Clustering, Ctd., and Dimensionality Reduction

https://tinyurl.com/cis5190-10-31-2022

Osbert Bastani and Zachary G. Ives CIS 4190/5190 – Fall 2022

https://unsplash.com/photos/UmmFQ1C5m8w

OWEEN



- Homework 4 due November 2, 8pm
- Quiz 8 due November 3, 8pm

Recall from Last Week: K-Means Clustering

K-Means (K, X)

- Randomly choose *K* cluster center locations (centroids)
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster



K-Means Clustering

K-Means (K, X)

- Randomly choose *K* cluster center locations (centroids)
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster



K-Means Clustering

K-Means (K, X)

- Randomly choose *K* cluster center locations (centroids)
- Loop until convergence, do:
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster

Optimizer: "Alternating Minimization"

K-means finds a local optimum of the following objective function:



$$\arg\min_{S} \sum_{k=1}^{K} \sum_{x \in S_{k}} \|x - \boldsymbol{\mu}_{K}\|_{2}^{2}$$

where $S = \{S_1, ..., S_K\}$ are sets corresponding to disjoint clusters, and the clusters together include all samples.

K-Means Clustering Convergence



https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/

Generalizing to Other Distance Functions

K-Means Objective Function:

$$\arg\min_{\boldsymbol{S}}\sum_{k=1}^{K}\sum_{\boldsymbol{x}\in S_{k}}\|\boldsymbol{x}-\boldsymbol{\mu}_{K}\|_{2}^{2}$$

But it is possible to define k-means with other notions of pairwise distance between samples too. For example:

$$\left(\sum_{d} |x_{1d} - x_{2d}|^{1}\right)^{\frac{1}{1}} Q: What would have to change in the algorithm?$$

$$\ell_{1} \text{ distance}$$

$$\sum_{d} |x_{1d} - x_{2d}|$$

$$\sum_{d} |x_{1d} - x_{2d}|$$

$$(2019-2x_{2d}^{x_{3}} \text{ arman, 0. Bastani, 2. Ives}$$

$$(K-medoids" clustering)$$

$$(K-me$$

Weaknesses of k-Means

Centroids can be heavily affected by outliers

- Remove (consistent) outliers, or
- Cluster over a random sample



k-means finds *spherical* clusters

k-means *converges* but can get trapped in local optima (thus isn't deterministic and depends on initialization...)

© 2019-22 D. Jayaraman, O. Bastani, Z. Ives

Figure: https://www.slideshare.net/AndresMendezVazquez/25-machinelearning-unsupervised-learaning-kmeans-kcenters

K-Means is Too Sensitive to Initialization

Alternative strategies:

- 1. Do many runs of K-Means, each with different initial centroids, and pick the best
- 2. Pick initial centroids using a better method than random choice

K-means+ + initialization

- Choose a data point uniformly at random as the first centroid
- Loop for 2: *K*, do:
 - Let D(x) be the distance from each point x to the closest centroid
 - Choose data point x randomly $\propto D(x)^2$ as the next centroid. Higher chance to pick points that are far from previous centroids.

K-means++ Illustrated

Place the initial centroids **far away from one another**:

- Initialize an empty set M (for storing selected centroids);
 Randomly select the first centroid from the input sample and assign it to M
- 2. For each x_i that is not in M, find the distance $D(x_i)$ to the closest centroid in M
- 3. Choose one new data point at random as a new centroid using probability distribution $\sim D(x)^2$
- 4. Repeat (2) and (3) until K centroids have been chosen

Then do "classic" k-means





- k-Means is highly scalable
- can compute distances from centroids in a highly parallel manner
- can compute centroids by grouping on the centroid-assignments

How Many Clusters?

Measuring the Performance of k-Means

How can we evaluate how good our clustering is? Some options:

- Evaluation using the k-means objective itself
- Comparing to class labels (for a subset of data) Sometimes possible
- Subjective evaluation by a human domain expert

• • •

"Knee Point" For Selecting K

Elbow Method For Optimal k



https://blog.cambridgespark.com/how-to-determine-the-optimal-number-of-clusters-for-k-means-clustering-14f27070048f

k-Means Clustering

- Non-deterministic (may get different output based on different start values) but guaranteed to converge
- Iterative algorithm with two sub-steps (after random cluster centroids chosen):
 - **1.** Assign points to nearest cluster
 - **2.** Recompute cluster centroid
- Select number of clusters by exploring error (distortion)

Q: how else might we formalize the problem of finding clusters?

Hierarchical Clustering

Hierarchical Clustering

Instead of repeating until convergence based on global measures, like k-Means – Let's consider a greedy *local* algorithm, that iteratively makes choices

- **1.** Start with **single-item clusters**, then build successively bigger and bigger clusters: *agglomerative*
- **2.** Or: **start with one cluster**, break into the most logical sub-clusters, repeat: *divisive*

These are called *hierarchical clustering* approaches...

Basic Intuition

- Agglomerative: In each iteration, we find the closest clusters and merge them
- Divisive: In each iteration, we find the two most distant sub-clusters and split there

But: we know how to compute differences in points – need to generalize this to computing distances among clusters

Potential Cluster Differences

Single Linkage: Compute distances between the **most similar** members for each pair of clusters

Merge the clusters with the smallest **min-distance**

Complete Linkage: Compute distance between the **most dissimilar** members for each pair of clusters

Merge the clusters with the smallest **max-distance**



How Do We Do this Efficiently? Considering Agglomerative Case

Really inefficient to iterate over each point and compute its distance to every other point

• $O(n^2)$ computations, which is bad for big data!

Idea: precompute and memorize:

- Compute a distance matrix where distance[i,j] is the distance between nodes i,j
- We'll update this matrix every time we merge

Pseudocode: Agglomerative Clustering with Complete Linkage single

Agglomerative Clustering:

1. Compute distance matrix **dist** between all pairs of points (a,b). [SciPy pdist]

2. Repeat:

Iterate over pairs of clusters A, B, compute their distance: Look at all pairwise distances dist[a,b] between $a \in A, b \in B$ Merge the pair of clusters with *min* distance between most distant members Update the distance matrix Look at all pairwise distance matrix

Until a single cluster remains

Example: Reconstructing Phylogenetic Trees

Q: Is a panda a bear?

Or is its closest relative the red panda, which is related to the raccoon?









https://www.nwf.org/Educational-Resources/Wildlife-Gu

https://www.smithsonianmag.com/science-nature/eight-amazing-facts-about-red-pandas-180979708/

https://towardsdatascience.com/hierarchical-clustering-and-its-applications-41c1ad4441a6

Example: Reconstructing Phylogenetic Trees



https://towardsdatascience.com/hierarchical-clustering-and-its-applications-41c1ad4441a6

Summary of Hierarchical Clustering

Hierarchical clustering is often easier to visualize and interpret with a taxonomy "Dendrogram" plots We don't need to specify the number of clusters up front!

Limitation: Doesn't scale well to big problems

General Notes on Distance Measures

- In studying Clustering techniques, we assume that we are given a matrix of distances between all pairs of data points.
- We assume that the input to the problem is:



- In studying clustering techniques, we assume that we are given a matrix of distances between all pairs of data points.
- A distance measure (*metric*) is a function $d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ that satisfies:

$$1. d(x, y) \ge 0, d(x, y) = 0 \iff x = y$$

$$2. d(x, y) + d(y, z) \ge d(x, z)$$

$$3. d(x, y) = d(y, x)$$

- For the purpose of clustering, sometimes the distance (similarity) is not required to be a metric
 - No Triangle Inequality
 - No Symmetry

Examples:

Euclidean Distance:

$$d(x,y) = \sqrt{(x-y)^2} = \sqrt{(x-y)^T (x-y)} = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

Manhattan Distance:

$$d(x, y) = |x - y| = \sum_{i=1}^{d} |x_i - y_i|$$

Infinity (Sup) Distance:

$$d(x, y) = \max_{1 \le i \le d} |x_i - y_i|$$

- Notice that if d(x, y) is the Euclidean metric, $d^2(x, y)$ is not a metric
- But can be used as a measure (no triangle inequality)

© 2019-22 D. Jayaraman, O. Bastani, Z. Ives

- Examples:
 - Euclidean Distance:

$$d(x,y) = \sqrt{(x-y)^2} = \sqrt{(x-y)^T (x-y)} = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

Manhattan Distance:

$$d(x, y) = |x - y| = \sum_{i=1}^{d} |x_i - y_i|$$

Infinity (Sup) Distance:

$$d(x, y) = \max_{1 \le i \le d} |x_i - y_i|$$

Euclidean: $(4^2 + 2^2)^{\frac{1}{2}} = 4.47$ Manhattan: 4 + 2 = 6Sup: Max(4,2) = 4



Comparing clustering algorithms



https://scikit-learn.org/stable/modules/clustering.html#clustering

Summary of Clustering

- Critical to understanding the structure of our data
- Often useful for creating high-level features useful for supervised learning
- We saw two approaches: k-Means vs hierarchical clustering

Dimensionality Reduction

Osbert Bastani and Zachary G. Ives CIS 4190/5190 – Fall 2022

The Next Key Question: How Do We Get the "Best" Features?

For a variety of reasons, we may want to reduce the number of dimensions:

- Reduce the complexity of our learning problem
- Remove multicollinearity / correlated features
- Remove less informative features (results in simpler model)
- Visualize the features

Key problem: mapping from D-dimensions down to a D'-dimensional subspace (D' << D)

Consider: Visualizing High-Dimensional Data

	LotFrontage	LotArea	Street	LotShape Utilitie	s LandSlop	e OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	ExterQual	ExterCond	BsmtQual	BsmtExposure	BsmtFinType1	BsmtFinSF1	BsmtFinType2	Sa	aleCondition_Abnorml
0	65.0	8450	2	4	4	3 7	5	2003	2003	196.0	4	3	4	0	6	706	1		0
1	80.0	9600	2	4	4	3 6	8	1976	1976	0.0	3	3	4	3	5	978	1		0
2	68.0	11250	2	3	4	3 7	5	2001	2002	162.0	4	3	4	1	6	486	1		0
3	60.0	9550	2	3	4	3 7	5	1915	1970	0.0	3	3	3	0	5	216	1		1
4	84.0	14260	2	3	4	3 8	5	2000	2000	350.0	4	3	4	2	6	655	1		0
5	85.0	14115	2	3	4	3 5	5	1993	1995	0.0	3	3	4	0	6	732	1		0
6	75.0	10084	2	4	4	3 8	5	2004	2005	186.0	4	3	5	2	6	1369	1		0
7	0.0	10382	2	3	4	3 7	6	1973	1973	240.0	3	3	4	1	5	859	4		0
8	51.0	6120	2	4	4	3 7	5	1931	1950	0.0	3	3	3	0	1	0	1		1
9	50.0	7420	2	4	4	3 5	6	1939	1950	0.0	3	3	3	0	6	851	1		0
10	70.0	11200	2	4	4	3 5	5	1965	1965	0.0	3	3	3	0	3	906	1		0
11	85.0	11924	2	3	4	3 9	5	2005	2006	286.0	5	3	5	0	6	998	1		0
12	0.0	12968	2	2	4	3 5	6	1962	1962	0.0	3	3	3	0	5	737	1		0
13	91.0	10652	2	3	4	3 7	5	2006	2007	306.0	4	3	4	2	1	0	1		0
14	0.0	10920	2	3	4	3 6	5	1960	1960	212.0	3	3	3	0	4	733	1		0
15	51.0	6120	2	4	4	3 7	8	1929	2001	0.0	3	3	3	0	1	0	1		0
16	0.0	11241	2	3	4	3 6	7	1970	1970	180.0	3	3	3	0	5	578	1		0
17	72.0	10791	2	4	4	3 4	5	1967	1967	0.0	3	3	0	0	0	0	0		0
18	66.0	13695	2	4	4	3 5	5	2004	2004	0.0	3	3	3	0	6	646	1		0
19	70.0	7560	2	4	4	3 5	6	1958	1965	0.0	3	3	3	0	2	504	1		1
20	101.0	14215	2	3	4	3 8	5	2005	2006	380.0	4	3	5	2	1	0	1		0
21	57.0	7449	2	4	4	3 7	7	1930	1950	0.0	3	3	3	0	1	0	1		0
22	75.0	9742	2	4	4	3 8	5	2002	2002	281.0	4	3	4	0	1	0	1		0
23	44.0	4224	2	4	4	3 5	7	1976	1976	0.0	3	3	4	0	6	840	1		0



Data Visualization

Maybe it isn't necessary to visualize all 227 dimensions

Is there a representation better than the raw features?

Idea: find a lower-dimensional subspace that retains most of the information about the original data

There are many methods; our focus will be on Principal Components Analysis



Image : https://arxiv.org/pdf/1703.08893.pdf

34

Dimensionality Reduction Objective

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{ND} \end{bmatrix}_{N \times D}$$

We can write each row (each data sample) x_i as:

$$x_{i} = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix}_{D} = x_{i1} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ D \end{bmatrix} + x_{i2} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ D \end{bmatrix} + \dots + x_{iD} \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}_{D}$$
Projections Original axes

We are looking for a new coordinate system to (approximately) express \mathbf{x}_i : $\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix} \approx f_1(\mathbf{x}_i)\mathbf{v}_1 + f_2(\mathbf{x}_i)\mathbf{v}_2 + \dots + f_{D'}(\mathbf{x}_i)\mathbf{v}_{D'}$ where the new axes \mathbf{v}_d 's are all *D*-dimensional unit norm, and $D' \ll D$

Principal Component Analysis



(Fig: stats.stackexchange)

Orthogonal projection of data onto lower-dimension linear space :

- maximizes variance of projected data (purple line)
- minimizes mean squared distance between data point and projections (sum of blue lines)

Reduce to D' = 1 **dimension?**

We are looking for a new coordinate system to (approximately) express
$$x_i$$
:
 $x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix} \approx (x_i \cdot v_1) v_1 + (x_i \cdot v_2) v_2 + \dots + (x_i \cdot v_{D'}) v_{D'}$
where the new axes v_d 's are all *D*-dimensional unit norm, and $D' \ll D$

Simplest case:
$$D' = 1$$
?
We want to find \boldsymbol{v}_1 and $f_1(\boldsymbol{x}_i)$ such that:
 $\boldsymbol{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix}$ best approximates $f_1(\boldsymbol{x}_i)\boldsymbol{v}_1$
For a given \boldsymbol{v}_1 , $f_1(\boldsymbol{x}_i)$ should be the projection $\boldsymbol{x}_i \cdot \boldsymbol{v}_1 / \|\boldsymbol{v}_1\|_2$
So, only need to find \boldsymbol{v}_1



Objective Function: Maximizing Variance

Find unit vector
$$\boldsymbol{v}_1$$
 (with $\|\boldsymbol{v}_1\|_2 = 1$), to optimize:
Reconstruction
MSE
$$\begin{aligned}
\min_{\|\boldsymbol{v}_1\|_2=1} \frac{1}{N} \sum_i \|(\boldsymbol{x}_i \cdot \boldsymbol{v}_1) \boldsymbol{v}_1 - \boldsymbol{x}_i\|_2^2 \\
\|\boldsymbol{v}_1\|_2=1
\end{aligned}$$
Projection error
$$\max_{\|\boldsymbol{v}_1\|_2=1} \operatorname{variance}(\boldsymbol{x}_i \cdot \boldsymbol{v}_1)$$

Intuitively, if the variance of the projection on v_1 was low, then v_1 would not be very informative about samples x_i . Conversely, directions with high variance projections preserve the most information.



⁽Fig: stats.stackexchange)

So, how to find this direction of maximum variance?

Covariance Matrix



For zero-centered data,

Covariance =
$$C = \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T] = \mathbb{E}\begin{bmatrix} x_{i1}x_{i1} & \cdots & x_{i1}x_{iD} \\ \vdots & x_{ij}x_{ik} & \vdots \\ x_{iD}x_{i1} & \cdots & x_{iD}x_{iD} \end{bmatrix}$$

For any unit vector v_1 ,

variance $(\boldsymbol{x}_i \cdot \boldsymbol{v}_1) = \boldsymbol{v}_1^T C \boldsymbol{v}_1$

To maximize $v_1^T C v_1$, we can set $v_1 = e_1(C)$, the first unit eigenvector of C

More than 1 dimension?

Repeat for d = 1, ..., D'

- Subtract means of all dimensions of *X*
- Compute $C_d = E[x_i x_i^T]$
- Set $\boldsymbol{v}_d = \boldsymbol{e}_1(C_d)$
- Set $x_i = x_i (x_i \cdot v_d) v_d$ (i.e., subtract current reconstructions to compute residuals... like gradient boosting!)

Equivalent to simply: Repeat for d = 1, ..., D'

• Set $v_d = \boldsymbol{e}_d(C_1)$

We are looking to find a new way to express
$$\mathbf{x}_i$$
: "Reconstruction"
 $\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix} \approx f_1(\mathbf{x}_i)\mathbf{v}_1 + f_2(\mathbf{x}_i)\mathbf{v}_2 + \cdots + f_{D'}(\mathbf{x}_i)\mathbf{v}_{D'}$



$$\boldsymbol{x}_{i} = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix} \approx \sum_{d=1}^{D'} (\boldsymbol{x}_{i}, \boldsymbol{v}_{d}) \boldsymbol{v}_{d}$$

So, the new low-dimensional representation is: $f(x_i) = [x_i \cdot v_1, x_i \cdot v_2, ..., x_i \cdot v_{D'}]$

PCA on a 2D Gaussian Dataset



© 2019-22 D. Jayaraman, O. Bastani, Z. Ives

By Nicoguaro - Own work, CC BY 4.0, https://commons.wikimedia.org/w/index.php?curid=46871195

PCA Algorithm

Given data { x_1 , ... x_n }, compute covariance matrix C

- **X** is the *N*×*D* data matrix
- Compute data mean (average over all rows of **X**)
- Subtract mean from each row of *X* (centering the data)
- Compute covariance matrix $C = X^T X$ (*C* is $D \times D$)

PCA basis vectors (new coordinate axes) are given by the eigenvectors of C

- $Q, \Lambda =$ numpy.linalg.eig(C)
- { q_d , λ_d } $_{d=1,...,D}$ are the eigenvectors/eigenvalues of *C* ($\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_D$)

But there are *D* eigenvectors, so where is the dimensionality reduction? • A: Larger eigenvalue \Rightarrow "more important" eigenvectors • 2019-22 D. Jayaraman, O. Bastani, Z. Wes

Dimensionality Reduction

• Can *ignore* the components of lesser significance



- PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10
 You do lose some information, but if the eigenvalues are small, you don't lose much
 - -choose only the first D' eigenvectors, based on their eigenvalues
 - -final data set has only D' dimensions

43

Recap

- Want to reconstruct data approximately in a new coordinate space
- Must find axes of this coordinate space, because the weights on those axes are just projections
- Objective: axes with lowest reconstruction error
 Same as axes with high variance projections
- Solution straight from linear algebra. Axes are eigenvectors of covariance matrix



• Each column of **Q** gives weights for a linear combination of the original features



PCA

Compute $x. e_d$ to get the new representation for each instance x

$$X = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 & \cdots \\ 1 & 1 & 0 & 1 & 1 & \cdots \\ 0 & 0 & 1 & 1 & 1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \\ 1 & 0 & 1 & 0 & 1 & \cdots \end{bmatrix} x_{3} \qquad \hat{Q} = \begin{bmatrix} 0.34 & 0.23 & 0.13 \\ 0.04 & 0.13 & 0.13 \\ 0.93 & \vdots \\ \vdots & 0.23 & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.04 & 0.13 \\ 0.93 & 0.93 \\ \vdots & 0.23 \\ 0.93 & 0.9$$

The new 2D representation for x_3 is given by $[\widehat{x_{31}} = x_3, e_1, \widehat{x_{32}} = x_3, e_2]$:

$$\widehat{x_{31}} = 0.34(0) + 0.04(0) - 0.64(1) + \cdots$$

 $\widehat{x_{32}} = 0.23(0) + 0.13(0) + 0.93(1) + \cdots$

The re-projected data matrix can be conveniently computed as $\hat{X} = X\hat{Q}$

Using PCA for Feature Reduction / Compression

We could use the PCA transformation of the data instead of the original features Keep only D' < D PCA features

PCA tries to retain most of the variance in the data

 So, we're reducing the dataset to features that retain meaningful variations of the data set

Eigenfaces

What happens when you compute the principal components of face images?







Michael Jackson



Hillary Clinton





David Beckham



Dwayne Johnson



Oprah Winfrey





Michael Jordan









Lindsay Davenport



George W Bush

Colin Powell











Vin Diesel







Mary Carey

Surakait Sathirathai



Dean Barkley



Colin Powell



(1000 64×64 images)







Frank Taylor

Eigenfaces

What happens when you compute the principal components of face images?

"Eigenfaces": main directions of deviation from the mean face



Figure #5: mean face



Eigenfaces

Vin Diesel

Rubens Barrichello

Let's try reconstructing these faces with the eigenfaces now!







Richard Myers



Frank Taylor



George W Bush



Colin Powell



Yasser Arafat















Mary Carey



Dean Barkley



Colin Powell



... with 1000 eigenvectors







Richard Myers





George W Bush



Colin Powell



Yasser Arafat





Vin Diesel



Rubens Barrichello

Sarah Price

Noah Wyle







Colin Powell



Surakait Sathirathai







Dean Barkley

© 2019-22 D. Jayaraman, O. Bastani, Z. Ives

52



... with 250 eigenvectors

Lindsay Davenport



Billy Crystal



Richard Myers





George W Bush





Yasser Arafat





Sarah Price

Vin Diesel







Surakait Sathirathai



Rubens Barrichello





Colin Powell





... with 100 eigenvectors

Lindsay Davenport



Billy Crystal



Richard Myers





George W Bush





Yasser Arafat





Vin Diesel



Rubens Barrichello

Sarah Price

Noah Wyle









Mary Carey



Dean Barkley



Colin Powell

... with 50 eigenvectors







Dean Barkley







55

PCA Visualization of Digits



Utility of PCA

- PCA is often used as a preprocessing step for supervised learning
 reduces dimensionality
 - eliminates multicollinearity
- Can also be used to aid in visualization

Beyond PCA: Non-linear dimensionality reduction



Recap: Unsupervised Learning

Basic idea: reduce feature space to a much lower set of dimensions

- Clustering: find structural similarity, return one k-valued higher-level feature
- PCA: find orthonormal dimensions in order of most to least variance
- Can be useful for human inspection (visualization) as well as supervised ML
- Next time: a very prominent datatype text and documents!

- Notice that:
 - Infinity (Sup) Distance < Euclidean <u>Distance < Manhattan Distance</u>:

$$L_{\infty} = \max_{1 \le i \le d} |x_i - y_i| \quad L_2 = \sqrt{(x - y)^2} = \sqrt{\sum_{i=1}^d (x_i - y_i)^2} \quad L_1 = |x - y| = \sum_{i=1}^d |x_i - y_i|$$

• But different distances do not induce same order on pairs of points

