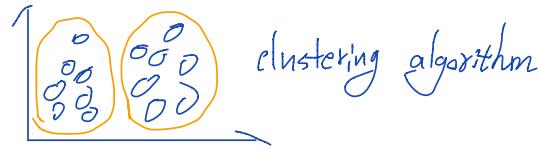


clustering — Unsupervised Learning Intro

Unsupervised Learning



Training set $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

Application of clustering.

- Market segmentation
- Social Network analysis
- Organize Computing clusters
- Astronomical data analysis
-

clustering is an example of unsupervised Learning

o Clustering — k-means algorithm

- Input:

- K (number of clusters)
- Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- $x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

- Procedure

- Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$
 $\in \mathbb{R}^n$
- Repeat {
 - for $i=1$ to m (cluster assignment)

$c^{(i)}$: = index (from 1 to K) of cluster centroid
 closest to $x^{(i)}$

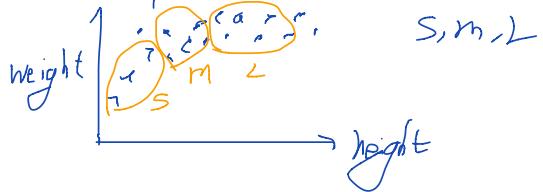
for $k=1$ to K : (update centroid)

μ_k : = average (mean) of points assigned to
cluster k

}

If no points assigned to a certain centroid during a certain step, then eliminate the corresponding centroid.
or randomly re-initialize that centroid

K-means for non-separated clusters



clustering — optimization objective

$- c^{(i)}$ = index of cluster (1, 2, ..., K) to which example $x^{(i)}$ is
currently assigned

$- \mu_k$ = cluster centroid k ($\mu_k \in \mathbb{R}^n$)

$- \mu_{c(i)}$ = cluster centroid of cluster to which example $x^{(i)}$ has
been assigned

Optimization objective

$$T(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_r) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)} \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

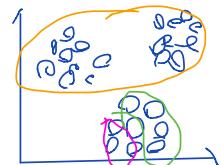
cost function / distortion of k-mean algorithm

- clustering – Random initialization

- should have $k < m$
- randomly pick K training examples
- set μ_1, \dots, μ_K equal to these K examples

Problem of local optima

Assume $K=3$, could end up with:



Solution: Run k-means clustering several times

for $i=1$ to 100 { → 50 - 100

randomly initialize k-means

Run k-means, Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$

Compute cost function (distortion)

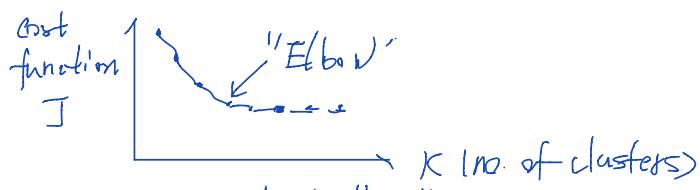
$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

}

Pick clustering that give lowest cost $J(\dots)$

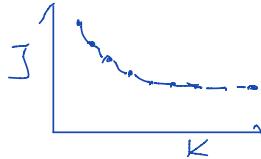
- clustering – choosing the number of clusters

- Elbow method



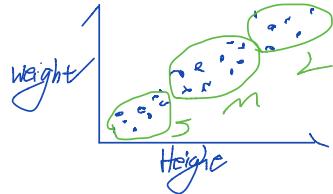
choose K as where corresponding to the "Elbow"

* but "Elbow" is not always obvious, e.g.

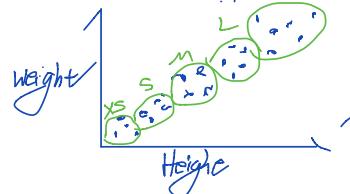


- Base on later purpose

Sometimes you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose, e.g. T-shirt size:



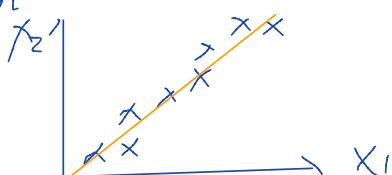
size: S, M, L



size: XS, S, M, L, XL

Dimensionality Reduction

◦ Motivation 1: Data Compression



$$x^{(1)} \in \mathbb{R}^n \rightarrow z^{(1)} \in \mathbb{R}$$

$$x^{(2)} \in \mathbb{R}^n \rightarrow z^{(2)} \in \mathbb{R}$$

$$x^{(3)} \in \mathbb{R}^n \rightarrow z^{(3)} \in \mathbb{R}$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots$$

$$\Rightarrow \longrightarrow 1 \triangleright$$



* This is purely a data projection

$3D \rightarrow 2D$: Project 3D data to 2D plane

o Motivation II: Data visualization

reduce dimensionality to view high dimension data on 2D or 3D or 1D

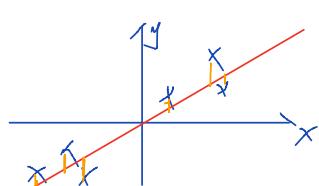
Principal Component Analysis (PCA)

o Problem formulation

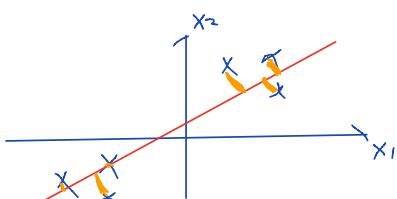
Reduce $2D \rightarrow 1D$: find a direction (a vector $u^{(1)} \in \mathbb{R}^2$) on to which to project the data so as to minimize the projection error.

Reduce from $nD \rightarrow kD$: find k vectors $u^{(1)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

* PCA is not Linear Regression



Linear regression
 $\sum x_i \text{ vs } y$



PCA
 $\sum x_i \text{ vs } x_j, \text{ no } y$

Vertical distance
Projection distance

o PCA Algorithm

1) Data preprocessing (feature scaling / mean normalization)

- $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$

- Replace each $x_j^{(i)}$ with $x_j^{(i)} - \mu_j$

- If different features on different scales, scale features to have comparable ranges

Problem: how to compute $\mu^{(i)}$ and Σ

2) Algorithm: reducing data from $n \rightarrow k \rightarrow$:

- Compute "Covariance matrix"

$$\Sigma = \frac{1}{m} \sum_{i=1}^n \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^T}_{1 \times n} \Rightarrow n \times n$$

- Compute "eigenvectors" of matrix Σ :

$$[U, S, V] = \text{svd}(\Sigma) \quad (\text{In mathematica})$$

svd: Singular Value Decomposition

$$U(n \times n) = \left[\begin{array}{cccc} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & \dots & | \end{array} \right] \text{ eigenvectors}$$

Take the first k eigenvectors

- Take first k eigenvectors $\dots \rightarrow \begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}$

$$U \quad \text{U reduce} = U_{n \times k}$$

$$\tilde{Z} = U_{\text{reduce}}^T X$$

PCA actually minimizes squared projection error

• Reconstruction from compressed representation

How do you go back from $\tilde{Z}^{(i)}$ to $X^{(i)}$?

$$\tilde{Z} \in \mathbb{R}^k \rightarrow X \in \mathbb{R}^n$$

$$\tilde{Z} = U_{\text{reduce}}^T X \Rightarrow X_{\text{approx}}^{(i)} = U_{\text{reduce}}_{n \times k} \cdot \tilde{Z}^{(i)}$$

• Choosing the number of principal component K

Averaged squared projection error: $\frac{1}{m} \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2$

Total variation in the data: $\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2$

Typically, choose K to be smallest value so that:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

("99% of variance is retained")

Algorithm:

try $k = 1, 2, \dots$

Compute $U_{\text{reduce}}, \tilde{Z}^{(1)}, \dots, \tilde{Z}^{(m)}, X_{\text{approx}}^{(1)}, \dots, X_{\text{approx}}^{(m)}$

check if $\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2} \leq 0.01$

if true, break. find good K value

if True, return \hat{y}

if False, continue.

$$[V, S, Y] = \text{Svd}(\text{Sigma})$$

$$S = \begin{bmatrix} s_{11} & s_{12} & \dots \\ & s_{22} & \\ & & \ddots & s_{nn} \end{bmatrix}$$

for a given k :

$$\frac{\frac{1}{m} \left(\sum_{i=1}^m \|x^{(i)} - x^{(i)}_{\text{approx}}\|^2 \right)^{1/2}}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} = 1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.9$$

only need to update numerator for different K

o Advice for applying PCA

Use PCA to speed up learning algorithm

Supervised learning Speed up

$$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$$

Extract inputs:

unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{m \times n}$

\Downarrow PCA

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{m \times k}$$

New training set:

$$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$$

New test example

$$x \xrightarrow{\text{PCA}} z$$

$$y = h_n(z) \rightarrow \text{predictions}$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets

Application of PCA:

- Reduce memory/disk needed to store data
- Speed up learning algorithm
(choose K by % of variance retain)
- visualization
 $K=2$ or $K=3$ for visualize high dimensional data

Bad use of PCA I:

Bad thinking: PCA \rightarrow reduce num of features
 \rightarrow prevent overfitting

Actually: This might be ok, but isn't a good way to address overfitting. Use regularization instead

Because: PCA doesn't information of y involved, it might through away some valuable information.
while regularization has all information retained.

Bad use of PCA II:

Design of ML system (Bad)

- 1) Get training set $(x^{(1)}, y^{(1)})$, ..., $(x^{(m)}, y^{(m)})$
- 2) Run PCA to reduce $x^{(i)}$ to $z^{(i)}$
- 3) Train
- 4) Test on test set: map $x_{\text{test}}^{(i)}$ to $z_{\text{test}}^{(i)}$, run $h(z^{(i)})$

Suggestion:

Before implementing PCA, first try running ML on the raw data. Only if that doesn't do what you want, then implement PCA.