

Review:

② Mean and variance and covariance:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

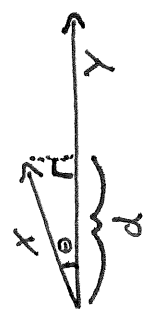
$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \bar{x})^2$$

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \bar{x})(y^{(i)} - \bar{y})$$

④ Dot product, orthogonal projections, ~~orthogonal~~ and basis, orthogonal matrix:

$$X \cdot Y = \sum_{i=1}^p x_i y_i$$

Recall:
 $\cos(\theta) = \frac{x \cdot y}{\|x\| \|y\|}$
 $\cos(\theta) = \frac{d}{\|x\|}$



$$d = \frac{x \cdot y}{\|y\|} = x \cdot \hat{y}$$

Let V be an orthogonal basis. That means the cols of V are all orthogonal to one another and have unit length.

Note:

$$V^T V = I$$

$$\therefore V^T = V^{-1}$$

$$\therefore V V^T = I$$

put an orthonormal basis into a matrix to make an orthonormal matrix

PCA method:

① Our dataset:

n samples $x^{(1)}, x^{(2)}, \dots, x^{(n)}$
 p features $x^{(i)} \in \mathbb{R}^p$

③ Step 1: Create the "centered" design matrix "M".

$$M = \begin{bmatrix} x^{(1)} - \bar{x} \\ x^{(2)} - \bar{x} \\ \vdots \\ x^{(n)} - \bar{x} \end{bmatrix}_{n \times p}$$

⑤ Step 2: Compute the covariance matrix

$$\Sigma = M^T M / n$$

⑦ Step 3: Compute the eigenvectors/values of Σ .

$$\text{evecs, evals} = \text{eig}(M^T M)$$

The event w/ the largest eval is the direction of most variance in the data.
 ... second most...

2d example:



event 1 is the "first" principal component
 all events are orthogonal.

PCA derivation:

→ What will Σ look like?
 ... It is like there is a little covariance between all pairs, and some pairs will have (possibly) a lot of covariance.

Crazy idea: Could we find a new orthogonal basis, project our data to it, and have no covariance?
 Let's see... (no covariance)

$$(MV)^T(MV) = \begin{bmatrix} r_1 & r_2 & \emptyset \\ \emptyset & r_2 & \dots \\ \emptyset & \dots & r_p \end{bmatrix}$$

$$(MV)^T(MV) = V^T M^T M V$$

$$V V^T M^T M V = V \begin{bmatrix} r_1 & r_2 & \emptyset \\ \emptyset & r_2 & \dots \\ \emptyset & \dots & r_p \end{bmatrix}$$

$$M^T M V = V \begin{bmatrix} r_1 & \emptyset \\ \emptyset & r_2 \end{bmatrix}$$

$$\begin{bmatrix} M^T M \end{bmatrix} \begin{bmatrix} r_1 \\ r_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} r_1 & \emptyset \\ \emptyset & r_2 \end{bmatrix} \begin{bmatrix} r_1 \\ r_2 \\ \vdots \end{bmatrix}$$

consider only one column of result on each side

$$\begin{bmatrix} M^T M \end{bmatrix} \begin{bmatrix} r_1 \\ r_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \end{bmatrix} \rightarrow Ax = \lambda x \rightarrow \text{Eigenvectors!}$$

9

The relative size of the eigenvalues tells us the % explained variance of each eigenvector.

Let $\lambda_1, \lambda_2, \dots, \lambda_p$ be the decreasing list of eigenvalues. ^{also called} the "loadings".

$$\text{total power} = \sum_{i=1}^p \lambda_i$$

$$\text{power of } k \text{ eigenvectors} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^p \lambda_i}$$

Typically, set k s.t. 90% of the variance is explained.

10 Keep k eigenvectors \rightarrow let $U_{p \times k}$ = $\begin{bmatrix} | & | & | \\ u_1 & u_2 & \dots & u_k \\ | & | & | \end{bmatrix}$
Edit: (I should have called this U_i)

where u_i are the principal components.
i.e. $\text{eig}(M^T M)$

11 Reduce the dims by $R = MU$ i.e. $R_{n \times k} = M_{n \times p} U_{p \times k}$
 \hookrightarrow reduced to k dims.

12 Reconstruct \hat{M} an estimate of M by "undoing" the dim reduction:

$$R = MU \rightarrow \hat{M} = RU^T = RU^T$$

Singular Value Decomposition (SVD)

① Every matrix M has a unique decomposition in the following form:

$$M = U S V^T$$

where:

$M_{n \times p}$ is any matrix

$U_{n \times n}$ is an orthogonal matrix

$S_{n \times p}$ is a diagonal matrix with positive dec. real values on the diagonal
 these values are the "singular values" of M .

$V^T_{p \times p}$ is an orthogonal matrix

Some Properties:

The columns of U are the eigenvectors of $M M^T$

The columns of V are the eigenvectors of $M^T M$

If the rank of M is k , the SVD can look like this:

$$M_{n \times p} = U_{n \times k} S_{k \times k} V_{k \times p}^T$$

Note: $k \leq \min(n, p)$

② We can already see the SVD is computing the same thing as PCA from our description on the left: both compute $\text{eig}(M^T M)$

Here's another way to see the same thing (that SVD and PCA are doing the same thing):

Recall: (from PCA)

$\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_p \end{bmatrix}$

~~$M^T M = U \Lambda U^T$~~

~~$(M U)^T M U = \Lambda$~~

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

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

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

New: (in SVD world)

$M = U S V^T$

$M^T M = (U S V^T)^T U S V^T$

$= V S^T U^T U S V^T$

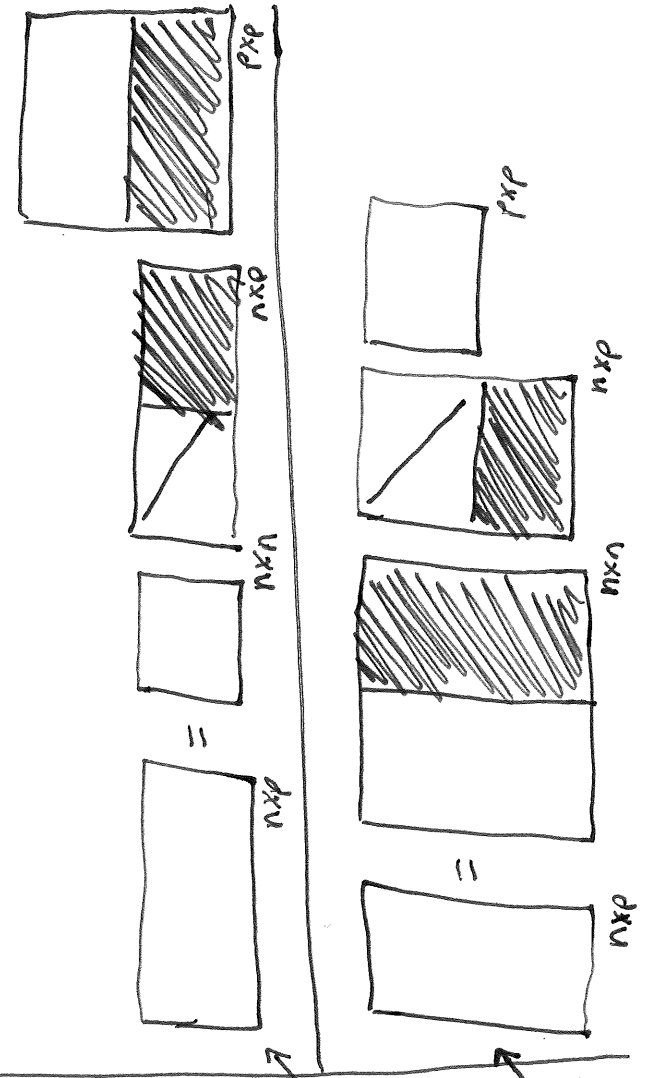
$= V S^T S V^T$

$M^T M = V S^2 V^T$

$\Lambda = S^2$

In SVD world, V is the matrix where the cols are $\text{eig}(M^T M)$

\rightarrow the evals are the squared singular vals.



So... SVD and PCA do the same thing... but there are two ways SVD wins.

SVD FTW #1:

③ When we calculate eigenfaces...

M is $N \times P$ where

$$N=105$$

$$P = 320 \times 243 = 77,760$$

$M^T M$ is $P \times P$...

... so $M^T M$ contains

6,046,617,600

real values (floats in python)

... a float is 8 bytes, so we need ~ 48 GB of RAM to store $M^T M$.

We can shortcut this with SVD.

SVD allows us to ~~calculate~~ ^{compute}

eig($M^T M$) without actually

~~calculating~~ ^{computing} $M^T M$.

Using "compact SVD": $M_{n \times p} = U_{n \times r} S_{r \times r} V_{p \times r}^T$

where $r \leq N$ ($r = \text{rank of } M$)

SVD FTW #2:

④ SVD can reveal ~~hidden~~ "latent features" in your data. This can be used for "topic modeling" in datasets of user ratings.

... Switch to IPython notebook

because we don't want

to be writing a ton of

real values on the board...

... we'll do a real example

of this (w/ real numbers)

in the notebook...

Takeaway: SVD is more computationally efficient than PCA when $N \ll P$.

PCA analogy to shadows cast by my body!

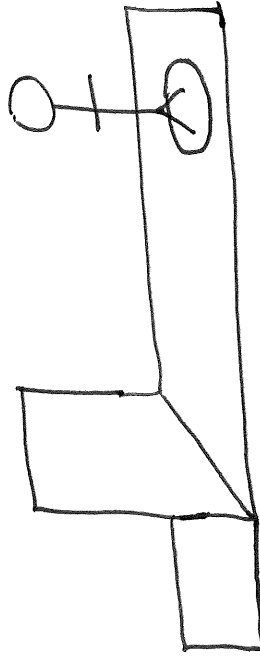
Projection from

3d to 2d via

the sun.



↳ Not in
the direction
of my body's most
Variance.



In the direction
of my body's most
variance.

