

Machine Learning & Data Mining CS/CNS/EE 155

Lecture 9: Clustering & Dimensionality Reduction

Lecture 9: Clustering & Dimensionality Reduction

Kaggle Competition

• Released soon

• Teams of 2-3

- Competition closes Tuesday Feb 12th, 2pm
 - Winners announced in class
 - Report due Feb 14th, 9pm

Topic Overview

Supervised Learning

Linear Models	Overfitting	Loss Functions
Non-Linear Models	Learning Algorithms & Optimization	Probabilistic Modeling

Unsupervised Learning

Lecture 9: Clustering & Dimensionality Reduction

Today (Unsupervised Learning)

• Clustering

- Dimensionality Reduction
 - Matrix Factorization

What is Clustering?

• Clustering is the process of grouping data points into "clusters".

- High intra-cluster similarity
- Low inter-cluster similarity

Example



Example



Unsupervised Learning

- Given: unlabeled data:
 - Only input features
 - No labels

$$S = \left\{ x_i \right\}_{i=1}^N$$

- **Goal:** find hidden structure/patterns
 - E.g., hidden structure is a clustering of data
 - A generative model of data P(x)
 - Discussed further in future lectures
 - I.e., a low dimensional summary of the data

Why is Clustering Useful?

- Clustering is a "summary" of data
 - Can just inspect cluster centers
 - Or inspect a few data points per cluster

Images Related to "Pluto"



Image Source: http://research.microsoft.com/en-us/people/jrwen/mm04.pdf

Lecture 9: Clustering & Dimensionality Reduction

Why is Clustering Useful?

- Clustering is a "summary" of data
 - Can just inspect cluster centers
 - Or inspect a few data points per cluster
- Compact pre-processing of data before supervised training































K-Means Objective



EM Algorithm for K-Means (Expectation/Maximization)



- Estimate C_k
- Estimate cluster membership
- M-Step
 - Estimate c_k
 - Estimate model parameters

$$\begin{array}{l} \text{E-Step} \\ \text{argmin} \\ S = C_1 \cup \dots \cup C_K, \ \{c_1, \dots, c_K\} \\ k \\ k \\ x \in C_k \end{array} \| x - C_k \|^2 \qquad S = \{x_i\}_{i=1}^N
\end{array}$$

• For each x:

– Assign to cluster C_k with smallest distance to c_k



M-Step

$$\underset{S=C_1\cup\ldots\cup C_K, \{c_1,\ldots,c_K\}}{\operatorname{argmin}} \sum_{k} \sum_{x\in C_k} \|x-c_k\|^2 \qquad S = \{x_i\}_{i=1}^N$$

• For each c_k:

- Compute $c_k = mean(C_k)$



Interpretation

- Summarize data by cluster membership
- Learn clustering to minimize intra-cluster variance
 - "Best reconstruction of the data"



Recap: K-Means

- Centroid-based Clustering
 - Defines clusters using a notional of centrality
 - E.g., all items in the cluster must be close to each other
- Solve using EM algorithm
 - Also probabilistic variant (Gaussian Mixture Models)
- Useful when centrality assumption is good
 - But bad when centrality assumption is bad...

Thought Experiment

What is good clustering?



Linkage Based Clustering (Hierarchical Clustering)

K-Means used centroid clustering structure
 – Clustered data points are "close" to cluster center

- Sometimes a linkage structure is better...
 - Employ hierarchical clustering
 - E.g., agglomerative clustering

Agglomerative Clustering



Agglomerative Clustering

- Equivalent to finding minimum spanning tree
 - Kruskal's Algorithm
 - <u>http://en.wikipedia.org/wiki/Kruskal%27s_algorithm</u>
- Order that edges are added defines the cluster hierarchy

• Equivalent to finding a binary tree partitioning with progressively smaller partition distances

Recap: Clustering

- Unsupervised learning
 - Finds the clustering structure of input features
- Centroid based
 - Clusters should be clumped together
 - K-Means
- Linkage Based
 - Clusters can be organized hierarchically
 - Agglomerative Clustering
- Works great when clustering assumption is good!

Limitations of Clustering





Summarizing Data

- Summarize data using smaller #attributes $S = \{x_i\}_{i=1}^N$
- Clustering: summarize data via clusters
 - K-Means: summarize via cluster membership
 - Gaussian Mixture Model: Summarize via distribution over K clusters
- PCA: summarize via orthogonal projections
 - Define new feature representation
 - Rotation + Projection





New Feature Representation!

Lecture 9: Clustering & Dimensionality Reduction

Orthogonal Matrix

- A matrix U is orthogonal if $UU^T = U^TU = I$
 - For any column u: $u^{T}u = 1$
 - For any two columns u, u': $u^Tu' = 0$
 - U is a rotation matrix, and U^T is the inverse rotation
 - If $x' = U^T x$, then x = Ux'



Properties of Orthogonal Matrices

- $x' = U^T x$, x = Ux'
- Norm preserving:

$$x'^{T} x' = \left(U^{T} x\right)^{T} \left(U^{T} x\right) = x^{T} U U^{T} x = x^{T} x$$

• Preserves Total Variance:

$$\sum_{d=1}^{D} \sum_{i=1}^{N} \left(x_{i}^{(d)} \right)^{2} = \sum_{d=1}^{D} \sum_{i=1}^{N} \left(x_{i}^{(d)} \right)^{2}$$

Assuming zero mean



Summarize Using 1 Feature?



Summarize Using 1 Feature?

Summarize Using 1 Feature?



PCA Formal Definition

• Define M=matrix of all data:

$$X = [x_1, \dots, x_N] \in \operatorname{Re}^{D \times N}$$

• Mean center:

$$\overline{X} = X - \left[\overline{x}, ..., \overline{x}\right]$$

• PCA:

$$\overline{X}\overline{X}^{T} = U\Lambda U^{T}$$
Symmetric
Orthogonal
Diagonal

Properties of PCA



- Each column of U is an Eigenvector
- Each λ is an Eigenvalue

 $-\lambda_1 \geq \lambda_2 \geq ... \geq \lambda_D$

$$\left(XX^{T}\right)u_{d} = \lambda_{d}u_{d}$$

Interpretation

Feature Covariance Matrix:

$$\Sigma = XX^T = U\Lambda U^T$$

Assuming zero mean

PCA Solution

- $\Sigma_{dd'}$ is the covariance of features d & d' in training data.
- The first column u₁ is the single direction of greatest variation

 $-\lambda_1$ is the total variation along u_1 :

$$\lambda_{1} = \sum_{i=1}^{N} \left(u_{1}^{T} x_{i} \right)^{2} = \sum_{i=1}^{N} \left(x_{i}^{\prime(1)} \right)^{2}$$



Interpretation Continued

- The first column u₁ is the single direction that minimizes the squared loss of reconstructing the original x's
 - I.e., minimizes the amount of residual variation
- One can prove that:

$$u_{1} = \underset{u: u^{T}u=1}{\operatorname{argmin}} \sum_{i=1}^{N} \left\| x_{i} - uu^{T} x_{i} \right\|^{2}$$

"Residual"

(From definition in previous slide)



Definition: u_1 is the direction that captures the most variation $u_1 = \underset{u: u^T u=1}{\operatorname{arg\,max}} \sum_{i=1}^N \left\| u^T x_i \right\|^2$

Step 1: for any x, its residual direction is orthogonal to u₁

Residual:
$$x - u_1 u_1^T x$$

 $(x - u_1 u_1^T x)^T u_1 = x^T u_1 - x^T u_1 u_1^T u_1 = x^T u_1 - x^T u_1 = 0$

Step 2: establish relationship and complete proof

$$\sum_{i=1}^{N} \left\| x_{i} - uu^{T} x_{i} \right\|^{2} = \sum_{i=1}^{N} \left(x_{i} - uu^{T} x_{i} \right)^{T} \left(x_{i} - uu^{T} x_{i} \right) = \sum_{i=1}^{N} \left(x_{i}^{T} x_{i} - 2x_{i}^{T} uu^{T} x_{i} + x_{i}^{T} uu^{T} uu^{T} x_{i} \right)$$
$$= \sum_{i=1}^{N} \left(x_{i}^{T} x_{i} - x_{i}^{T} uu^{T} x_{i} \right) \qquad = \sum_{i=1}^{N} \left(x_{i}^{T} x_{i} \right) - \sum_{i=1}^{N} \left(x_{i}^{T} uu^{T} x_{i} \right)$$

Interpretation Continued

Find the u₁ that minimizes the residual squared norm:



Solving PCA (Iterative Algorithm)

- Given: $X = [x_1, ..., x_N] \in \operatorname{Re}^{D \times N}$ Assuming zero mean
- Init: $X_1 = X$

• For d=1,...,D
- Solve:
$$u_d = \underset{u: \ u^T u=1}{\operatorname{arg\,min}} \|X_d - uu^T X_d\|_{Fro}^2$$

$$X_{d+1} = X_d - u_d u_d^T X_d$$

Property of PCA $XX^T = U\Lambda U^T$

 The first K columns of U are guaranteed to be the K-dimensional subspace that captures the most variability of X

• We just proved K=1 a few slides ago

Dimensionality Reduction

- Solve PCA: $XX^T = U\Lambda U^T$
- Use first K columns of U to create K-dim representation:

$$x' = U_{1:K}^T x$$

• This creates a compact summary of original dataset

- E.g., K = 50, D = 1,000,000

Example: Eigenfaces



PCA on a corpus of faces.Every pixel is a "feature"Visualizing the top Eigenvectors of U

http://www.cs.princeton.edu/~cdecoro/eigenfaces/

Lecture 9: Clustering & Dimensionality Reduction

Example: Eigenfaces



Visualizing Projection Visualizing Projection Using top K Eigenvectors: $U_{1\cdot K}U_{1\cdot K}^T X$

http://www.cs.princeton.edu/~cdecoro/eigenfaces/

Lecture 9: Clustering & Dimensionality Reduction

CS 155 Eigenfaces



Avg Face

























5 eigenfaces

10 eigenfaces





5 eigenfaces

10 eigenfaces









30 eigenfaces

75 eigenfaces









150 eigenfaces







75 eigenfaces









50 eigenfaces



150 eigenfaces



Singular Value Decomposition $X = U \Sigma V^{T}$ $M = U \Sigma V^{T}$

- SVD operates on X, as opposed to XX^T
- Equivalence between SVD & PCA

$$XX^{T} = (U\Sigma V^{T})(U\Sigma V^{T})^{T} = U\Sigma V^{T}V\Sigma U^{T} = U\Sigma^{2}U^{T}$$

• V corresponds to new representation x'

• Flatten each image into vector



Each Column is Image

Mean center



Mean



• Singular Value Decomposition: $X' = U\Sigma V^T$



• Merging Σ into U and V: $X' = U\Sigma V^T = U'V'^T$



Interpreting U & V

- Each col of U' is an "Eigenface"
- Each col of V'^{T} = coefficients of a student



24.138 , -29.3105



-24.9924 , -2.3168

-50.2606 , -16.9522



15.4892, 46.3845

6.1785 , 3.4943



13.5041 , 22.731





-51.484, 8.5238





6.7881 , -2.3789



6.5127 , 7.5933



36.4135, -3.6669





-20.3981 , -16.4748



Limitations of Eigenfaces

- Each dimension is a pixel (& color channel)
 - Not semantically meaningful
 - Squared reconstruction error in pixel space



- Suppose each dimension had more meaning
 - E.g., dim 1 = location of left eye
 - Then U components would have cleaner visualization

Summary

- Clustering & PCA (and SVD) reduce the dimensionality of data representation.
- For each data point
 - Store K numbers
 - Cluster membership probabilities
 - Coefficients in K-dimensional projection
- Nice visualization & interpretation?
 Depends on semantics of raw dimensions...

Next 2 Lectures

• Latent Factor Models

- Matrix Factorization with Missing Values
 E.g., the "Netflix Problem"
- Embeddings