmin}} g(a_2) \end{aligned}$$

$$\begin{aligned} \text{a) } \vec{c} &= \underset{\vec{c}}{\operatorname{argmin}} \sum_{i=1}^N \|\vec{x}^{(i)} - \vec{c}_{z^{(i)}}\|^2 \quad (\text{Min over centers, w/ assignments fixed}) \\ \text{b) } \vec{z} &= \underset{\vec{z}}{\operatorname{argmin}} \sum_{i=1}^N \|\vec{x}^{(i)} - \vec{c}_{z^{(i)}}\|^2 \quad (\text{Min over assignments, w/ centers fixed}) \end{aligned}$$

3a) decomposes:

$$\sum_{i=1}^N \|\vec{x}^{(i)} - \vec{c}_{z^{(i)}}\|^2 = \sum_{j=1}^K \sum_{i: z^{(i)}=j} \|\vec{x}^{(i)} - \vec{c}_j\|^2$$

$$\begin{aligned} \vec{c}_1 &= \underset{\vec{c}_1}{\operatorname{argmin}} \sum_{i: z^{(i)}=1} \|\vec{x}^{(i)} - \vec{c}_1\|^2 \\ \vec{c}_2 &= \underset{\vec{c}_2}{\operatorname{argmin}} \dots \\ \dots \\ \vec{c}_k &= \underset{\vec{c}_k}{\operatorname{argmin}} \sum_{i: z^{(i)}=k} \|\vec{x}^{(i)} - \vec{c}_k\|^2 \end{aligned}$$

Each is just Easy Case #1:

$$\begin{aligned} \vec{c}_j &= \underset{\vec{c}_j}{\operatorname{argmin}} \sum_{i: z^{(i)}=j} \|\vec{x}^{(i)} - \vec{c}_j\|^2 \\ &= \text{mean of points in cluster } j \\ &= \frac{1}{N_j} \sum_{i: z^{(i)}=j} \vec{x}^{(i)} \end{aligned}$$

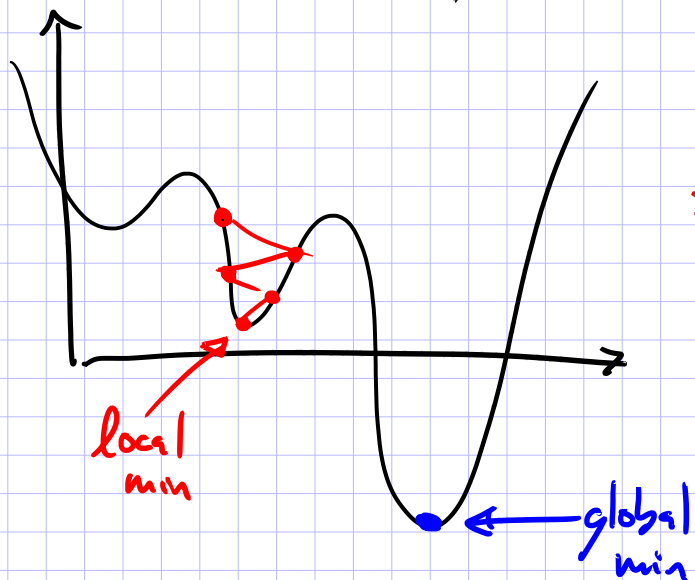
3b) decomposes

$$\begin{aligned} z^{(i)} &= \underset{j}{\operatorname{argmin}} \|\vec{x}^{(i)} - \vec{c}_j\|^2 \\ \dots \\ z^{(N)} &= \underset{j}{\operatorname{argmin}} \|\vec{x}^{(N)} - \vec{c}_j\|^2 \end{aligned}$$

just find the closest cluster center for each point $\vec{x}^{(i)}$!

Initialization

This is a local opt. alg. for a nonconvex objective.



K-Means just finds a local min.
 \Rightarrow How we initialize \tilde{c}, \tilde{z} is critical.

Three Options:

① Random: select points uniformly at random (w/o replacement) from dataset D

② Furthest Traversal:

select points in D s.t. c_j is as far as possible from c_1, \dots, c_{j-1}

③ K-Means++:

interpolate between ① and ②

★ good theoretical guarantees.