

FoDA

L18

•
•

Dimensionality

Reduction

SVD

Dimensionality Reduction

Input Data

$$A \in \mathbb{R}^{n \times d}$$

$$A \subset \mathbb{R}^d$$

$$A = \{a_1, a_2, \dots, a_n\}$$
$$a_i \in \mathbb{R}^d$$

- too high for vis
- runtime depend on d dimensionality d
- statistical : signal vs. noise
 k -dim $(d-k)$ -dim

Data as matrix

$$A \in \mathbb{R}^{n \times d}$$

$d=10,000$ →

$$k = \begin{matrix} 2, 3 \\ 10 \\ 100 \end{matrix}$$

→ SVD

→ map to each $f(a_i) \rightarrow b_i \in \mathbb{R}^k$

Unsupervised Learning

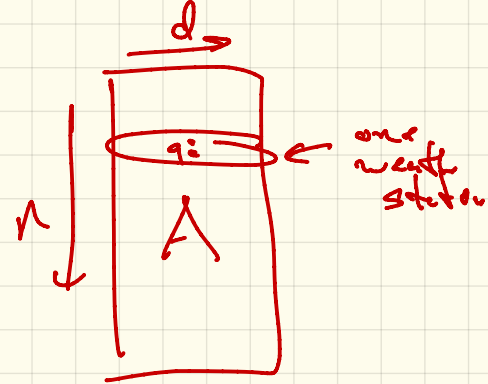
Input A not (x, y)
no labels

↳ Goal: Simpler Representation
of Data.

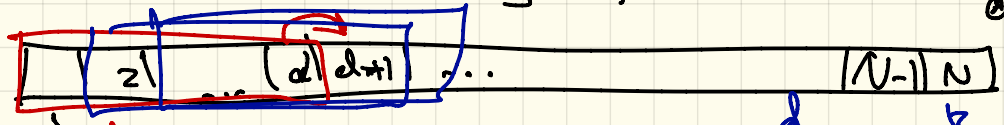
optimize over $A \rightarrow B$

Data Matrix

$$A \in \mathbb{R}^{n \times d}$$

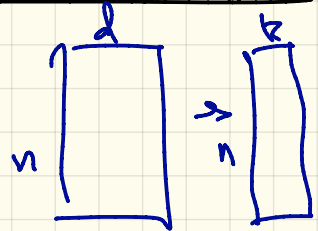


- n weather stations measure temp at d time points
- n users who rode d movies
- stock prices: closing price for N days



$$d=25$$

$$N \rightarrow n = N - d + 1$$
$$d = 25$$



All coordinates have the
same units!

$\|a_i - a_j\| = \text{Euclidean Dist}$

$$= \sqrt{\sum_{l=1}^d (a_{il} - a_{jl})^2}$$

Rth coord.

adding coordinates

Projections

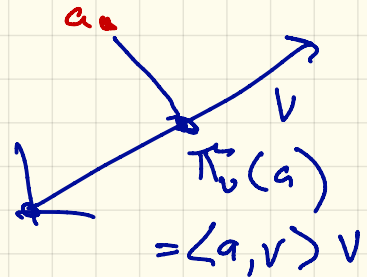
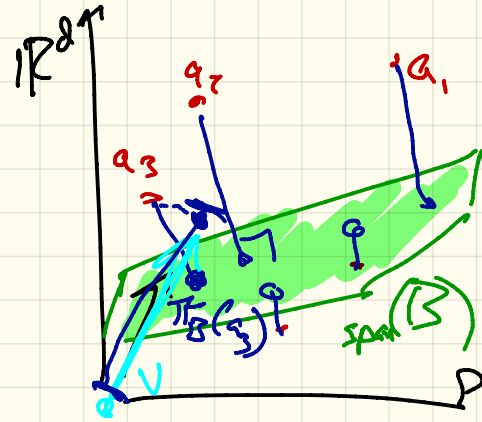
Projecting $a_i \in \mathbb{R}^d$
to $\text{span}(B)$

$$\Pi_B(a_i)$$

B is 1-dimensional $v \in \mathbb{R}^d$
 $\|v\|=1$

$$\Pi_v(a_i) = \langle v, a_i \rangle v$$

Goal: $A \rightarrow B \subset \mathbb{R}^d$
 $\text{rank}(B)=k$



k -dimensional subspace B

orthogonal basis $V_B = \{v_1, v_2, \dots, v_k\}$

• $\|v_i\| = 1$

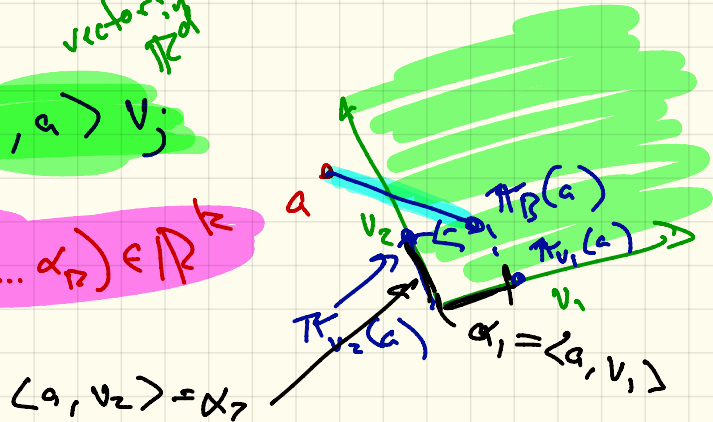
• pair v_i, v_j : $\langle v_i, v_j \rangle = 0$



• For any $x \in B$: $x = \sum_{j=1}^k \alpha_j v_j$ $\alpha_j = \langle x, v_j \rangle$

$\pi_B(a) = \sum_{j=1}^k \langle v_j, a \rangle v_j$

also $\downarrow (\alpha_1, \alpha_2, \dots, \alpha_k) \in \mathbb{R}^k$



Sum of Squared Errors

$$SSE(A, B) = \sum_{a_i \in A} \|a_i - \pi_B(a_i)\|^2$$

Euclidean dist

Goal: best k -dimensional subspace B
(represented by $V_B = \{v_1, v_2, \dots, v_k\}$)

to minimize

$$B^* = \arg \min_B (SSE(A, B))$$

SVD \Rightarrow solve for B s.t. B contains
Origin $(0, 0, \dots, 0)$

PCA \Rightarrow has no origin restriction.
Principal Component Analysis.

Singular Value Decomposition (SVD)

Fortran code

LAPACK

Input

$$A \in \mathbb{R}^{n \times d}$$

(ex. $d < n$)

output

3 matrices

$$U \in \mathbb{R}^{n \times n}$$

$$S \in \mathbb{R}^{n \times d}$$

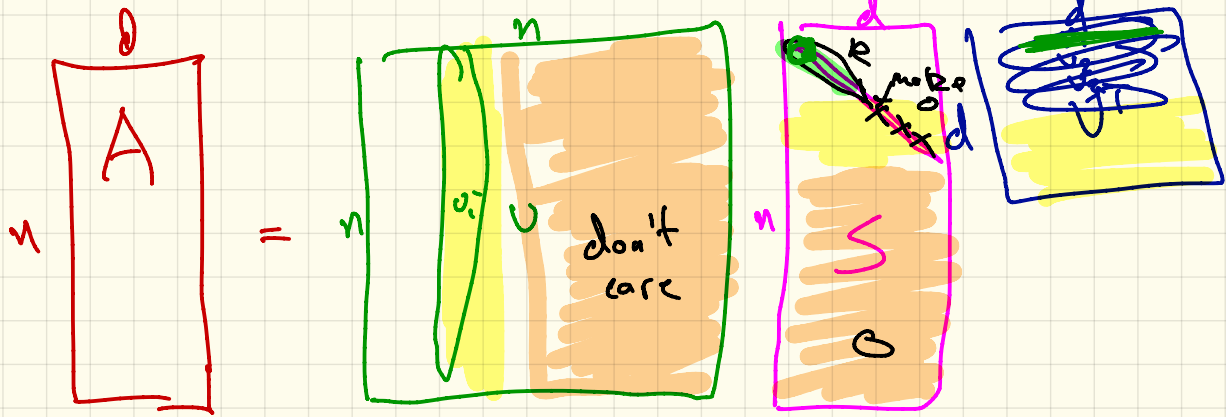
$$V \in \mathbb{R}^{d \times d}$$

$$A = USV^T$$

orthogonal

diagonal $S = \text{diag}(\sigma_1, \sigma_2, \dots)$





$$\|v_i\| = 1$$

$$\langle v_j, v_i \rangle = 0$$

$$v_1 > v_2 > \dots > v_k$$

$$v_1 \quad v_2$$

Best subspace $V_B = \{v_1, v_2, \dots, v_k\}$