Applied Machine Learning

Dimensionality reduction

Reihaneh Rabbany



COMP 551 (winter 2023) ¹

Admin

- Projects:
 - Project 4 is now optional (Project 3 is due this Friday)
 - Project grading is now: 50% = P1 (10+2%), P2 (15+3%), P3 (15+5%)
 - P4 worth 10% and will be added to the overall grade
 - Groups of 2-3
 - Needs a paragraph on contributions
- Quizzes:
 - Quiz 11 is the last quiz
 - Lowest grade will be dropped

• no need to take it if you already have the full grade for quizzes

Learning objectives

What is dimensionality reduction?What is it good for?Linear dimensionality reduction:

- Principal Component Analysis
- Relation to Singular Value Decomposition

Motivation

Scenario: we are given high dimensional data and asked to make sense of it!

Real-world data is high-dimensional

- Visualization: we can't visualize beyond 3D
- Compression: processing and storage is costly
- Downstrean analysis, e.g. clustering or classification
 - features may not have any semantics (value of the pixel vs happy/sad)
 - many features may not vary much in our dataset (e.g., background pixels in face images)

Dimensionality reduction: faithfully represent the data in low dimensions

- We can often do this with real-world data (manifold hypothesis)
- finding meaningful low-dimensional structures in high-dimensional observations

Dimensionality reduction

Dimensionality reduction: faithfully represent the data in low dimensions

• learn a mapping between (coordinates) at low-dimension and high-dimensional data



some methods give this mapping in both directions and some only in one direction.

Dimensionality reduction

Dimensionality reduction: faithfully represent the data in low dimensions

• learn a mapping between (coordinates at) low-dimension and high-dimensional data



6

Principal Component Analysis (PCA)

PCA is a **linear** dimensionality reduction method



where Q has orthonormal columns $Q^ op Q = I$ it follows that the pseudo-inverse of Q is $Q^\dagger = (Q^ op Q)^{-1}Q^ op = Q^ op$

PCA: optimization objective

PCA is a linear dimensionality reduction method



PCA: optimization objective

PCA is a **linear** dimensionality reduction method faithfulness is measured by the reconstruction error $\min_{Q} \sum_{n} ||x^{(n)} - x^{(n)^{\top}Q}Q^{\top}||_{2}^{2} \quad s.t. \qquad Q^{\top}Q = I$

strategy: find $D \times D$ matrix Q, and only use D' columns Since Q is orthogonal we can think of it as a change of coordinates



PCA: a change of coordinates

strategy: find $D \times D$ matrix Q, and only use D' columns

Since Q is orthonormal we can think of it as a change of coordinates



we want to change coordinates such that coordinates 1,2,...,D' best explain the data for any given D'



example D = 2 (0, 1, 0) q_1 (1, 0, 0)

 q_2

PCA preserves variance

Find a change of coordinate using orthonormal matrix

first new coordinate has maximum **variance** (lowest reconstruction error) second coordinate has the next largest variance

$$P = \begin{bmatrix} Q_{1,1}, \dots, Q_{1,D} \\ \vdots, \ddots, \vdots \\ Q_{D,1}, \dots, Q_{D,D} \end{bmatrix}$$

along which one of these directions the data has a higher variance? more spread out?



. . .

this direction is the vector q_1 projection is given by $\frac{x^{(n)^\top q_1}}{||q_1||_2} = x^{(n)^\top q_1}$ projection of the whole dataset is $Xq_1 = z_1$ $z_1^\top = [z_1^{(1)}, z_1^{(2)}, \dots, z_1^{(N)}]$

PCA preserves variance

Find a change of coordinate using *orthonormal matrix*

first new coordinate has maximum variance

projection of the whole dataset is $z_1 = Xq_1$ $Var(z_1) = \frac{1}{N}\sum_n (z_1^{(n)} - 0)^2$ assuming features have zero mean, maximize the variance of the projection: $\frac{1}{N}z_1^\top z_1$

$$egin{aligned} \max_{q_1} rac{1}{N} z_1^ op z_1^ op z_1 &= \max_{q_1} rac{1}{N} q_1^ op X^ op X q_1 &= \max_{q_1} q_1 \Sigma q_1^ op u_1^ op u_1$$

$$\Sigma = rac{1}{N} X^ op X = rac{1}{N} \sum_n (x^{(n)} - 0) (x^{(n)} - 0)^ op$$

because the mean is zero

 $\Sigma_{i,j}~~$ is the sample covariance of feature i and j

$$\Sigma_{i,j} = \operatorname{Cov}[X_{:,i},X_{:,j}] = rac{1}{N}\sum_n x_i^{(n)} x_j^{(n)}$$

Covariance matrix

variance of a random variable $\operatorname{Var}(x) = \mathbb{E}[(x - \mathbb{E}[x])^2] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$ covariance of two random variable $\operatorname{Cov}(x, y) = \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$ for $x \in \mathbb{R}^D$ we have covariance matrix

$$\mathrm{Cov}(x_1,x_1) = \mathrm{Var}(x_1) \sum_{\mathbf{1},\mathbf{1}} \sum_{\mathbf{1},\mathbf{1}$$

given a dataset $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$ sample covariance matrix

 $\hat{\Sigma}^{MLE} \hat{\Sigma} = \mathbb{E}_{\mathcal{D}}[(x - \mathbb{E}_{\mathcal{D}}[x])(x - \mathbb{E}_{\mathcal{D}}[x])^{ op}]$ $egin{array}{c} ext{the empirical estimate} \ x - ig(rac{1}{N}\sum_{x \in \mathcal{D}}xig) \end{array}$

Correlation and dependence

correlation is normalized covariance

$$\operatorname{Corr}(x_i,x_j) = rac{\operatorname{Cov}(x_i,x_j)}{\sqrt{\operatorname{Var}(x_i)\operatorname{Var}(x_j)}} \ \in [-1,+1]$$

two variables that are independent are uncorrelated as well

$$p(x_i,x_j) = p(x_i)p(x_j)
ightarrow \mathbb{E}[x_ix_j] = \mathbb{E}[x_i]\mathbb{E}[x_j]
ightarrow \mathrm{Cov}(x_ix_j) = 0$$

the inverse is generally not true (zero correlation doesn't mean independence)



in each example above correlation between two coordinates is zero, but they are not independent

Decomposing the covariance matrix

covariance matrix is symmetric positive semi definite

• symmetric

•
$$\Sigma_{d,d'} = \operatorname{Cov}(x_d, x_{d'}) = \operatorname{Cov}(x_{d'}, x_d) = \Sigma_{d',d}$$

- positive semi definite
 - for any $y \in \mathbb{R}^D$ we have $y^\top \Sigma y = (y^\top \mathbb{E}[(x \mathbb{E}[x])(x \mathbb{E}[x])^\top]y) = \operatorname{Var}(y^\top x) \geq 0$

any symmetric positive semi-definite matrix can be decomposed as

$$\Sigma = Q \Lambda Q^{ op}$$
 Spectral Decomposition
 $\begin{vmatrix} I \\ diagonal \ D imes D \end{vmatrix}$
orthogonal $QQ^{ op} = Q^{ op}Q = I$ (rotation and reflection)

PCA with Eigenvalue decomposition

find a change of coordinate using *an orthogonal matrix*

first new coordinate has maximum variance

$$\max_{q_1} q_1 \Sigma q_1^\top \quad s.t. \quad ||q_1|| = 1$$

covariance matrix is **symmetric** and **positive semi-definite** $(X^{\top}X)^{\top} = X^{\top}X$ $a^{\top}\Sigma a = \frac{1}{N}a^{\top}X^{\top}Xa = \frac{1}{N}||Xa||_2^2 \ge 0 \quad \forall a$

any symmetric matrix has the following decomposition

$$\Sigma = Q \Lambda Q^ op$$
 (as we see shortly using Q here is not a co-incidence)

 $QQ^{\top} = Q^{\top}Q = I$ dxd orthogonal matrix each column is an eigenvector

diagonal and sorted ($\lambda_1 > \lambda_2 > \lambda_3 > ...$) corresponding eigenvalues are on the diagonal

positive semi-definiteness means these are non-negative

PCA: Principal Component Analysis

find a change of coordinate using *an orthogonal matrix*

first new coordinate has maximum variance

$$q_1^* = rg\max_{q_1} {q_1}^{ op} \Sigma q_1 \qquad s.t. \quad ||q_1|| = 1$$

 $\max_{q_1} q_1^ op Q \Lambda Q^ op q_1 = \lambda_1$ using eigenvalue decomposition

maximizing direction is the eigenvector with the largest eigenvalue (first column of Q)

 $oldsymbol{q_1} = Q_{:,1}$ first principal direction

second eigenvector gives the $q_2=Q_{:,2}$ second principal direction

so for PCA we need to find the eigenvectors of the covariance matrix

Reducing dimensionality

projection into the principal direction q_i is given by Xq_i



think of the projection XQ as a change of coordinates

we can use the first D' coordinates $Z = XQ_{:,:D'}$ to reduce the dimensionality while capturing a lot of the variance in the data we can project back into original coordinates using $\tilde{X} = ZQ_{::D'}^{\top}$

reconstruction

18

Example: digits dataset

let's only work with digit 2! $x^{(n)} \in \mathbb{R}^{784}$



center the data and form the covariance matrix $\sum_{784 imes 784}$

find the **eigenvectors** of the covariance matrix, the principal directions



use the first 20 directions to reduce dimensionality from 784 to 20!

using 20 numbers we can represent PC coefficient $x^{+}q_{i}$ (the new coordinates) each image with a good accuracy 2022× 1790× -602× -203× 195× -134× -146× -356× 294× -397× -330× input 34 × 52 × -12 × 114× 34 × -50 × -167> rec

19

example: digits dataset

3D embedding of MNIST digits

(https://projector.tensorflow.org/)

$$x^{(n)} \in \mathbb{R}^{784}$$

the embedding 3D coordinates are

 Xq_1, Xq_2, Xq_3



example: text dataset



example: face dataset

eigenfaces for face recognition read more here

 $x^{(n)}$ Lindsay Davenport George W Bush Vin Diesel Surakait Sathirathai Billy Crystal Rubens Barrichello Mary Carey Richard Myers Yasser Arafat Sarah Price Dean Barkley Frank Taylor Sheryl Crow Noah Wyle Colin Powell





Lindsay Davenport

Billy Crystal

Richard Myers

Frank Taylor

Sheryl Crow

Noah Wyle

 $q_1, q_2, \ldots q_{15}$ eigenface 0 eigenface 4 $x^{(n)} \in \mathbb{R}^{64 imes 64}$ $z^{(n)} = x^{(n)} \, {}^{^{+}} Q_{:,:250}$ eigenface 8 Surakait Sathiratha Vin Diese eigenface 12 Rubens Barrichello Mary Carey Yasser Arafat Sarah Price Dean Barkley

eigenface 1 eigenface 2 eigenface 5 eigenface 6 eigenface 9 eigenface 10 eigenface 13 eigenface 14 Figure #6: Bunch of ghost shaped images. Look at them in the eyes.

> mean face used for centring the data



eigenface 3

eigenface 7

eigenface 11

eigenface 15

there is another way to do PCA

without using the covariance matrix

Singular Value Decomposition (SVD)

any N x D real matrix has the following decomposition



Singular value & eigenvalue decomposition

recall that for PCA we used the eigenvalue decomposition of $\ \ \Sigma = rac{1}{N} X^ op X$

how does it relate to SVD?

$$X^ op X = (USV^ op)^ op (USV^ op) = VS^ op U^ op USV^ op = VS^2V^ op$$
 compare to $rac{1}{N}X^ op X = Q\Lambda Q^ op$

 \rightarrow

eigenvectors of Σ are right singular vectors of X $\ Q = V$ for PCA we could use SVD

• this is the standard computation which works directly with data matrix instead of the covariance matrix

Picking the number of PCs

number of PCs in PCA is a hyper-parameter, how should we choose this?

each new principle direction explains some variance in the data

such that we have $a_1 \geq a_2 \geq \ldots \geq a_D$ (by definition of PCA) we can divide by total variance to get a ratio $r_i = rac{a_i}{\sum_d a_d}$

$$a_d = rac{1}{N}\sum_n {z_d^{(n)}}^2$$



first few principal directions explain most of the variance in the data!

Picking the number of PCs

recall that for picking the principal direction we maximized the variance of the PC

$$\begin{array}{ll} \max_{q} \frac{1}{N} q X^{\top} X q^{\top} &= \max_{q} q \Sigma q^{\top} &= \max_{q_{1}} q^{\top} Q \Lambda Q^{\top} q = \lambda_{1} \\ ||q|| = 1 & ||q|| = 1 & ||q|| = 1 \end{array}$$
so the variance ratios are also given by $r_{i} = \frac{\lambda_{i}}{\sum_{d} \lambda_{d}}$
so we can also use eigenvalues to pick the number of PCs



digits **example**:

two estimates of variance ratios do match

Matrix factorization

PCA and SVD perform matrix factorization



this gives a row-rank approximation to our original matrix X

- we can use this to compress the matrix
- we can find give a "smooth" reconstruction of X (remove noise or fill missing values)

example

200

600 29

Matrix factorization





changing the rank D' gives different amount of compression



Matrix factorization

relationship to **K-means**

K-means also can be seen as matrix factorization



matrix product simply equates each row of X with one row of the factor matrix

- instead of principal components cluster centers
 factor loading matrix one nonzero per row of Z (each node belongs to one cluster)

Autoencoders

a feed-forward neural net which predicts its input

- can be trained with reconstruction loss
 - e.g. mean squared error: $\sum_n ||x^{(n)} \hat{x}^{(n)}||_2^2$

dimensionality reduction with a **bottleneck layer** much smaller than input



Autoencoders

a feed-forward neural net which predicts its input

- can be trained with reconstruction loss
 - e.g. reconstruction loss: $||x \psi(\phi(x))||_2^2$

dimensionality reduction with a **bottleneck layer** much smaller than input

- optimal weights for linear autoencoder are the principal components
- **nonlinear** dimensionality reduction if activations are not all linear
 - projecting the data on a non-linear manifold
 - deep autoencoders are very powerful



https://www.cs.toronto.edu/~hinton/science.pdf

Autoencoders: example

MNIST digits





33

To sign attendance go to: http://www.reirab.com/Teaching/AML23/link44.html course website

Summary

Dimensionality reduction helps us:

- visualize our data
- compress it
- simplify the computational need of further analysis (clustering, supervised learning etc.)
- also can be used for anomaly detection (not discussed)

PCA is a linear dimensionality reduction method

- projects the data to a linear space (spanned by D' principal directions)
 - directions are eigenvectors of the covariance matrix
 - the projection has maximum variance (minimum reconstruction error)
 - eigenvalues tell us about the contribution of each new principal direction
- PCA using Singular Value Decomposition
- Model selection for PCA
- PCA as matrix factorization and its relationship to k-means
- practical note: don't forget to subtract the mean!