

Lecture 10

SVD in Least Squares

Principal Components Regression

- 1 Reduce dimension of each sample from p to $k < p$.

Let $X \in \mathbb{R}^{n \times p}$ = training features and $U \Sigma V^T = \text{svd}(X)$. Let $V_k = 1^{\text{st}} k$ columns of V .

Let $\underline{z}_i = V_k^T \underline{x}_i = (i^{\text{th}} \text{ col of } \Sigma_k U_k^T)$

- 2 Train ML model on \underline{z}_i 's:

$$\hat{\underline{b}} = \underset{\underline{b}}{\text{argmin}} \|\underline{y} - \underline{Z} \underline{b}\|_2^2 = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y} \quad \text{where } \underline{Z} = \begin{bmatrix} -\underline{z}_1^T- \\ -\underline{z}_2^T- \\ \vdots \\ -\underline{z}_n^T- \end{bmatrix} \in \mathbb{R}^{n \times k}$$

- 3 Predict for new sample $\underline{x}_{\text{new}}$:

$$\underline{z}_{\text{new}} = V_k^T \underline{x}_{\text{new}} \\ \hat{y}_{\text{new}} = \underline{z}_{\text{new}}^T \hat{\underline{b}}$$

usually k is selected so k features are linearly independent (i.e. all k σ_i 's in Σ_k are > 0)

- 4 (Optional) Find equivalent \underline{w} so that $\underline{x} \underline{w} \approx \underline{Z} \underline{b}$

$$\underline{z}_i^T = i^{\text{th}} \text{ row of } U_k \Sigma_k \iff \underline{Z} = U_k \Sigma_k \Rightarrow \underline{Z} \underline{b} = U_k \Sigma_k \underline{b} = U_k \Sigma_k V_k^T V_k \underline{b} = \underline{X}_k V_k \underline{b}$$

$$\Rightarrow \underline{Z} \hat{\underline{b}} = \underline{X}_k \hat{\underline{w}} \quad \text{where } \hat{\underline{w}} = V_k \hat{\underline{b}}$$

$$\begin{aligned} \text{also, } \hat{\underline{b}} &= (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y} \\ &= (\Sigma_k U_k^T U_k \Sigma_k)^{-1} \Sigma_k U_k^T \underline{y} \\ &= (\Sigma_k \Sigma_k^T)^{-1} \Sigma_k U_k^T \underline{y} \\ &= \Sigma_k^+ U_k^T \underline{y} \end{aligned}$$

$$\Rightarrow \hat{\underline{w}} = V_k \Sigma_k^+ U_k^T \underline{y} \iff \text{PCR same as projecting each } \underline{x}_i \text{ onto best } k\text{-dim subspace, then taking pseudoinverse.}$$

Principal Components Crime

Broke

1. Run PCA on all n training samples,
mapping $\underbrace{\underline{x}_i}_{\in \mathbb{R}^p} = V_k^T \underbrace{\underline{z}_i}_{\in \mathbb{R}^k}$, $i=1, \dots, n$

Specifically, V_k is 1st k cols of V matrix from SVD of $X \in \mathbb{R}^{n \times p}$

2. Split data into train set and test set.
3. Train ML model on (\underline{z}_i, y_i) , $i=1, \dots, n_{\text{train}}$
4. Measure accuracy on $\underbrace{(\underline{z}_i, y_i), i=n_{\text{train}}+1, \dots, n}_{\text{test set.}}$

Problem: V_k depends on test data.

Learned model depends on test data.

Estimated accuracy artificially large

↑
This really happens, especially in
large organizations with different
teams doing data prep vs. training

Woke

1. Split data into train and test sets
2. Run PCA on n_{train} training samples,
mapping $\underline{x}_i \rightarrow \underline{z}_i = V_k^T \underline{x}_i$ for $i=1, \dots, n_{\text{train}}$

Now V_k is 1st k cols of V matrix from SVD
of 1st n_{train} rows of X only

3. Train ML model on (\underline{z}_i, y_i) , $i=1, \dots, n_{\text{train}}$
4. Measure accuracy on $\underbrace{(\underline{z}_i, y_i), i=n_{\text{train}}+1, \dots, n}_{\text{test set.}}$

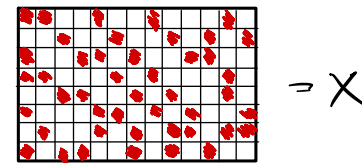
where $\underline{z}_i = V_k^T \underline{x}_i$

⇒ much more reliable predictor of accuracy

Matrix Completion

$X \in \mathbb{R}^{n \times p}$ (eg. n movies, p customers) assume X is low rank

only observe subset of entries of X , want to fill in remainder



$\mathcal{R} = \{(i,j) \text{ s.t. } X_{ij} \text{ is observed}\}$

"Ideal":

$$\hat{X} = \arg \min_M \text{rank}(M) \quad \text{s.t.} \quad M_{ij} = X_{ij} \quad \forall (i,j) \in \mathcal{R}$$

- find lowest rank matrix that fits observed entries

intractable

$\text{rank}(X) = \#\{i : \sigma_i > 0\}$
 $\|X\|_* = \sum_i \sigma_i$
= trace norm
= nuclear norm

Tractable alternative:

$$\hat{X} = \arg \min_M \|M\|_* \quad \text{s.t.} \quad M_{ij} = X_{ij} \quad \forall (i,j) \in \mathcal{R}$$

or, if data is noisy

$$\hat{X} = \arg \min_M \|X_{\mathcal{R}} - M_{\mathcal{R}}\|_F^2 + \lambda \|M\|_*$$

Algorithm:

Iterative Singular Value Thresholding

initialize: $\hat{X} = \text{zeros}(n, p)$

$\hat{X}_{\mathcal{R}} = X_{\mathcal{R}} \leftarrow$ fill in obs. entries

set threshold

for $k = 1, 2, \dots$

$\hat{X}_{\text{old}} = \hat{X}$

$[U, S, V] = \text{svd}(\hat{X})$

$\hat{S} = S \cdot (S \geq \text{threshold})$

$\hat{X} = U \cdot \hat{S} \cdot V^T$

$\hat{X}_{\mathcal{R}} = X_{\mathcal{R}}$

if $\|\hat{X} - \hat{X}_{\text{old}}\|_F < \epsilon$, stop

end

Eigendecomposition and Page Rank

Given a matrix A , we say a non zero vector v is an eigenvector of A if

$$Av = \lambda v \quad \text{where } \lambda \text{ is a scalar}$$

Let $V = [v_1, v_2, \dots, v_n]$ be a matrix whose columns are all eigenvectors of A .

$$AV = A[v_1, v_2, \dots, v_n] = [Av_1, Av_2, \dots, Av_n] = [\lambda_1 v_1, \lambda_2 v_2, \dots, \lambda_n v_n]$$

$$\Rightarrow AV = V\Lambda$$

$$[v_1 \dots v_n] \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_n \end{bmatrix}$$

If the eigenvectors are Linearly independent
then V is invertible

$$\Rightarrow AVV^{-1} = V\Lambda V^{-1} \Rightarrow A = V\Lambda V^{-1}$$

Take a matrix X with SVD $X = U\Sigma V^T$

$$\text{Then } A = \underline{X^T X} = (U\Sigma V^T)^T (U\Sigma V^T) = V \underbrace{\Sigma U^T U \Sigma}_{I} V^T = V \Sigma^2 V^T$$

Recall that V is orthogonal $\Rightarrow V^T = V^{-1}$

$$\Rightarrow A = X^T X = V \underbrace{\Sigma^2}_{\Lambda} V^{-1}$$

$$\Rightarrow A = V \Lambda V^T = \text{eigendecomposition of } A$$

\Rightarrow eigenvectors of A are the right singular vectors of X
eigenvalues of A are the squared singular values of X

$$\left\{ \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \right\}$$

for general A , if A is real and symmetric, then we can ensure the eigenvectors are real and orthonormal $\Rightarrow \underline{V^{-1} = V^T}$

If we can write A as $X^T X$ for some X , then A is real + symmetric
 \downarrow real

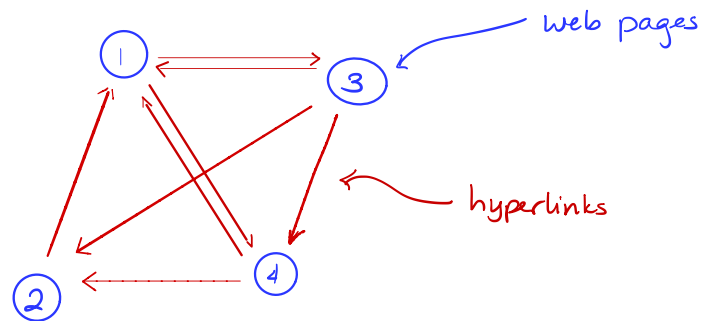
Eigenvalues + Eigenvectors

A vector \underline{v}_k (with $\|\underline{v}_k\|=1$) is an eigenvector of A if there is a scalar λ_k such that
 $A\underline{v}_k = \lambda_k \underline{v}_k$

Let $V = [\underline{v}_1 \ \underline{v}_2 \ \dots \ \underline{v}_n]$ be the n eigenvectors of A . Then

$$\begin{aligned} AV &= A[\underline{v}_1 \ \underline{v}_2 \ \dots \ \underline{v}_n] = [\lambda_1 \underline{v}_1 \ \lambda_2 \underline{v}_2 \ \dots \ \lambda_n \underline{v}_n] \\ &= \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} V = \Lambda V \end{aligned}$$

PageRank



Imagine surfing web by randomly clicking on links at each page at which you arrive.

You will visit more "important" web pages more frequently.

If you do this long enough, you'll reach a steady state where π_i is the probability you're at page i at any given time.

Let M be adjacency matrix of links and A its column normalized version (Markov matrix)

$$M = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 & \frac{1}{3} & \frac{1}{2} \\ 0 & 0 & \frac{1}{3} & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{3} & 0 \end{bmatrix}$$

Think about

A_{ij} = Prob (visit site i next given we're at site j now)

$$\text{SVD: } X = U \Sigma V^T \in \mathbb{R}^{p \times n}$$

$$\begin{aligned} \text{Let } A &= X^T X = V \Sigma U^T U \Sigma V^T \\ &= V \Sigma^2 V^T =: \Lambda V^T \end{aligned}$$

\Rightarrow eigenvalues are real and eigenvectors are orthonormal

$$\Leftrightarrow M_{ij} = \begin{cases} 1 & \text{if page } j \text{ links to page } i \\ 0 & \text{otherwise} \end{cases}$$

We want to find a vector of probabilities $[\pi_1, \pi_2, \pi_3, \pi_4]^T$ so that $\underline{\pi} = A \underline{\pi}$

$\Rightarrow \underline{\pi}$ is simply the leading eigenvector of A !

Markov matrices are special: we know $\lambda_1 = 1$ and $\lambda_1 > \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_n$

To see this, let $\pi = v$,
(1st col of V). Then

$$\pi = A \pi$$

$$\Rightarrow v_i^T \pi = v_i^T A \pi = v_i^T V \Lambda V^T \pi$$

$$\Rightarrow v_i^T v_i = 1 = e_i^T \Lambda e_i = \lambda_i$$

We can find this vector using Power Iterations:

Let $\underline{\pi}^{(0)}$ be an initial guess of π .

for $k = 1, 2, \dots$

$$\text{compute } \underline{\pi}^{(k)} = \frac{A \underline{\pi}^{(k-1)}}{\|A \underline{\pi}^{(k-1)}\|_2}$$

using the 2-norm in the denominator is standard, it ensures the final $\pi^{(k)}$ has norm $\|\pi^{(k)}\|_2 = 1$

for page rank, we might prefer $\|\pi^{(k)}\|_1 = \sum_i \pi_i^{(k)} = 1$, we can manage this 3 ways

$$1. \quad \underline{\pi}^{(k)} = A \underline{\pi}^{(k-1)} / \|A \underline{\pi}^{(k-1)}\|_1$$

2. use 2-norm in denominator, at end, take $\frac{\pi^{(k)}}{\|\pi^{(k)}\|_1}$ as final val.

3. start with $\|\pi^{(0)}\|_1 = 1$ and A column normalized, then don't need denominator at all

Why does this work? We show this under the assumption that the eigenvectors are orthonormal.
(Otherwise proof is more complex)

We can write $\underline{\pi}^{(0)} = c_1 \underline{v}_1 + c_2 \underline{v}_2 + \dots + c_n \underline{v}_n$ for some c_1, \dots, c_n , with $c_i \neq 0$, where $\underline{v}_i = i^{\text{th}}$ col of V

now

$$\underline{\pi}^{(k)} \propto A \underline{\pi}^{(k-1)} \propto A^2 \underline{\pi}^{(k-2)} \propto \dots \propto A^k \underline{\pi}^{(0)}$$

$$= (V \Lambda V^T)^k \underline{\pi}^{(0)} = V \Lambda^k V^T \underline{\pi}^{(0)} = V \Lambda^k V^T (c_1 \underline{v}_1 + \dots + c_n \underline{v}_n) = V \Lambda^k (c_1 \underline{e}_1 + \dots + c_n \underline{e}_n)$$

$$= c_1 \lambda_1^k \underline{v}_1 + \dots + c_n \lambda_n^k \underline{v}_n = c_1 \lambda_1^k \left(\underline{v}_1 + \frac{c_2 \lambda_2^k}{c_1 \lambda_1^k} \underline{v}_2 + \dots + \frac{c_n \lambda_n^k}{c_1 \lambda_1^k} \underline{v}_n \right)$$

$$\text{as } k \rightarrow \infty, \frac{\lambda_i^k}{\lambda_1^k} \rightarrow 0 \text{ for } i \neq 1$$

$$\text{so } A^k \underline{\pi}^{(0)} \rightarrow c_1 \lambda_1^k \underline{v}_1$$

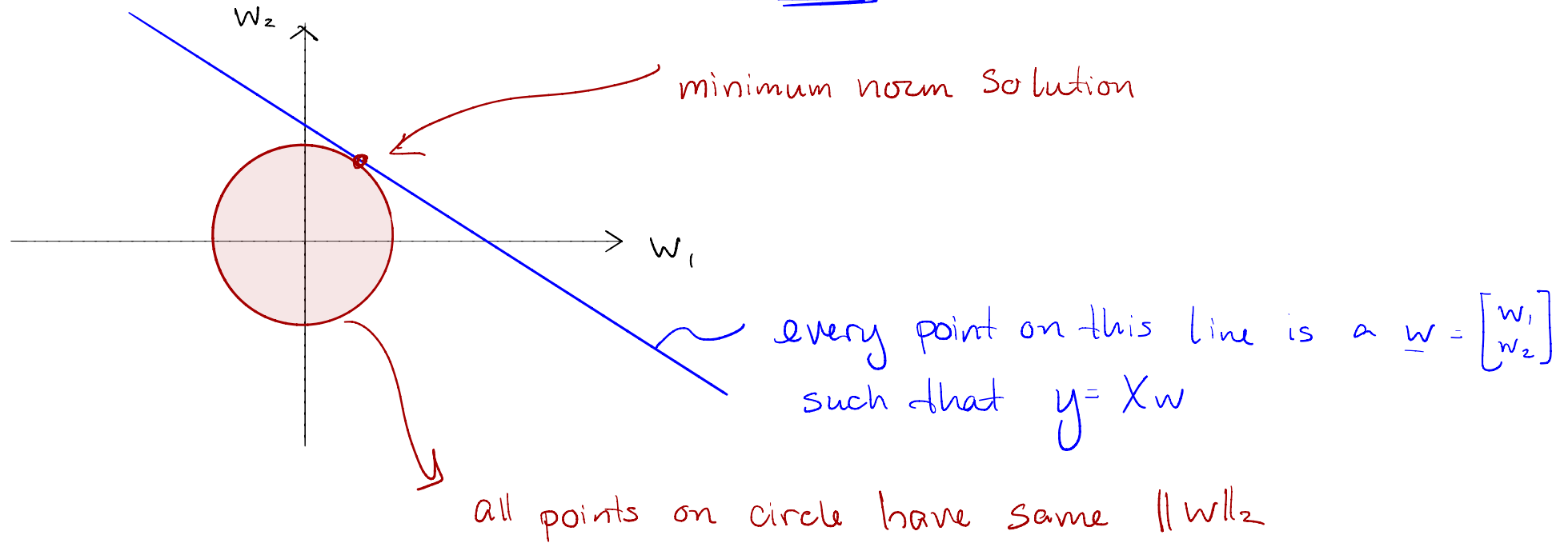
$$\underline{\pi}^{(k)} = \frac{A^k \underline{\pi}^{(0)}}{\|A^k \underline{\pi}^{(0)}\|_2} \rightarrow \frac{c_1 \lambda_1^k \underline{v}_1}{\|c_1 \lambda_1^k \underline{v}_1\|} = \underline{v}_1 \quad (\text{or } -\underline{v}_1)$$

Extra Notes
on Background
(not for in-class lectures)

We want to find w so that $y = Xw$. Imagine columns of X are linearly dependent $\Leftrightarrow X$ has some singular values $= 0$. $\Rightarrow \infty$ many w so that $y = Xw$.

\Rightarrow Minimum norm solution: $\min_w \|w\|_2^2$ such that $y = Xw$

(1)



Connection to Ridge Regression:

If data had a little noise, we might choose $\min_w \|w\|_2^2$ such that $\|y - Xw\| < \text{small threshold}$

(2)

Using optimization theory and Lagrange multipliers, we can show that

$$\min_w \|w\|_2^2 + \frac{1}{\lambda} \|y - Xw\|^2$$

(3)

(3) has the same solution as (2) (for the right choice of λ)

(3) is Ridge regression