# Lec 10

Tuesday, October 1, 2019    10:51

# Recap: PCA

## Auto-encoder view of PCA

encoder    $e : \mathbb{R}^p \longrightarrow \mathbb{R}^q$

$$z = e(x) = A^T x$$

decoder    $d : \mathbb{R}^q \rightarrow \mathbb{R}^p$

$$\hat{x} = d(z) = A z$$

Where

$A \in \mathbb{R}^{p \times q}$

$A^T A = I$

Soln : PCA : get $A = V_q$

↑

for minimizing

$\sum_i \| \hat{x}_i - x_i \|_2^2$

↓

$\sum_i \| d(e(x_i)) - x_i \|_2^2$

the first $q$ cols of $V$ in the SVD $X = U \Sigma V^T$

## Non-Centered ~~data~~ PCA for

decoder:   $\hat{x} = d(z) = A z + a$

encoder:   $z = e(x) = A^T(x - a)$

Best $A, a$ are:

$$a = \overline{X} = \frac{1}{n} \sum_i x_i$$

$A = V_q = $ first $q$ cols from the SVD

of   $\underline{X - \overline{x}} = U \Sigma V^T$

where $(X - \overline{x})_i \overset{i\text{th row}}{=} x_i - \overline{x}$

$= $ first $q$ eigenvectors from the eigendecomp of

$$\frac{1}{n} (X - \overline{x})^T (X - \overline{x})$$

# Clustering

Assign data points to finitely many, $k \in \mathbb{N}$, clusters

One perspective: find clusters in the data

Auto-encoder perspective:

dimensionality reduction w/

$$z = e(x) \in \{1, \cdots, k\}$$
$$\hat{x} = d(z) = \mu_z$$

with dictionary $\{\mu_1, \cdots, \mu_z\}$

# K-means

Tries to assign datapts to clusters s.t. the within-cluster distances are small

Data: $\underline{X} \in \mathbb{R}^{n \times p}$

# clusters: $k \in \mathbb{N}$

Let $C(i) = 1, \cdots, k$    $C: \{1, \cdots, n\} \rightarrow \{1, \cdots, k\}$
indicate the assignment of pt $i$ to a cluster

Quality of a Clustering $C$ is defined as the within-cluster diffs:

$$W(C) = \sum_{i=1}^{n} \| X_i - \mu_{C(i)} \|_2^2$$

Where $\mu_j = \frac{1}{n_j} \sum_{i: C(i) = j} X_i$

$n_j = \sum_{i: C(i) = j} 1$

$$= \sum_{i=1}^{k} \sum_{i: C(i) = i} \| X_i - \mu_i \|_2^2$$

$$= \frac{1}{2} \sum_{j=1}^{k} \frac{1}{n_j} \sum_{i: c(i)=j} \sum_{i': c(i')=j} \|x_i - x_{i'}\|_2^2$$

Want best $C$ interms of $W(c)$

How many $c$'s are there?

  # ways to assign $n$ things to $k$ buckets

  $\subseteq$ Stirling # of $2^{nd}$ kind

   $= HUGE!$

    e.g. for $k=4$,   $S(10,4) = 34105$

                 $S(19,4) > 10^{10}$

$k$-means : greedy, iterative approach to this hard optim problem

Start w/ some initial clustering $C_0$

For $t=1, 2, \cdots$ :

  1. Compute the cluster means for $C_{t-1}$

$$\mu_j^{(t-1)} = \frac{1}{n_j^{(t-1)}} \sum_{i: C_{t-1}(i)=j} x_i \qquad \forall_{j=1, \cdots, k}$$

  2. Reassign each pt $x_i$ to its closest center

$$C_t(i) = \operatorname*{argmin}_{j=1, \cdots, k} \|x_i - \mu_j^{(t-1)}\|_2^2$$

  3. Repeat

$\underline{Obs\ 1}$: Different initializations    (i.e. $C_0$)

lead to different solns.

Soln: Try diff (random) starts
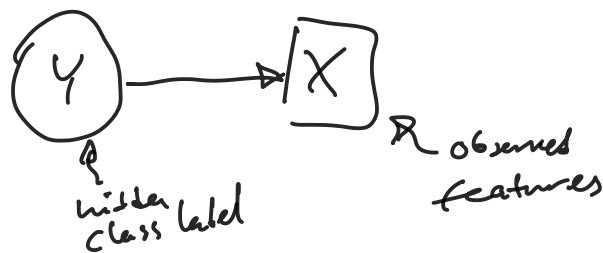
& pick the result

w/ best $W(c)$

__Obs 2:__ K-means will terminate/converge/stop moving

in a finite # of iterations.

Might not be the global opt (see obs 1)

## Soft clustering

K-means alg assign each pt to

exactly one cluster

— hard clustering

Sometimes not clear that there's

a clear cut distinction into clusters

Soft clustering: assign % membership



Y
p
hidden
class label

X
R observed
features

like we got a supervised learning data set

but Y col dropped (even in training)

# Fit Gaussian Mixture Model
# W/ EM algorithm

Suppose K=2

GMM: Let del distrib as follows

GMM:

$$X_0 \sim N(\mu_0, \sigma_0^2) \qquad X_1 \sim N(\mu_1, \sigma_1^2)$$

$$Y \sim Ber(\pi)$$

$$X = (1-Y)X_0 + YX_1 = \begin{cases} X_0 & Y=0 \\ X_1 & Y=1 \end{cases}$$

<u>Or:</u> $(X|Y=y) \sim N(\mu_y, \sigma_y^2)$

<u>Approach:</u> find $\Theta = \{\pi, \mu_0, \sigma_0^2, \mu_1, \sigma_1^2\}$

that best fit my data

(i.e., maximize the likelihood)

& then we'll use $\hat{\Theta}$ to

Compute $\hat{P}_{\hat{\Theta}}(Y=y|X)$ — our soft cluster members[hi]p in cluster $y$

$$P(Y=1|X=x) = \frac{P(X=x|Y=1)\,P(Y=1)}{P(X=x|Y=1)P(Y=1) + P(X=x|Y=0)\,P(Y=0)}$$

$$= \frac{\pi\,\varphi\left(\frac{x-\mu_1}{\sigma_1}\right)}{\pi\,\varphi\left(\frac{x-\mu_1}{\sigma_1}\right) + (1-\pi)\,\varphi\left(\frac{x-\mu_0}{\sigma_0}\right)}$$

$$= \text{cluster responsibility}$$
(soft membership to cluster 1)

<u>Fitting this w/ EM</u>

Given observation $X_1, \cdots, X_n$ the log-like is

$$\ell(\Theta; \underline{X}) = \sum_{i=1}^{n} \log\left((1-\pi)\,\varphi\left(\frac{x_i-\mu_0}{\sigma_0}\right) + \pi\,\varphi\left(\frac{x_i-\mu_1}{\sigma_1}\right)\right)$$

We want $\theta$ to max $l(\theta; \underline{X})$    $\left(\text{or } \min_{} -l(\theta; \underline{X})\right)$

Actually: hard optim problem