

# SDS 385: Stat Models for Big Data Lecture 12: PCA and random projections

Purnamrita Sarkar Department of Statistics and Data Science The University of Texas at Austin

https://psarkar.github.io/teaching

- Lots of high dimensional data:
  - Documents Can have about a million words.
  - Recommender systems If a datapoint is an user, there are around tens of thousands of movies
  - Images each image is represented using many pixel values more for higher resolution

- Interpretation/visualization is difficult
- Storage and computation is difficult
- Many features are useless, and may lead to bad generalization error

- Interpretation/visualization is difficult
- Storage and computation is difficult
- Many features are useless, and may lead to bad generalization error
- Solution:
  - Do feature selection
  - Represent data as a linear combination of *important* features

- Goal: Find the direction of the most variance.
- Say X is the data matrix
- The average is  $\bar{\mathbf{x}} = \frac{\sum_{i=1}^{n} \mathbf{x}_{i}}{n}$

• Let 
$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$$

- Goal: Find the direction of the most variance.
- Say X is the data matrix
- The average is  $\bar{\mathbf{x}} = \frac{\sum_{i=1}^{n} \mathbf{x}_{i}}{n}$
- Let  $\tilde{\mathbf{x}}_i = \mathbf{x}_i \bar{\mathbf{x}}$
- The sample variance of  $(\tilde{x}_1, \dots, \tilde{x}_n)$  along a direction w is give by:

$$\frac{1}{n}\sum_{i=1}^n (\tilde{\boldsymbol{x}}_i^T\boldsymbol{w})^2$$

• What is the sample variance of  $(x_1, \ldots, x_n)$  along a direction w?

# **Principal Component Analysis**



• So the first PC direction is:

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \frac{1}{n} \sum_{i=1}^n (\tilde{\boldsymbol{x}}_i^T \boldsymbol{w})^2$$

• And the first PC component of  $\tilde{x}_i$  is  $\tilde{x}_i^T w_1$ 

• So the  $k^{th}$  PC direction is:

$$\boldsymbol{w}_{k} = \arg \max_{\substack{\|\boldsymbol{w}\|=1\\ \boldsymbol{w} \perp \boldsymbol{w}_{1}, \dots, \boldsymbol{w}_{k-1}}} \frac{1}{n} \sum_{i=1}^{n} (\tilde{\boldsymbol{x}}_{i}^{T} \boldsymbol{w})^{2}$$

- And the  $k^{th}$  PC component of  $\tilde{\mathbf{x}}_i$  is  $\tilde{\mathbf{x}}_i^T \mathbf{w}_k$
- Note that  $w_1, \ldots, w_k$  form an orthogonal basis.

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a low dimensional representation of  $\tilde{X}$

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$w_1 = \arg \max_{\|w\|=1} w^T \tilde{X}^T \tilde{X} w$$

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a low dimensional representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \boldsymbol{w}^T \tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}} \boldsymbol{w}$$

• This is the first eigenvector of  $S = \tilde{X}^T \tilde{X}$ 

- Any square symmetrix matrix S has real eigenvalues
- The *i*<sup>th</sup> eigenvalue, vector pair satisfy  $S w_i = \lambda_i w_i$
- The eigenvectors are orthogonal to each other, and normalized to have length 1.

- Any square symmetrix matrix S has real eigenvalues
- The *i*<sup>th</sup> eigenvalue, vector pair satisfy  $Sw_i = \lambda_i w_i$
- The eigenvectors are orthogonal to each other, and normalized to have length 1.
- In matrix terms, we can write:

$$S = U \Sigma U^T$$
, where

- columns of U are the organal eigenvectors, and
- $\boldsymbol{\Sigma}$  is a diagonal matrix with eigenvalues on the diagonal
- The larger the magnitude of the eigenvalue, more important the eigenvector

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \boldsymbol{w}^T \tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}} \boldsymbol{w}$$

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \boldsymbol{w}^T \tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}} \boldsymbol{w}$$

• This is the first eigenvector of  $S = \tilde{X}^T \tilde{X}$ 

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \boldsymbol{w}^T \tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}} \boldsymbol{w}$$

- This is the first eigenvector of  $S = \tilde{X}^T \tilde{X}$
- What is S?

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \boldsymbol{w}^T \tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}} \boldsymbol{w}$$

- This is the first eigenvector of  $S = \tilde{X}^T \tilde{X}$
- What is S?
- Its the scalar multiple of the sample covariance matrix

$$\hat{\Sigma} = \frac{1}{n} \sum_{i} \tilde{\boldsymbol{x}}_{i} \tilde{\boldsymbol{x}}_{i}^{T} = \frac{S}{n}$$

- Let W is a matrix with  $w_k$  along its columns
- $\tilde{X}W$  gives a **low dimensional** representation of  $\tilde{X}$
- We can frame the optimization problem also as

$$\boldsymbol{w}_1 = \arg \max_{\|\boldsymbol{w}\|=1} \boldsymbol{w}^T \tilde{X}^T \tilde{X} \boldsymbol{w}$$

- This is the first eigenvector of  $S = \tilde{X}^T \tilde{X}$
- What is S?
- Its the scalar multiple of the sample covariance matrix

$$\hat{\Sigma} = \frac{1}{n} \sum_{i} \tilde{\boldsymbol{x}}_{i} \tilde{\boldsymbol{x}}_{i}^{T} = \frac{S}{n}$$

• So, all you have to do is to calculate eigenvectors of the covariance matrix.

- So, all you have to do is to calculate eigenvectors of the covariance matrix.
- But, do I even need to do that?

- So, all you have to do is to calculate eigenvectors of the covariance matrix.
- But, do I even need to do that?
- The right singular vectors of  $\tilde{X}$  is just fine.

- So, all you have to do is to calculate eigenvectors of the covariance matrix.
- But, do I even need to do that?
- The right singular vectors of  $\tilde{X}$  is just fine.
- How many PC's? (more of a dissertaiton question)

# Singular value decomposition



- The columns of U are orthogonal eigenvectos of  $AA^T$
- The columns of V are orthogonal eigenvectos of  $A^T A$
- $A^T A$  and  $AA^T$  have the same eigenvalues

# Second interpretation

$$\mathbf{x}_i$$
  $\mathbf{v}$   
 $(\mathbf{v}^T \mathbf{x}_i) \mathbf{v}$ 

• Minimum reconstruction error:

$$(\mathbf{x}_i - (\mathbf{x}_i^T \mathbf{w})\mathbf{w})^T (\mathbf{x}_i - (\mathbf{x}_i^T \mathbf{w})\mathbf{w}) = \mathbf{x}_i^T \mathbf{x}_i - (\mathbf{x}_i^T \mathbf{w})^2$$

• So, the first PC direction gives the direction projecting on which has the **minimum reconstruction error**.

• Take the centered data matrix  $\tilde{X}$  with SVD

 $\tilde{X} = USV^T$ 

• Project on the top k PC's  $W \in \mathbb{R}^{p \times k}$ 

• Take the centered data matrix  $\tilde{X}$  with SVD

$$\tilde{X} = USV^T$$

- Project on the top  $k \text{ PC's } W \in \mathbb{R}^{p \times k}$
- You get  $\tilde{X}W = US_kV^T$ , where  $S_k$  has zeroed out all singular values  $\leq \sigma_k$

• Take the centered data matrix  $\tilde{X}$  with SVD

$$\tilde{X} = USV^{T}$$

- Project on the top  $k \text{ PC's } W \in \mathbb{R}^{p \times k}$
- You get  $\tilde{X}W = US_kV^T$ , where  $S_k$  has zeroed out all singular values  $\leq \sigma_k$
- So  $W = \arg \min_{\substack{\operatorname{rank}(B)=k, B \in \mathbb{R}^{n \times p}} \|\tilde{X} B\|_F^2$  and the reconstruction error is  $\sum_{i=k+1}^p \sigma_i^2$

• Take the centered data matrix  $\tilde{X}$  with SVD

$$\tilde{X} = USV^{T}$$

- Project on the top  $k \text{ PC's } W \in \mathbb{R}^{p \times k}$
- You get  $\tilde{X}W = US_kV^T$ , where  $S_k$  has zeroed out all singular values  $\leq \sigma_k$
- So  $W = \arg \min_{\substack{\operatorname{rank}(B)=k, B \in \mathbb{R}^{n \times p}} \|\tilde{X} B\|_F^2$  and the reconstruction error is  $\sum_{i=k+1}^p \sigma_i^2$
- This explains why you want to take large k to reduce approx. error.

• We will make a covariance matrix and generate independent multivariate gaussian random variables

```
: d=1000
Sigma=np.zeros([d,d])
Sigma[0:200,0:200]=.3;
np.fill_diagonal(Sigma,1)
```

• We will make a covariance matrix and generate independent multivariate gaussian random variables

```
: d=1000
Sigma=np.zeros([d,d])
Sigma[0:200,0:200]=.3;
np.fill_diagonal(Sigma,1)
```

```
: plt.matshow(Sigma)
```

: <matplotlib.image.AxesImage at 0x188c373d0>



• Now lets compute eigenvectors of the covariance matrix.

```
X=np.random.multivariate_normal(np.zeros(d),Sigma,5000)
```

```
S=np.cov(np.transpose(X))
u,s,vt=svd(S)
```

• Now lets compute eigenvectors of the covariance matrix.

```
X=np.random.multivariate_normal(np.zeros(d),Sigma,5000)
```

```
S=np.cov(np.transpose(X))
u,s,vt=svd(S)
```

plot(u[:,0])

[<matplotlib.lines.Line2D at 0x167e5ae50>]



• Now lets do SVD on the data matrix X.

u,s,vt=svd(X)

```
plot(vt[0,:])
```

[<matplotlib.lines.Line2D at 0x1680810a0>]



- Erikki Oja wrote a seminal paper in 1982 about a simple neural network model.
- He was inspired by the Hebbian principle (1949, "The organization of behavior", Donald Hebb) which claims that the synaptic energy increases from presynaptic cells stimulating post-synaptic cells.



• For each data-point, you do:

$$w_{t+1} \leftarrow w_t + \eta_t (x_t^T w_t) x_t$$
$$w_{t+1} \leftarrow w_{t+1} / || w_{t+1} ||$$

• For each data-point, you do:

$$w_{t+1} \leftarrow w_t + \eta_t (x_t^T w_t) x_t$$
$$w_{t+1} \leftarrow w_{t+1} / || w_{t+1} ||$$

- Note that here, you do not need to construct the covariance matrix explicitly, which is extremely useful, when *d* is much larger than *n*, i.e. in high dimensional settings.
- Step size  $\eta_t$  can be set as  $c \log n/n$  or  $\eta_t \propto 1/t$
- Sharp error bounds show that the final solution converges to the principal component and the error has weak dependence on dimensionality *d*

- Now lets do Oja's algorithm for X.
- Set  $\eta = 0.001 \log n/n$



- Now lets do Oja's algorithm for X.
- Set  $\eta = 0.01 \log n/n$



20

- Now lets do Oja's algorithm for X.
- Set  $\eta = 0.1 \log n/n$



- Now lets do Oja's algorithm for X.
- Set  $\eta = 1 \log n/n$



Dot product with truth



- Oja's algorithm maximizes  $w^T \operatorname{cov}(X) w$  over ||w|| = 1
- Why is it nonconvex?
  - First, the constraint is non-convex.
  - But you can change the optimization to do max w<sup>T</sup> cov(X)w and cov(X) is PSD.
  - Sure, but even then, you are maximizing a convex function here, not minimizing it.

# **Random projections**

• What if we did something crazy and projected the data on a random vector.

$$\tilde{x}_i = x_i^T R$$

## **Random projections**

• What if we did something crazy and projected the data on a random vector.

$$\tilde{x}_i = x_i^T R$$

• *R* could be Gaussian random variables:

$$R_i \sim N(0,1)$$

• *R* could be centered bernoullis:

$$R_i \sim egin{cases} 1 & ext{with prob } 1/2 \ -1 & ext{with prob } 1/2 \end{cases}$$

• R could be sparse

$${\cal R}_i \sim egin{cases} 1/\sqrt{s} & ext{with prob } 1/2s \ 0 & ext{with prob } 1-1/s \ -1/\sqrt{s} & ext{with prob } 1/2s \end{cases}$$

24

- Why will this work?
- Lets take a unit vector u and see if RP preserves the norm.
- Take d dimensional standard normal vector R

- Why will this work?
- Lets take a unit vector u and see if RP preserves the norm.
- Take d dimensional standard normal vector R

• 
$$Y = u^T R = \sum_i u_i R_i \sim N(0,1)$$

- Why will this work?
- Lets take a unit vector u and see if RP preserves the norm.
- Take d dimensional standard normal vector R

• 
$$Y = u^T R = \sum_i u_i R_i \sim N(0,1)$$

• 
$$Y^2 \sim \chi^2$$
,  $EY^2 = 1$ 

- Why will this work?
- Lets take a unit vector u and see if RP preserves the norm.
- Take d dimensional standard normal vector R

• 
$$Y = u^T R = \sum_i u_i R_i \sim N(0,1)$$

• 
$$Y^2 \sim \chi^2$$
,  $EY^2 = 1$ 

• On an average the length is preserved

- Why will this work?
- Lets take a unit vector u and see if RP preserves the norm.
- Take d dimensional standard normal vector R

• 
$$Y = u^T R = \sum_i u_i R_i \sim N(0,1)$$

• 
$$Y^2 \sim \chi^2$$
,  $EY^2 = 1$ 

- On an average the length is preserved
- But how about the variance?

• Now take  $R \in N(0,1)^{d \times m} / \sqrt{m}$ 

- Now take  $R \in N(0,1)^{d \times m} / \sqrt{m}$
- $u^T R = [u^T R(:, 1), \dots u^T R(:, m)]/\sqrt{m}$

- Now take  $R \in N(0,1)^{d \times m} / \sqrt{m}$
- $u^T R = [u^T R(:, 1), \dots u^T R(:, m)]/\sqrt{m}$

$$||u^{T}R||^{2} = \sum_{i} (u^{T}R(:,i))^{2}/m$$

Average of m chi-squared RVs

concentrates around 1

- Take  $u = \frac{X_i X_j}{\|X_i X_j\|}$ , its norm will be preserved under a random projection, with high prob.
- We are saying, with high probability,

$$(1-\epsilon) \|X_i - X_j\|_2^2 \le \|X_i^T R - X_j^T R\|_2^2 \le (1+\epsilon) \|X_i - X_j\|_2^2$$

• If you pick  $m = \log n/\epsilon^2$ , then this will be satisfied for all pairs, with high probability.

• Draw 1000 dimensional Gaussians, 300 from origin and *S* as cov, 300 from all ones mean vector.



Pairwise distances



Pairwise distances with k = 20

• Draw 1000 dimensional Gaussians, 300 from origin and *S* as cov, 300 from all ones mean vector.



Pairwise distances with k = 70

Histogram of error

• Some pictures are borrowed from Brett Bernstein's notes from NYU.