Dimensionality Reduction PCA

Aarti Singh

Machine Learning 10-315 Nov 10, 2021

Slides Courtesy: Tom Mitchell, Eric Xing, Lawrence Saul



High-Dimensional data

High-Dimensions = Lot of Features

Document classification

Features per document =

thousands of words/unigrams
millions of bigrams, contextual
information



Surveys - Netflix

480189 users x 17770 movies

	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6
Tom	5	?	?	1	3	?
George	?	?	3	1	2	5
Susan	4	3	1	?	5	1
Beth	4	3	?	2	4	2

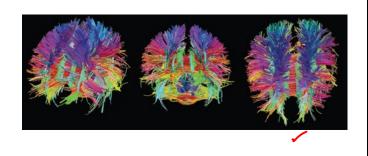
High-Dimensional data

High-Dimensions = Lot of Features

High resolution images millions of pixels

Diffusion scans of Brain 300,000 brain fibers





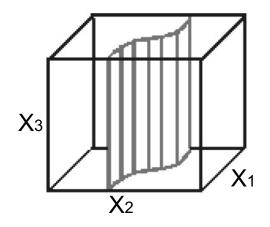


Curse of Dimensionality

- Why are more features bad?
 - Redundant features (not all words are useful to classify a document)
 more noise added than signal
 - Hard to interpret and visualize
 - Hard to store and process data (computationally challenging)
 - Complexity of decision rule tends to grow with # features. Hard to learn complex rules as it needs more data (statistically challenging)

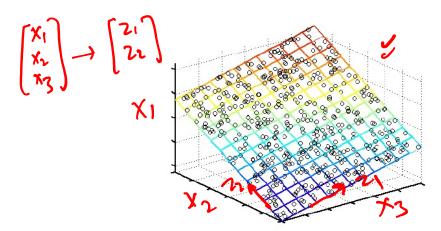
Dimensionality Reduction

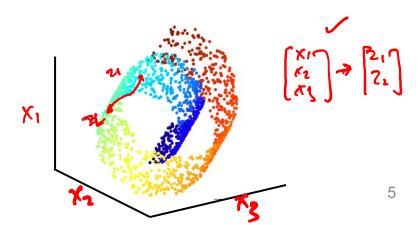
Feature Selection – Only a few features are relevant to the learning task



X₃ - Irrelevant

 Latent features – Some linear/nonlinear combination of features provides a more efficient representation than observed features



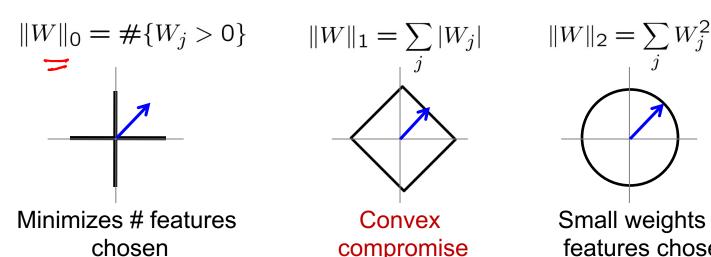


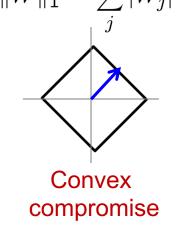
Feature Selection

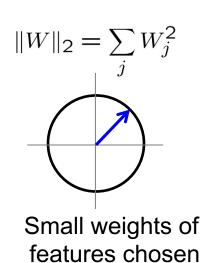


One Approach: Regularization (MAP) Integrate feature selection into learning objective by penalizing number of features with non-zero weights

$$\widehat{W} = \arg\min_{W} \sum_{i=1}^{n} -\log P(Y_{i}|X_{i};W) + \lambda \|W\|$$
 -ve log likelihood penalty







Latent Features

Combinations of observed features provide more efficient representation, and capture underlying relations that govern the data

E.g. Ego, personality and intelligence are hidden attributes that characterize human behavior instead of survey questions

Topics (sports, science, news, etc.) instead of documents -

Often may not have physical meaning

Linear

Principal Component Analysis (PCA)

Factor Analysis

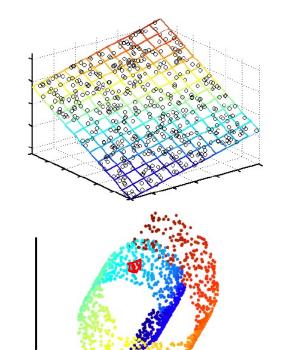
Independent Component Analysis (ICA)

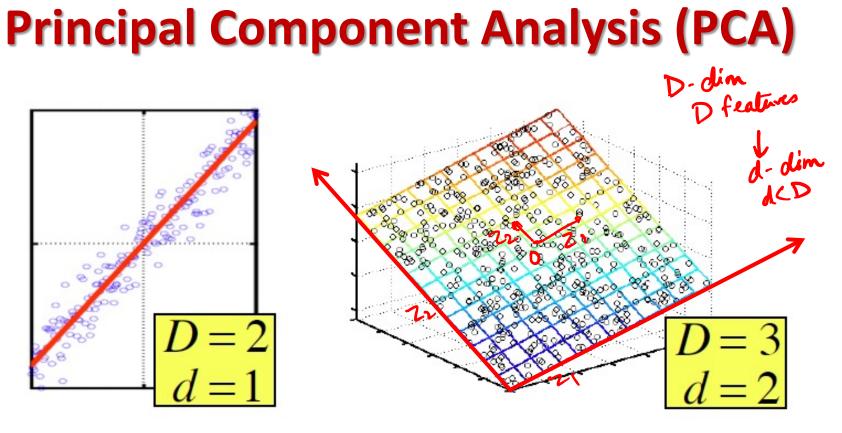
Nonlinear

Kernel PCA

Laplacian Eigenmaps, ISOMAP, LLE 👉

Autoencoders <





When data lies on or near a low d-dimensional linear subspace, axes of this subspace are an effective representation of the data

Identifying the axes is known as Principal Components Analysis, and can be obtained by Eigen or Singular value decomposition

Data for PCA

Data $X = [x_1, x_2, ..., x_n]$ where each data point x_i is D-dimensional vector

X is D x n matrix

Assume data are centered i.e. sample mean

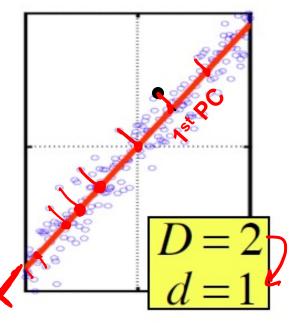
$$\frac{1}{n} \sum_{i=1}^{n} \mathbf{x_i} = 0 \quad \checkmark$$

What if data is not centered?

Subtract off sample mean from each data point

Since data matrix is centered, sample covariance matrix can be written as

$$S = \frac{1}{n} X X^{\top}$$



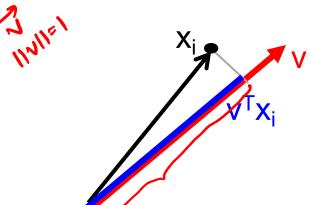
Principal Components (PC) are orthogonal directions that capture most of the variance in the data

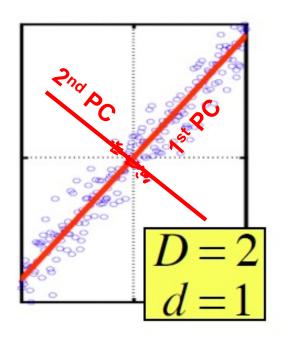
1st PC – direction of greatest variability in data

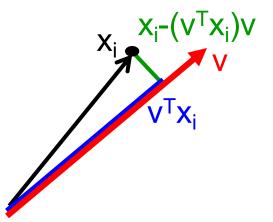
Projection of data points along 1st PC discriminate the data most along any one direction

Take a data point x_i (D-dimensional vector)

Projection of x_i onto the 1st PC v is v^Tx_i







Principal Components (PC) are orthogonal unit norm directions that capture most of the variance in the data

1st PC – direction of greatest variability in data VI (|VIII)=1

2nd PC – Next orthogonal (uncorrelated) direction of greatest variability

(remove all variability in first direction, then find next direction of greatest variability)

And so on ...

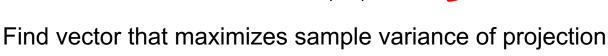
$$D \rightarrow d < D$$

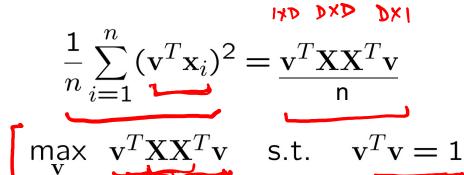
Let v₁, v₂, ..., v_d denote the principal components

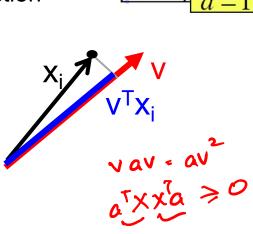
Orthogonal and unit norm $v_i^T v_i = 0$ $\underline{i} \neq \underline{j}$

$$v_i^T v_j = 0 \quad \underline{i} \neq \underline{j}$$

$$v_i^T v_i = 1 = ||V_i|| = ||V_i||$$







Poll:

Is this a convex optimization problem?

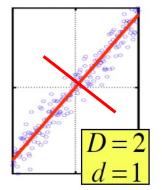
Let v₁, v₂, ..., v_d denote the principal components

Orthogonal and unit norm
$$v_i^T v_j = 0 \quad i \neq j$$

$$v_i^T v_j = 0 \quad i \neq j$$

$$v_i^T v_i = 1$$

Find vector that maximizes sample variance of projection



$$\frac{1}{n} \sum_{i=1}^{n} (\mathbf{v}^T \mathbf{x}_i)^2 = \underline{\mathbf{v}}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

$$\max_{\mathbf{v}} \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

$$\max_{\mathbf{v}} \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

$$\mathbf{v}^T\mathbf{v} = 1$$



Lagrangian: $\max_{\mathbf{v}} \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} - \lambda (\mathbf{v}^T \mathbf{v} - \mathbf{1})$ $2 \mathbf{X} \mathbf{v}^T - 2 \lambda \mathbf{v} = 0$

$$\partial/\partial \mathbf{v} = 0$$

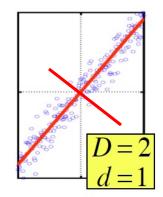
$$\partial/\partial \mathbf{v} = 0 \qquad (\mathbf{X}\mathbf{X}^T - \lambda \mathbf{I})\mathbf{v} = 0$$

$$\Rightarrow (\mathbf{X}\mathbf{X}^T)\mathbf{v} = \lambda\mathbf{v}$$

1) center data

$$(\mathbf{X}\mathbf{X}^T)\mathbf{v} = \lambda\mathbf{v}$$

Therefore, v is the eigenvector of sample covariance matrix XX^T



Sample variance of projection = $\mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda$

Thus, the eigenvalue λ denotes the amount of variability captured along that dimension (aka amount of energy along that dimension).

Eigenvalues $\lambda_1 > \lambda_2 > \lambda_3 \geq \dots$

The 1st Principal component v₁ is the eigenvector of the sample covariance matrix XX^T associated with the largest eigenvalue λ₁

The 2nd Principal component v₂ is the eigenvector of the sample covariance matrix XX^T associated with the second largest eigenvalue λ₂

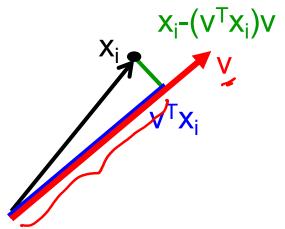
Another interpretation

Maximum Variance Subspace: PCA finds vectors v such that projections on to the vectors capture maximum variance in the data

$$\frac{1}{n} \sum_{i=1}^{n} (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

Minimum Reconstruction Error: PCA finds vectors v such that projection on to the vectors yields minimum MSE reconstruction

$$\frac{1}{n} \sum_{i=1}^{n} \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2$$



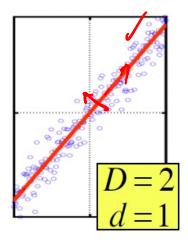
$$\begin{aligned} &\frac{1}{n}\sum_{i=1}^{n}\|\mathbf{x}_{i}-(\mathbf{v}^{T}\mathbf{x}_{i})\mathbf{v}\|^{2} = \frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{x}_{i}^{T}\mathbf{x}_{i}+\mathbf{v}^{T}(\mathbf{v}^{T}\mathbf{x}_{i})^{2}\mathbf{v}\right) \\ &= \frac{1}{n}\sum_{i=1}^{n}\|\mathbf{x}_{i}-(\mathbf{v}^{T}\mathbf{x}_{i})\mathbf{v}\|^{2} \\ &= \frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{x}_{i}^{T}\mathbf{x}_{i}-(\mathbf{v}^{T}\mathbf{x}_{i})\mathbf{v}^{T}\mathbf{x}_{i}\right) \\ &= \frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{x}_{i}^{T}\mathbf{x}_{i}-(\mathbf{v}^{T}\mathbf{x}_{i})^{T}\right) \\ &= \arg\max_{i} -\frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{v}^{T}\mathbf{x}_{i}\right)^{T} \\ &= \arg\max_{i} \frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{v}^{T}\mathbf{x}_{i}\right)^{T} \\ &= \operatorname{arg}\max_{i} \frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{v}^{T}\mathbf{x}_{i}\right)^{T} \end{aligned}$$

Dimensionality Reduction using PCA

The eigenvalue λ denotes the amount of variability captured along that dimension.

Zero eigenvalues indicate no variability along those directions => data lies exactly on a linear subspace

Only keep data projections onto principal components with non-zero eigenvalues, say $v_1, ..., v_d$ where $d = rank(XX^T)$

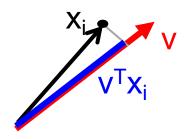


Original Representation data point

$$x_i = [x_i^1, x_i^2, \dots, x_i^D]^T$$

(D-dimensional vector)

Transformed representation projections

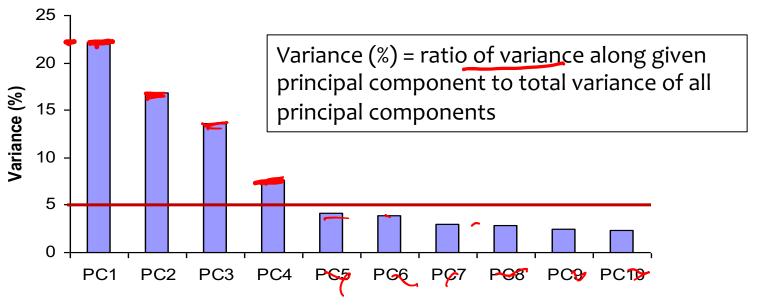


Dimensionality Reduction using PCA

In high-dimensional problem, data usually lies near a linear subspace, as noise introduces small variability

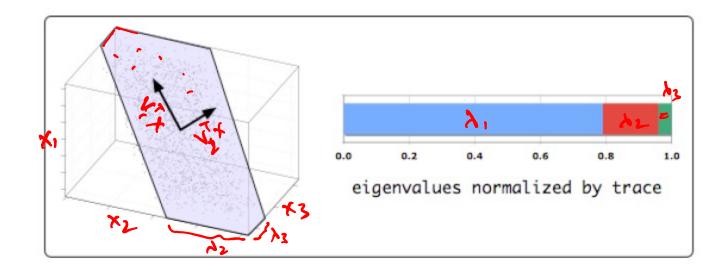
Only keep data projections onto principal components with large eigenvalues

Can *ignore* the components of lesser significance.



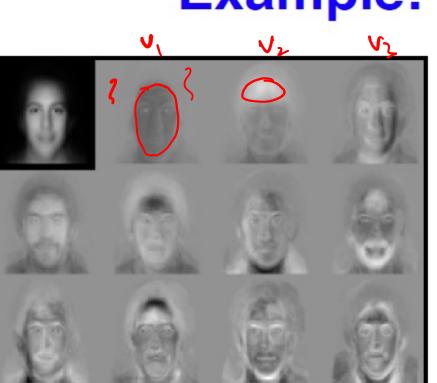
You might lose some information, but if the eigenvalues are small, you don't lose much

Example of PCA



Eigenvectors and eigenvalues of covariance matrix for n=1600 inputs in d=3 dimensions.





Figenfaces from 7562 images:

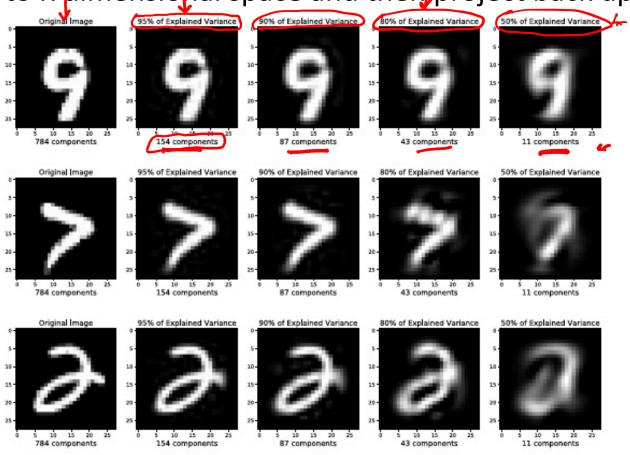
top left image is linear combination of rest.

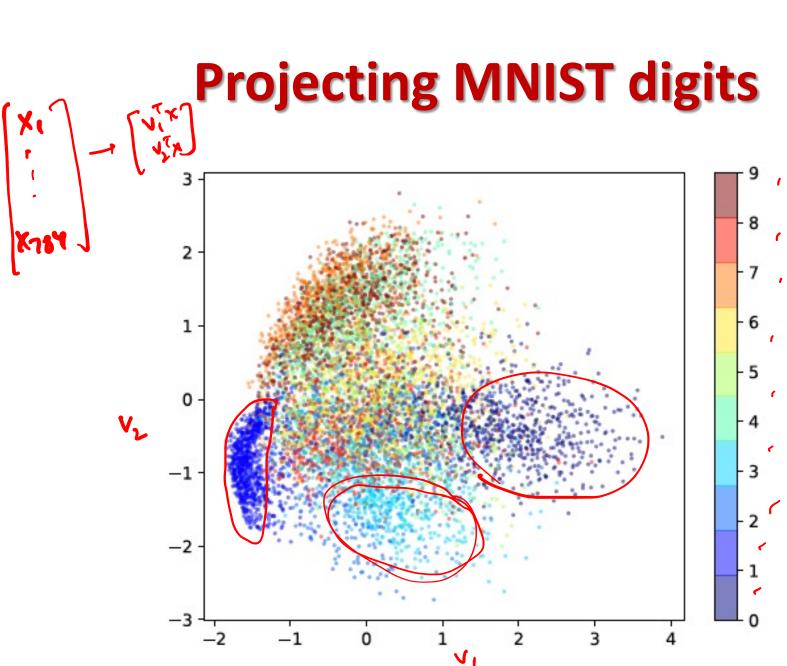
Sirovich & Kirby (1987) Turk & Pentland (1991)

Example: MNIST digits

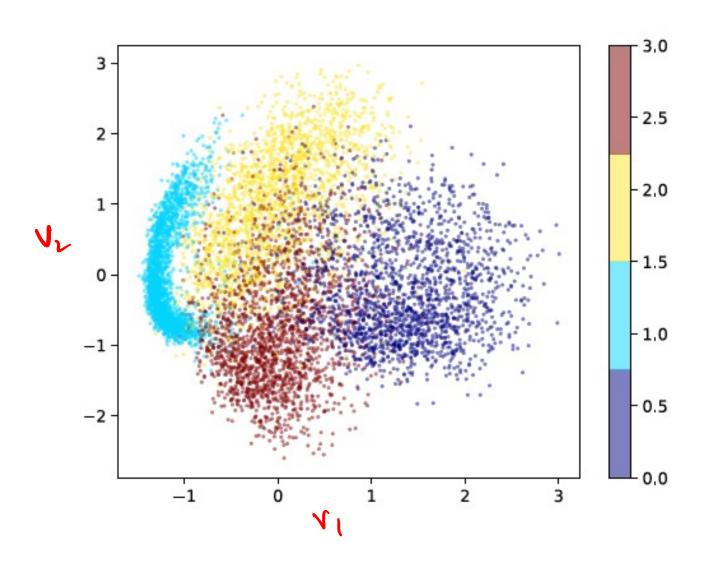
28x28 images = 784 PCA vectors

Project to K dimensional space and then project back up





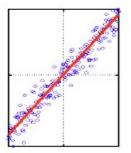
Projecting MNIST digits

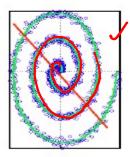


Properties of PCA

Strengths

- Eigenvector method
- No tuning parameters
- Non-iterative *
- –No local optima ⁻





Weaknesses

- -Limited to second order statistics
- Limited to linear projections

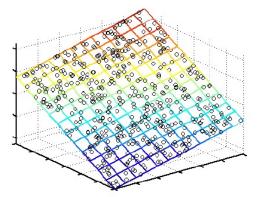
Unsupervised Dimensionality Reduction

Linear

Principal Component Analysis (PCA)

Factor Analysis

Independent Component Analysis (ICA)



Nonlinear

Kernel PCA



Laplacian Eigenmaps, ISOMAP, LLE Autoencoders

