

PCA

have data set $x_1, \dots, x_n \quad x_i \in \mathbb{R}^d$

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

assume mean 0 $(\sum_{i=1}^n x_i = 0)$

Examples

2

- n images & pixels each
- n measurements & sensors
- n docs & words
- n people & movies
- n customers & products

Fix k

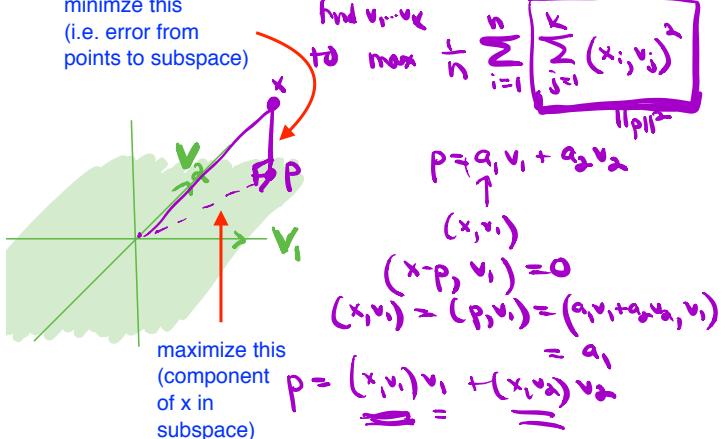
Find k-dimensional subspace (defined by orthonormal vectors $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_k$)

so as to minimize $\frac{1}{n} \sum_{i=1}^n \text{distance}(x_i \rightarrow \text{subspace S spanned by } \vec{v}_1, \vec{v}_2, \dots, \vec{v}_k)$

\equiv maximizing $\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (x_i, v_j)^2$
variance of projected pts

projecting $x_i \rightarrow S$
 $\sum_{j=1}^k (x_i, v_j) \vec{v}_j$

minimize this
(i.e. error from points to subspace)



$$\text{maximizing } \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (x_i \cdot v_j)^2 = \sum_{j=1}^k (x \cdot v_j)^T x \cdot v_j = \frac{1}{n} \sum_{j=1}^k v_j^T X^T X v_j$$

empirical covariance matrix $X^T X$

How to find $\vec{v}_1, \dots, \vec{v}_k$

① Compute eigendecomposition of empirical covariance matrix $X^T X$

$$A := X^T X = Q D Q^T$$

As A is symmetric eigendecomposition is special: Q is orthogonal! orthonormal eigenvectors

$$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
all eigenvalues ≥ 0

$$A q_i = \lambda_i q_i \quad \forall i$$

where

$$Q = \begin{pmatrix} | & | & | \\ q_1 & \dots & q_n \\ | & | & | \end{pmatrix} \quad \text{orthonormal eigenvectors}$$

$$\begin{aligned} q_1, \dots, q_n &\text{ orthonormal} \\ Q Q^T = Q^T Q &= I \\ \forall w \quad \|Qw\| &= \|w\| \end{aligned}$$

What happens when we multiply a vector by this matrix A? Its the projection of v onto the basis given by the eigenvectors

$Q D Q^T v = Q D (q_1 * v, q_2 * v, \dots, q_n * v)^T \rightarrow q_i * v$ is the component of v in the direction of q_i

$$\begin{aligned} \textcircled{2} \quad \text{Set } v_1 &:= q_1 \\ v_2 &:= q_2 \\ &\vdots \\ v_k &:= q_k \end{aligned}$$

New coordinates of x_i

in the basis of eigenvectors

are $(x_i \cdot v_1), (x_i \cdot v_2), \dots, (x_i \cdot v_k)$

low dimensional representation
of \vec{x}_i

1. Find the top component, \mathbf{v}_1 , using power iteration.

2. Project the data matrix orthogonally to \mathbf{v}_1 :

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_m \end{bmatrix} \mapsto \begin{bmatrix} (\mathbf{x}_1 - \langle \mathbf{x}_1, \mathbf{v}_1 \rangle \mathbf{v}_1) \\ (\mathbf{x}_2 - \langle \mathbf{x}_2, \mathbf{v}_1 \rangle \mathbf{v}_1) \\ \vdots \\ (\mathbf{x}_m - \langle \mathbf{x}_m, \mathbf{v}_1 \rangle \mathbf{v}_1) \end{bmatrix}.$$

$$\mathbf{X} - \mathbf{X} \mathbf{v} \mathbf{v}^\top$$

This corresponds to subtracting out the variance of the data that is already explained by the first principal component \mathbf{v}_1 .

3. Recurse by finding the top $k-1$ principal components of the new data matrix.

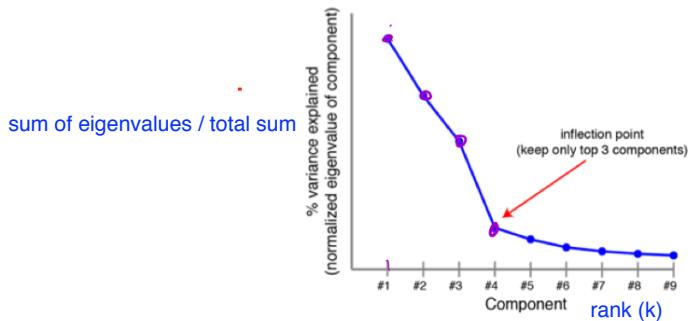
Applications

- ① Visualization ② Compression ③ Learning

How to choose k ?

How do we choose what low rank approximation to use? What rank?

- For visualization: a few
- compression: Look at eigenvalues.
As soon as small enough; happy.



Scree plot. Principal components are ranked by the amount of variance they capture in the original dataset, a scree plot can provide some sense of how many components are needed.

Singular Value Decomposition (SVD)

gives us best way to approximate a matrix with a "low rank" matrix

$$\begin{bmatrix} 7 & ? & ? \\ ? & 8 & ? \\ ? & 12 & 6 \\ ? & ? & 2 \\ 21 & 6 & ? \end{bmatrix}.$$

Motivation:

can we reconstruct missing entries?

construct the missing entries using as small as rank as possible (that matches the given entry)

1 2 3

$$\begin{array}{c|ccc} & 1 & 2 & 3 \\ \hline 1 & 7 & & \\ 2 & & 8 & \\ 3 & & 12 & 6 \\ 4 & & & 2 \\ 5 & 21 & 6 & \end{array}$$

Rank 0

all zeros matrix

Rank 1

rank 1 matrix; every row is a multiple of $v.T$; every column is a multiple of u

$$A = uv^T = \begin{bmatrix} - & u_1v^T & - \\ - & u_2v^T & - \\ \vdots & & \\ - & u_mv^T & - \end{bmatrix} = \begin{bmatrix} | & | & & | \\ v_1u & v_2u & \cdots & v_nu \\ | & | & & | \end{bmatrix}$$

Rank 2

$$A = uv^T + wz^T = \begin{bmatrix} - & u_1v^T + w_1z^T & - \\ - & u_2v^T + w_2z^T & - \\ \vdots & & \\ - & u_mv^T + w_mz^T & - \end{bmatrix} = \begin{bmatrix} | & | \\ u & w \\ | & | \end{bmatrix} \cdot \begin{bmatrix} - & v^T & - \\ - & z^T & - \end{bmatrix}$$

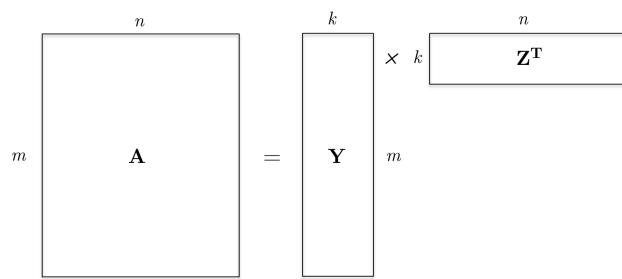


Figure 1: Any matrix \mathbf{A} of rank k can be decomposed into a long and skinny matrix times a short and long one.

Why might a matrix be approximately low rank?

Example: movie ratings

- Suppose each movie characterized by relatively small # of attributes
 - e.g. romance, violence, comedy, ...
- and each person characterized by their preferences on each of these

Singular Value Decomposition (SVD)

$$\begin{matrix} d \\ n \end{matrix} \begin{matrix} A \\ = \end{matrix} \begin{matrix} n \\ n \end{matrix} \begin{matrix} U \\ \cdot \quad \dots \end{matrix} \begin{matrix} d \\ n \end{matrix} \begin{matrix} S \\ \cdot \quad \dots \end{matrix} \begin{matrix} d \\ d \end{matrix} \begin{matrix} V^T \\ \vdots \end{matrix}$$

If the matrix was square and symmetric this would be the eigendecomposition?

Figure 2: The singular value decomposition (SVD). Each singular value in \mathbf{S} has an associated left singular vector in \mathbf{U} , and right singular vector in \mathbf{V}^T .

Running time $O(\min(n^2d), O(d^2n))$

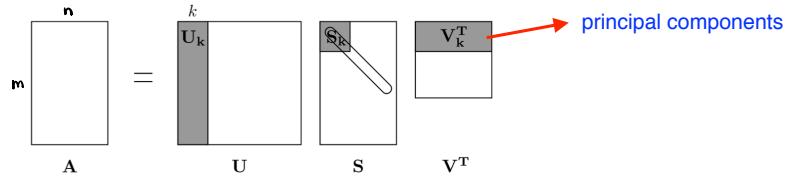


Figure 3: Low rank approximation via SVD. Recall that S is non-zero only on its diagonal, and the diagonal entries of S are sorted from high to low. Our low rank approximation is $A_k = U_k S_k V_k^\top$.

$$\|A - A_k\|_F^2 \leq \|A - B\|_F^2 \quad \text{for any } k \text{ rank matrix } B$$

i.e. in the least squared sense, the SVD is the best rank k approximation possible.

exercise: Knowing that v_1 is the first principal component, prove this theorem for $k = 1$ (i.e. via the definition of PCA)

Relationship between SVD and PCA

$$\begin{matrix} d \\ n \end{matrix} \begin{matrix} A \\ = \end{matrix} \begin{matrix} d \\ \vdots \\ \dots \\ d \end{matrix} \begin{matrix} U \\ \dots \\ \dots \\ V^T \end{matrix}$$

Figure 2: The singular value decomposition (SVD). Each singular value in S has an associated left singular vector in U , and right singular vector in V .

PCA

$$X^T X = Q D Q^T$$

SVD

$$X = U S V^T$$

$$X^T X = V S^T U^T U S V^T = V S^2 V^T$$

Thus $Q = V$ and $D = S^2$ so the eigenvalues are just the squares of the singular values from SVD!

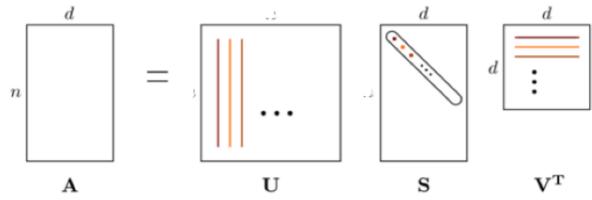


Figure 2: The singular value decomposition (SVD). Each singular value in \mathbf{S} has an associated left singular vector in \mathbf{U} , and right singular vector in \mathbf{V}^T .

Application 1 : Denoising

Suppose a matrix \mathbf{A} is a noisy version of a rank k matrix

$$\mathbf{A} = \mathbf{C} + \mathbf{N} \quad \text{where } \mathbf{C} \text{ is the true rank } k \text{ structure and } \mathbf{N} \text{ is noise}$$

We can “reconstruct” \mathbf{C} by considering a rank k approximation of \mathbf{A}

Collaborative Filtering

recommendations: which movies to see, which products to buy

Model:

each person's ratings for a set of movies

movies

The diagram illustrates the process of matrix completion. On the left, labeled "people" vertically and "movies" horizontally, is a matrix R with dimensions $n \times d$. The matrix contains numerical values (e.g., 2.5, 4, 1) in some entries and question marks in others, representing "ground truth". An arrow points from this matrix to the right, where another matrix $R^{\hat{}}$ is shown. This second matrix has the same structure but with more question marks, indicating missing data. Some specific values are highlighted in blue.

$$R \underset{n \times d}{\text{ground truth}}$$

$$R^{\hat{}}$$

only have some of the data (maybe the person did not leave a review or didn't see the movie). How can we infer the structure of the missing entries?

Assumptions:

R is low rank
e.g. rank k

example:
humor, violence,
romance

The matrix completion problem:

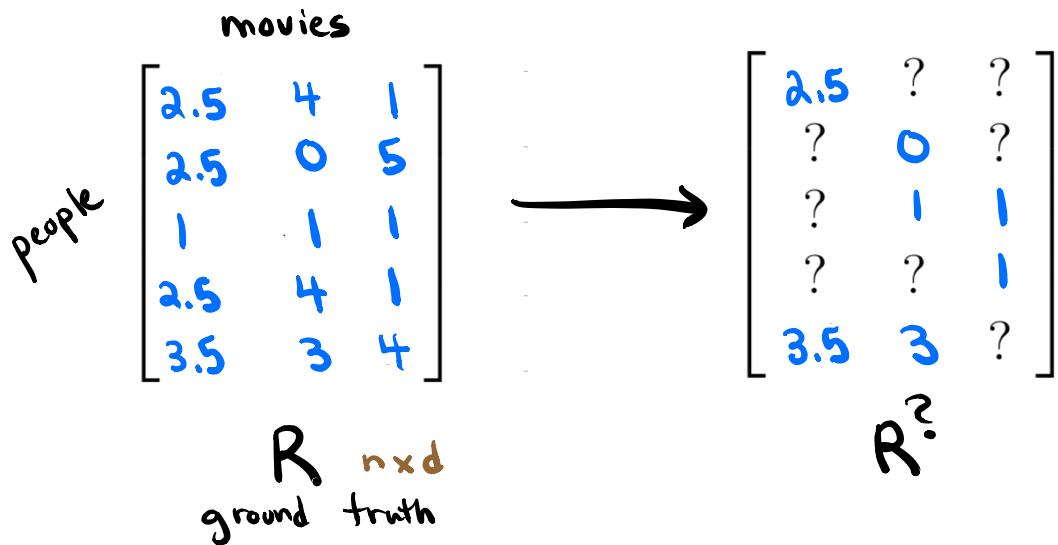
$$R = h \left(\begin{array}{c} k \\ \diagdown \\ - a_i - \end{array} \right) \left(\begin{array}{c} k \\ \diagup \\ b_j \end{array} \right)$$

Find a_i with $i = 1, \dots, n$ and b_j with $j = 1, \dots, d$
such that R captures all the known entries as good as possible

Collaborative Filtering

recommendations: which movies to see, which products to buy

Model:



Assumptions:

- ① R is low rank
e.g. rank k

example:
humor, violence,
romance

The matrix completion problem:

fill in the missing entries.

Theorem \hat{A} nxd matrix of indep. r.v.s, each with variance bounded by σ^2 each entry is a random variable; all we see is some sample

If $A = E(\hat{A})$ is rank k

then with high probability

$$\|A - \hat{A}_k\|_F^2 \text{ is } O(k\sigma^2(n+d))$$

k rank approximation of \hat{A} is quite close to A

average per element error

avg per elt error

$$O\left(\frac{k\sigma^2(n+d)}{nd}\right)$$

$$= o(1)$$

average error per element is constant?

Suppose we have a matrix A that is rank k and that entry (i,j) is deleted with probability $1-p_{ij}$ (i.e. goes to a ? in example below ; retained with probability p_{ij})

All we see are the nondeleted entries according to this random process (R? below) Then I can essentially reconstruct A

Let $\hat{A}_{ij} = \begin{cases} A_{ij} / p_{ij} & \text{if entry is present and 0 otherwise} \\ 0 & \text{otherwise} \end{cases}$

$E[\hat{A}_{ij}] = A_{ij}$ the true value! Thus \hat{A} is an unbiased estimator of our true matrix A. $A = E[\hat{A}]$. By the theorem above then, taking the low-rank approximation (of rank k) of \hat{A} we get a good approximation of A!

$$\begin{bmatrix} 2.5 & ? & ? \\ ? & 0 & ? \\ ? & 1 & 1 \\ ? & ? & 1 \\ 3.5 & 3 & ? \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} & & \\ & & \\ & & \\ & & \\ & & \end{bmatrix}$$

$R^?$

M