

Dimension Reduction

Adam Richards

Galvanize, Inc

Last updated: 2. April 2018

Why do PCA in the first place?

High dimensional data causes many problems. Here are a few:

- 1 The Curse of Dimensionality
- 2 It's hard to visualize anything with more than 3 dimensions.
- 3 Points are “far away” in high dimensions, and it's easy to overfit small datasets.
- 4 Often (especially with image/video data) the most relevant features are not explicitly present in the high dimensional (raw) data.
- 5 Remove Correlation (e.g. neighboring pixels in an image)

Singular Value Decomposition (SVD)

- So we can use a technique called SVD for more efficient computation
- It is not always easy to directly compute eigenvalues and eigenvectors
- SVD is also useful for discovering hidden topics or latent features

Every matrix has a unique decomposition in the following form

$$M = U\Sigma V^T$$

where

- U is column orthogonal: $U^T U = I$
- V is column orthogonal: $V^T V = I$
- Σ is a diagonal matrix of positive values, where the diagonal is ordered in decreasing order

We can reduce the dimensions by sending the smaller of the diagonals to 0.

SVD and PCA

In PCA we had

$$M^T M V = V \Lambda$$

where Λ is the diagonal matrix of eigenvalues

According to SVD we have

$$M = U \Sigma V^T$$

$$\begin{aligned} M^T M &= (U \Sigma V^T)^T U \Sigma V^T \\ &= V \Sigma^T U^T U \Sigma V^T \\ &= V \Sigma^2 V^T \end{aligned}$$

This is the same equation as with PCA we just have $\Lambda = \Sigma^2$

Movie ratings

	Matrix	Alien	Serenity	Casablanca	Amelie
Alice	1	1	1	0	0
Bob	3	3	3	0	0
Cindy	4	4	4	0	0
Dan	5	5	5	0	0
Emily	0	2	0	4	4
Frank	0	0	0	5	5
Greg	0	1	0	2	2

```
import numpy as np
from numpy.linalg import svd

M = np.array([[1, 1, 1, 0, 0],
               [3, 3, 3, 0, 0],
               [4, 4, 4, 0, 0],
               [5, 5, 5, 0, 0],
               [0, 2, 0, 4, 4],
               [0, 0, 0, 5, 5],
               [0, 1, 0, 2, 2]])

u, e, v = svd(M)
print M
print "="
print(np.around(u, 2))
print(np.around(e, 2))
print(np.around(v, 2))
```

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} -0.14 & -0.02 & -0.01 \\ -0.41 & -0.07 & -0.03 \\ -0.55 & -0.09 & -0.04 \\ -0.69 & -0.12 & -0.05 \\ -0.15 & 0.59 & 0.65 \\ -0.07 & 0.73 & -0.68 \\ -0.08 & 0.3 & 0.33 \end{bmatrix} \begin{bmatrix} 12.48 & 0.0 & 0.0 \\ 0.0 & 9.51 & 0.0 \\ 0.0 & 0.0 & 1.35 \end{bmatrix} \begin{bmatrix} -0.56 & -0.59 & -0.56 & -0.09 & -0.09 \\ -0.13 & 0.03 & -0.13 & 0.7 & 0.7 \\ -0.41 & 0.8 & -0.41 & -0.09 & -0.09 \end{bmatrix}$$

With $M = U\Sigma V^T$, U is the user-to-topic matrix and V is the movie-to-topic matrix.

Science Fiction

- First singular value (12.4)
- First column of the U matrix (note: the first four users have large values)
- First row of the V matrix (note: the first three movies have large values)

Romance

- Second singular value (9.5)
- Second column of the U matrix (note: last three users have large values)
- Second row of the V matrix (note: the last two movies have large values)

The third singular value is relatively small, so we can exclude it with little loss of data. Let's try doing that and reconstruct our matrix

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} -0.14 & -0.02 \\ -0.41 & -0.07 \\ -0.55 & -0.09 \\ -0.69 & -0.12 \\ -0.15 & 0.59 \\ -0.07 & 0.73 \\ -0.08 & 0.3 \end{bmatrix} \begin{bmatrix} 12.48 & 0.0 \\ 0.0 & 9.51 \end{bmatrix} \begin{bmatrix} -0.56 & -0.59 & -0.56 & -0.09 & -0.09 \\ -0.13 & 0.03 & -0.13 & 0.7 & 0.7 \end{bmatrix}$$

$$= \begin{bmatrix} 0.99 & 1.01 & 0.99 & -0.0 & -0.0 \\ 2.98 & 3.04 & 2.98 & -0.0 & -0.0 \\ 3.98 & 4.05 & 3.98 & -0.01 & -0.01 \\ 4.97 & 5.06 & 4.97 & -0.01 & -0.01 \\ 0.36 & 1.29 & 0.36 & 4.08 & 4.08 \\ -0.37 & 0.73 & -0.37 & 4.92 & 4.92 \\ 0.18 & 0.65 & 0.18 & 2.04 & 2.04 \end{bmatrix}$$