

# APS1070

Foundations of Data Analytics and  
Machine Learning

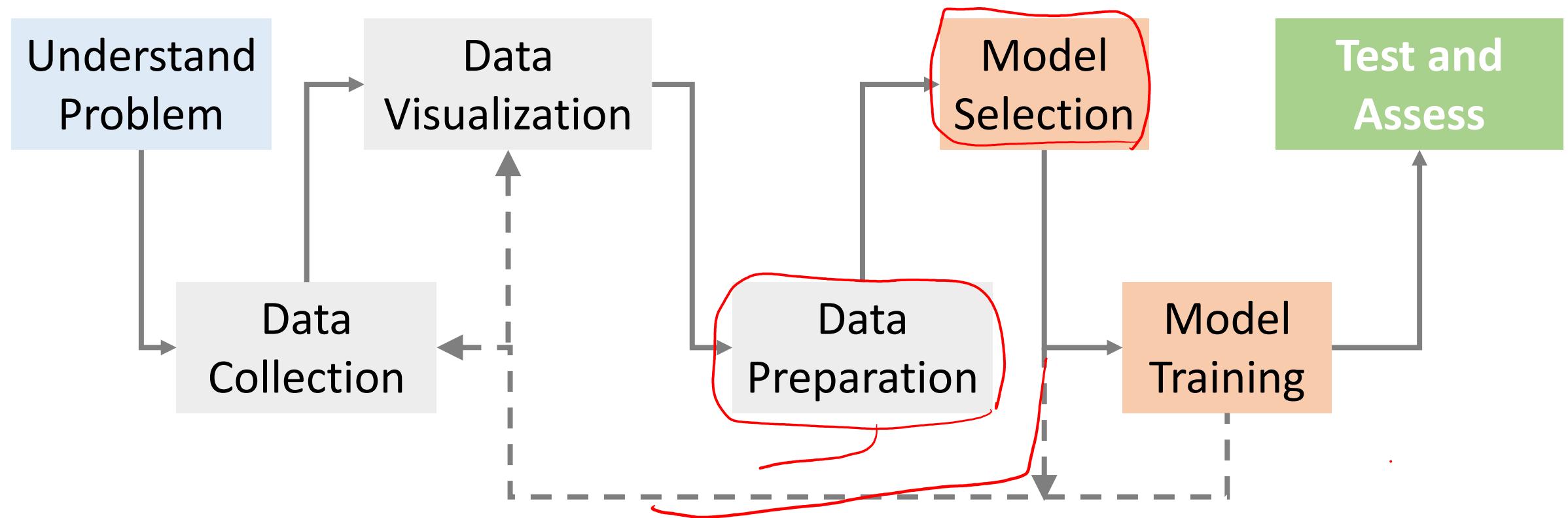
Winter 2022

## Week 8:

- *SVD*
- *Applications*
- *Vector Calculus*



# Course Theme



# Slide Attribution

These slides contain materials from various sources. Special thanks to the following authors:

- Marc Deisenroth
- Jason Riordon

# Last Time

- Looked into PCA for dimensionality reduction
  - Determinant
  - Trace
  - Eigendecomposition
  - PCA Applications
- Today we will introduce SVD for dimensionality reduction and interpretation of data.

# Rectangular Matrices

- Eigendecomposition (and PCA) is limited to square matrices
- Q: How then do we achieve matrix decomposition for rectangular matrices?

# Agenda

- Recap from last time —
- SVD
  - SVD vs Eigendecomposition
  - Dimensionality Reduction
  - Interpretations
- Applications
- Vector Calculus
  - Matrix Differentiation



Theme:  
**Dimensionality Reduction  
and Interpretations**

# Matrix Decompositions Continued

Readings:

- MML Chapter 4.5-8

# Recap: Principal Component Analysis

Data Matrix (rectangular)

	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1

- Taking the square (and symmetrical) covariance matrix we obtain:

$$A = PDP^T \quad \text{Spectral Theorem}$$

- To obtain principal components and associated scores:

Covariance Matrix (square)

	feature 2	
feature 1	val1	val2
feature 2	val2	val4

$$\begin{array}{c} \text{Covariance Matrix (square)} \\ \text{---} \\ \boxed{\begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix}} = \begin{matrix} & & \\ | & | & | \\ v_1 & v_2 & v_3 \\ | & | & | \end{matrix} X \begin{matrix} & & \\ a_{11} & 0 & 0 \\ 0 & a_{22} & 0 \\ 0 & 0 & a_{33} \end{matrix} X \begin{matrix} & & \\ v_1 & v_2 & v_3 \\ | & | & | \end{matrix}^T \end{array}$$

↓ ↓

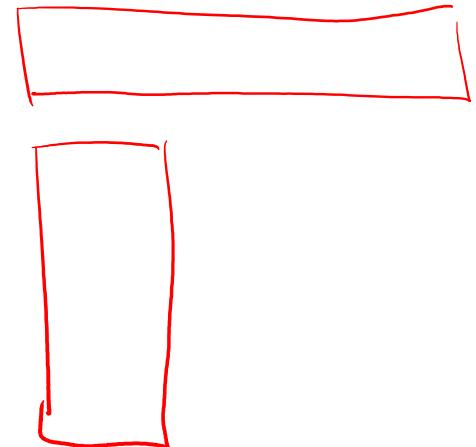
eigenvector matrix      eigenvalue matrix

# Recap: Matrix Decomposition

- We would like a general approach for **decomposing rectangular matrices**.
- PCA is already applicable to rectangular matrices...
- Mathematically PCA is:

$$\rightarrow A = \boxed{XX^T} = PDP^T$$

when the mean/expectation across each  
variable/feature is zero



# Singular Value Decomposition

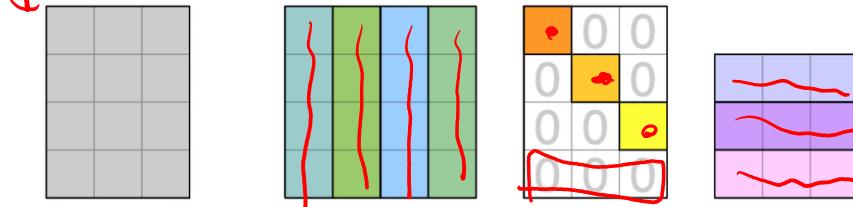
- We would like to obtain the following decomposition for a **rectangular matrix**:

$$A = \underbrace{U}_{m \times n} \underbrace{\Sigma}_{m \times m} \underbrace{V^T}_{n \times n}$$

The diagram illustrates the dimensions of the matrices in the SVD decomposition. Matrix  $A$  is  $m \times n$ . Matrix  $U$  is  $m \times m$ . Matrix  $\Sigma$  is  $m \times n$ . Matrix  $V^T$  is  $n \times z$ .

where  $U$  contains the left-singular vectors,  $V$  has the right-singular vectors and  $\Sigma$  are the singular values.

# Singular Value Decomposition



$$\underset{n \times m}{A} = \underset{n \times n}{U} \underset{n \times m}{\Sigma} \underset{m \times m}{V^*}$$

- Singular values matrix  $\Sigma$  has additional zero padding of rows or columns:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \vdots & & \vdots \\ 0 & 0 & \sigma_n & 0 & \dots & 0 \end{bmatrix} \quad \text{when } n < m$$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_m \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix} \quad \text{when } m < n$$

# Singular Value Decomposition

$$A = \begin{matrix} n \times m \\ \mathbf{U} & \Sigma & \mathbf{V}^T \end{matrix}$$

$A =$

Star Wars	5	4	1	Sci-fi
Blade Runner	5	5	0	
Amelie	0	0	5	French
Delicatessen	1	0	4	

$= \begin{bmatrix} \text{Sci-fi} & \text{French} & & \\ -0.6710 & 0.0236 & 0.4647 & -0.5774 \\ -0.7197 & 0.2054 & -0.4759 & 0.4619 \\ -0.0939 & -0.7705 & -0.5268 & -0.3464 \\ -0.1515 & -0.6030 & 0.5293 & -0.5774 \end{bmatrix}$

$\begin{bmatrix} 9.6438 & 0 & 0 \\ 0 & 6.3639 & 0 \\ 0 & 0 & 0.7056 \\ 0 & 0 & 0 \end{bmatrix}$

$\begin{bmatrix} \text{Sci-fi} & & & \\ -0.7367 & -0.6515 & -0.1811 & \\ 0.0852 & 0.1762 & -0.9807 & \\ 0.6708 & -0.7379 & -0.0743 & \end{bmatrix}$

# Implications

- If we can decompose our matrix, then that can give us insights into our data.
- SVD is behind some of the many machine learning applications
  - recommender systems
  - word embeddings
  - image compressions
  - background removal

# SVD Algorithm

- Finding the left and right singular vectors shares some similarities with PCA:

- Right-singular vectors:

- Left-singular vectors:

$$\begin{aligned} A^T A &= V D V^T \\ A A^T &= U D U^T \end{aligned}$$

$A = U \Sigma V^T$

The diagram illustrates the Singular Value Decomposition (SVD) of a matrix  $A$ . On the left, two equations are shown:  $A^T A = V D V^T$  and  $A A^T = U D U^T$ . The term  $V D V^T$  is annotated with a red bracket labeled  $P D [P^T]$ , where  $P$  is the right singular vector matrix and  $D$  is the diagonal matrix of singular values. The term  $U D U^T$  is annotated with a red bracket labeled  $(P) D P^T$ , where  $U$  is the left singular vector matrix. On the right, the full SVD decomposition is shown as  $A = U \Sigma V^T$ , where  $U$  is  $m \times m$ ,  $\Sigma$  is  $m \times n$ , and  $V^T$  is  $n \times n$ .

- $AA^T$  and  $A^T A$  are symmetrical (which makes eigendecomposition easier)

# SVD Algorithm

➤ Right-Singular:

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

➤ Left-Singular:

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

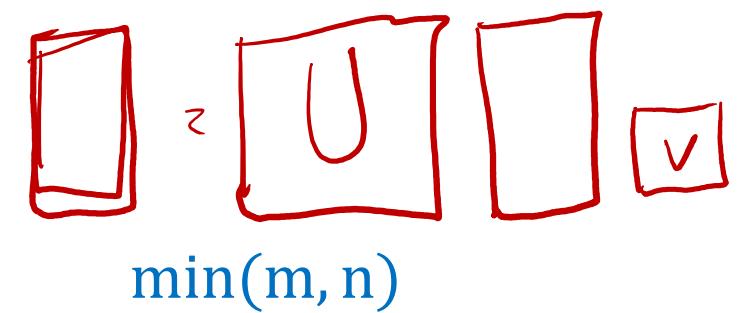
eigendecomposition  
of  $\textcolor{red}{A^T A}$ , with  $\underline{\Sigma} = \underline{\sqrt{D}}$

eigendecomposition  
of  $\textcolor{red}{A A^T}$  with  $\underline{\Sigma} = \underline{\sqrt{D}}$

# SVD Algorithm

- The right-singular vectors ( $V$ ) and left-singular vectors ( $U$ ) (which we know to be orthonormal) are connected through the singular value matrix:

$$AV = U\Sigma$$



$$Av_i = \sigma_i u_i \quad i = 1, \dots, r$$

$$\boxed{\frac{1}{\sigma_i} Av_i = u_i}$$

# Example

Compute the singular value decomposition for matrix A.

$$\underline{A} = \begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix}$$

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

$$\underline{A} = \underbrace{U}_{2 \times 3} \sum_{2 \times 2} \underbrace{\sqrt{\Sigma}}_{2 \times 3}^T$$

$$\boxed{\sqrt{\Sigma} = ?} \rightarrow \det(A^T A - \lambda I) = 0$$

$$\boxed{A^T A} = \begin{bmatrix} 1 & -2 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \underbrace{\begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix}}_A = \underbrace{\begin{bmatrix} 5 & -2 & 1 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}}_{V^T} = B = \boxed{P D P^T}$$

$$\begin{bmatrix} 5-\lambda & -2 & 1 \\ -2 & 1-\lambda & 0 \\ 1 & 0 & 1-\lambda \end{bmatrix} \rightarrow \det \begin{bmatrix} -\lambda^3 + 7\lambda^2 - 6\lambda = 0 \end{bmatrix}$$

$P D P^{-1}$

$$\lambda(-\lambda^2 + 7\lambda - 6) = 0$$

$$\begin{cases} \lambda_1 = 0 \\ \lambda_2 = 1 \\ \lambda_3 = 6 \end{cases}$$

$$\lambda = 6 \rightarrow \begin{bmatrix} -1 & -2 & 1 \\ -2 & -5 & 0 \\ 1 & 0 & -5 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} = 0 \rightarrow \text{span} \left( \begin{bmatrix} 1 \\ -2 \\ 5 \end{bmatrix} \right) \xrightarrow{\text{Normalize}} \sqrt{1 + \frac{4}{25} + \frac{1}{25}} = \frac{\sqrt{30}}{5}$$

$$\lambda = 1 \rightarrow \begin{bmatrix} 4 & -2 & 1 \\ -2 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} = 0 \rightarrow \text{span} \left( \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \right) \xrightarrow{\sqrt{1+4+4}=3\sqrt{5}} \begin{bmatrix} 0 \\ 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix}$$

$$\lambda = 0 \rightarrow \begin{bmatrix} 5 & -2 & 1 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} = 0 \rightarrow \text{span} \left( \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \right) \xrightarrow{\sqrt{1+4+1}=\sqrt{6}} \begin{bmatrix} 1/\sqrt{6} \\ 2/\sqrt{6} \\ -1/\sqrt{6} \end{bmatrix}$$

$$A^T A = \begin{bmatrix} \frac{5}{\sqrt{30}} & 0 & -\frac{1}{\sqrt{6}} \\ -\frac{2}{\sqrt{30}} & \frac{1}{\sqrt{5}} & -\frac{2}{\sqrt{6}} \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{6}} \end{bmatrix} \begin{bmatrix} 6 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} P^T \\ \dots \\ \dots \end{bmatrix}$$

This gives  $V^T$  for  $A = U\Sigma V^T$  because  $A^T A = V D V^T$

↑↑↑ Step 1 ✓

We get  $\Sigma$  from  $D$ :

$$\Sigma = \sqrt{D} \rightarrow \Sigma = \begin{bmatrix} \sqrt{6} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}_{2 \times 3}$$

$$\begin{bmatrix} 6 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$U = \boxed{?}$$

$$\boxed{AA^T} = \boxed{P_2} D_2 P^T$$

$$\boxed{U = P}$$

$$\frac{\sqrt{6}}{1}$$

$$\boxed{u_i = \frac{1}{\|A\|_2} A v_i}$$

$$\left\{ \begin{array}{l} u_1 = \frac{1}{\|A\|_2} A v_1 = \frac{1}{\sqrt{6}} \\ u_1 = \begin{bmatrix} \frac{1}{\sqrt{6}} \\ -2/\sqrt{6} \end{bmatrix} \end{array} \right.$$

$$u_2 = \frac{1}{\|A\|_2} A v_2 \Rightarrow$$

$$u_2 = \begin{bmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix}$$

OR

$$\begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} \frac{5}{\sqrt{30}} \\ -2/\sqrt{30} \\ 1/\sqrt{30} \end{bmatrix}$$

$$U = \frac{1}{\sqrt{5}} \begin{bmatrix} 1 & 2 \\ -2 & 1 \end{bmatrix}$$

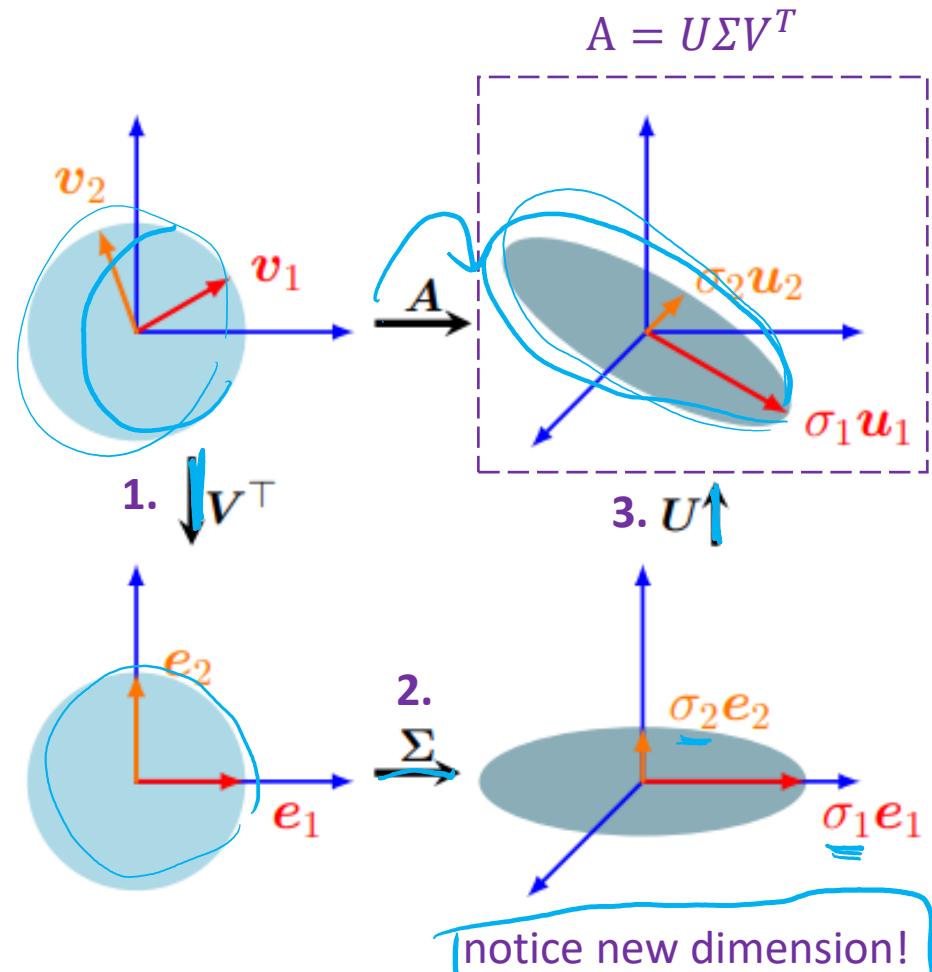
# Geometric Interpretation of SVD

- A **rectangular matrix** ( $n \times m$ ) can be factored into:

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

- These can be seen as a sequence of transformations:

1. rotation by right-singular  $V$
2. scaling by singular values  $\Sigma$
3. rotation by left-singular again  $U$



# Example

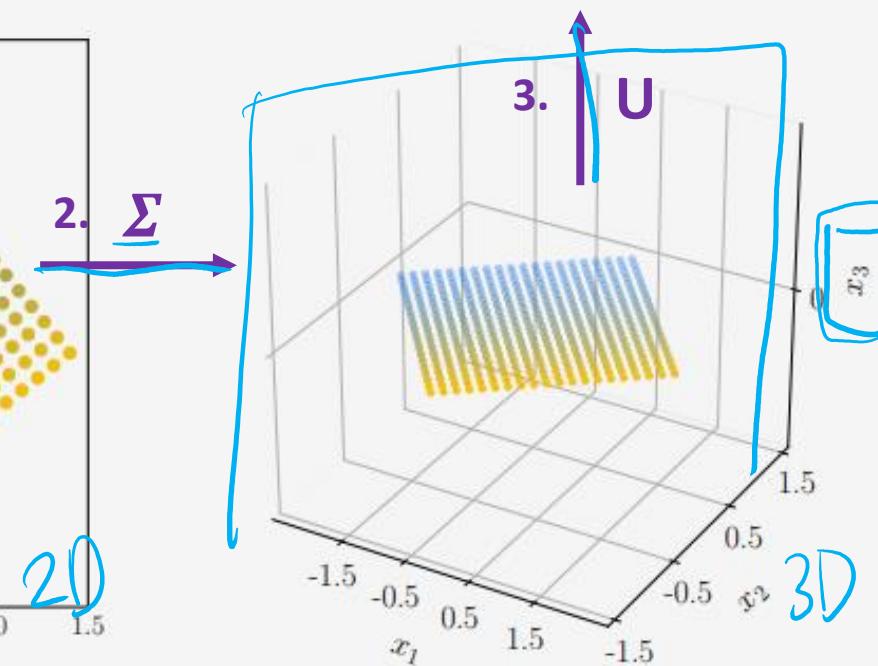
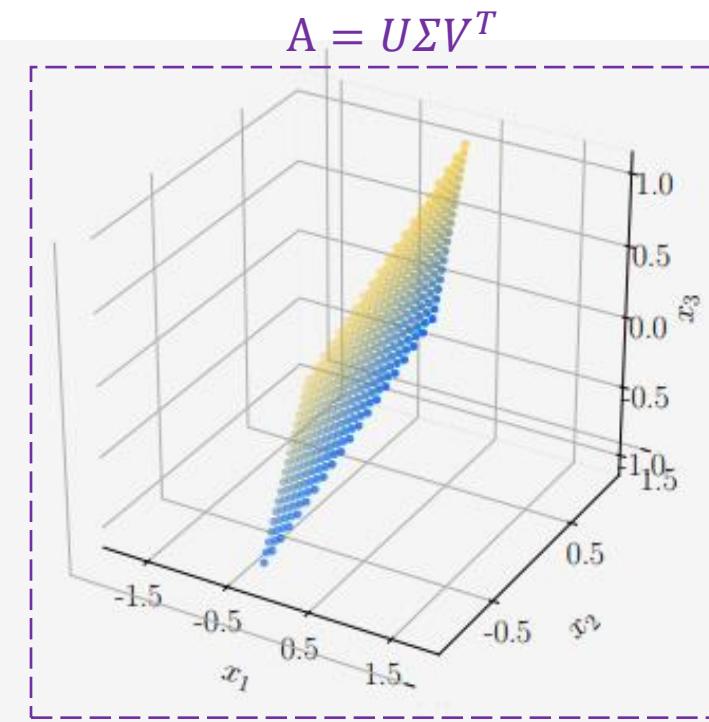
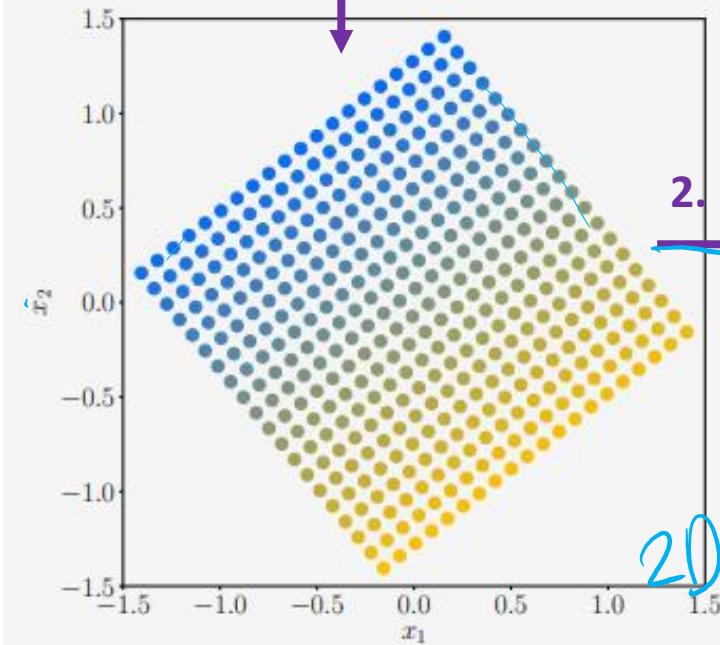
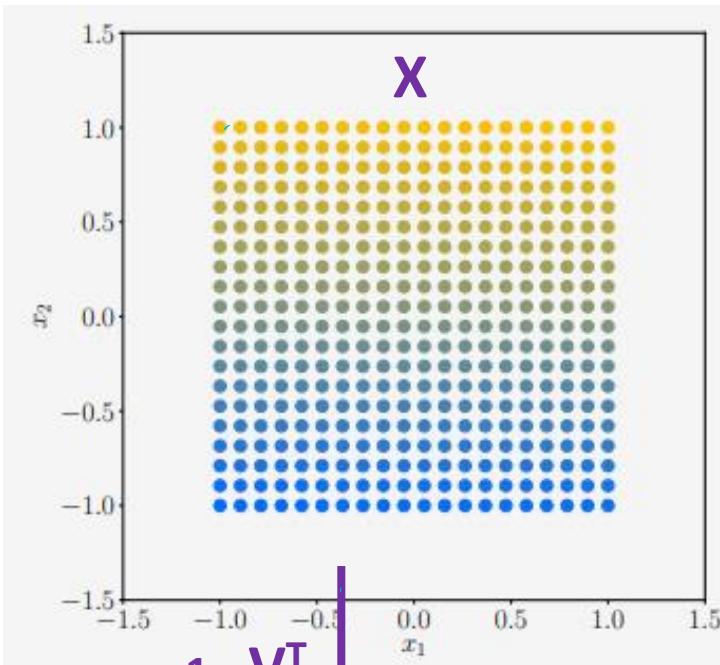
Given: A set of data points X  
and transformation A.

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & -0.8 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} = \mathbf{U}\Sigma\mathbf{V}^T \\ &= \begin{bmatrix} -0.79 & 0 & -0.62 \\ 0.38 & -0.78 & -0.49 \\ -0.48 & -0.62 & 0.62 \end{bmatrix} \begin{bmatrix} 1.62 & 0 \\ 0 & 1.0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} -0.78 & 0.62 \\ -0.62 & -0.78 \end{bmatrix} \end{aligned}$$

$\checkmark$        $\sum$        $\checkmark^T$

X -> set of vectors

1.  $\mathbf{V}^T$  -> rotates data X
2.  $\Sigma$  -> singular matrix maps onto  $\mathbb{R}^3$  (3<sup>rd</sup> dimension is 0, stretched in a plane by singular values)
3.  $\mathbf{U}$  -> rotation within  $\mathbb{R}^3$



# Similarities to Eigendecomposition -1

## ➤ Eigenvalue Decomposition vs Singular Value Decomposition:

- SVD exists for any matrix ( $m \times n$ ) whether square or rectangular
- SVD left and right singular vectors can be made orthonormal
- Both are a composition of three linear mappings:
  1. Change of basis in the domain
  2. Independent scaling of each new basis vector from domain to codomain
  3. Change of basis in codomain
- SVD domain and codomain can be vector spaces of different dimensions.

Eigendecomposition

$$A = PDP^{-1}$$

Singular Value Decomposition

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

# Similarities to Eigendecomposition -2

- SVD singular values are real and non-negative
- SVD and eigendecomposition are closely related through their projections:
  - The left-singular vectors of A are eigenvectors of  $AA^T$
  - The right-singular vectors of A are eigenvectors of  $A^TA$
  - The nonzero singular values of A are the square roots of nonzero eigenvalues of  $AA^T$  and  $A^TA$
- For symmetric square matrices ( $n \times n$ )  
eigendecomposition and SVD are the same  
when the features are normalized.

Eigendecomposition

$$A = PDP^{-1}$$

Singular Value  
Decomposition

$$A = \underline{U} \Sigma \underline{V}^T$$

$4 \times 4$        $3 \times 3$

$$\rightarrow P = U = V$$

# Relationship between PCA and SVD

- PCA and SVD are closely related:

$$\begin{aligned} \underline{\text{Cov}(A)} &= \frac{A^T A}{n - 1} \\ &= \frac{(U\Sigma V^T)^T (U\Sigma V^T)}{n - 1} \\ &= \frac{V\Sigma U^T U\Sigma V^T}{n - 1} \\ &= V \frac{\Sigma^2}{n - 1} V^T \end{aligned}$$

$$\begin{matrix} n \\ m \end{matrix} \boxed{A} = \begin{matrix} m \\ m \end{matrix} \boxed{U} \begin{matrix} n \\ m \end{matrix} \boxed{\Sigma} \begin{matrix} n \\ m \end{matrix} \boxed{V^T}$$

SVD columns of  $V$  are principal direction in PCA

- The result is the same form as eigendecomposition of  $A$  and hence:

$$D = \frac{\Sigma^2}{n - 1}$$

eigenvalues    singular values

# Applications of Matrix Decompositions

Dimensionality Reduction  
Data Interpretation

# Recap: Eigenfaces

- Last time we saw how PCA could be used to obtain a compressed version of face images.

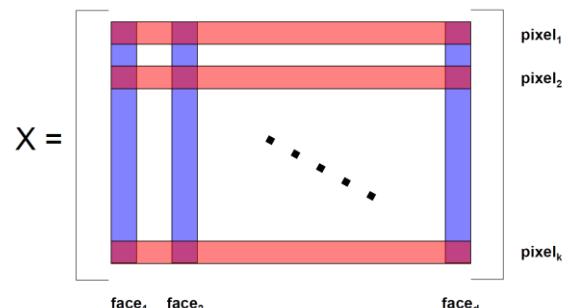
Normalized Face Data



1.  $f_{260}$

The Data Matrix X

- Matrix with columns as faces, rows as pixels

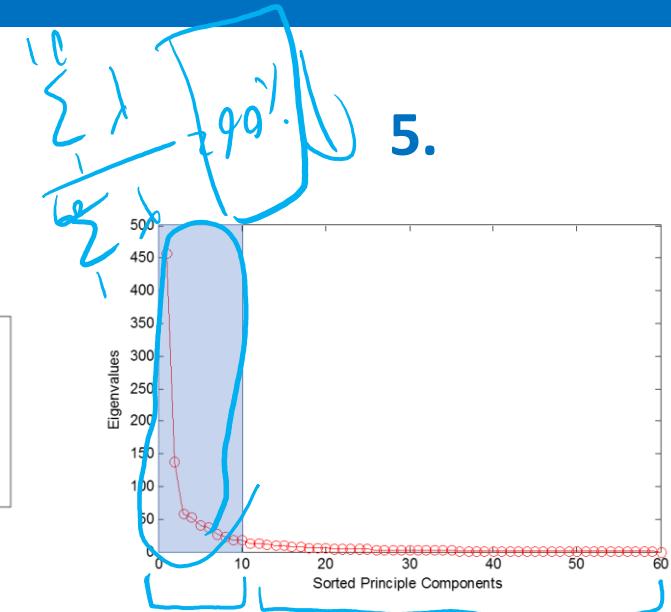


2.

The Covariance Matrix

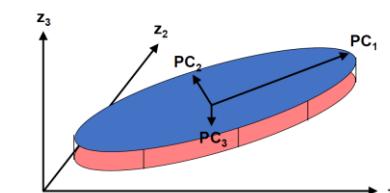
$$\text{Cov}(f) = f f^T = \begin{bmatrix} f_1^2 & f_{1j} \\ f_{ij} & \ddots \end{bmatrix}$$

3.



4. Principle Components

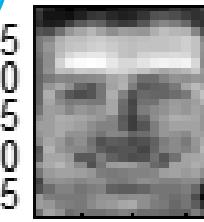
- Why look at eigenvectors of covariance?



- If data lives in linear subspace...

- Covariance indicates principle data dimensions
- Then eigenvectors = 'principle data components'

5.



5 10 15

6.  $f_{2/10}$



27

# Algorithmic complexity

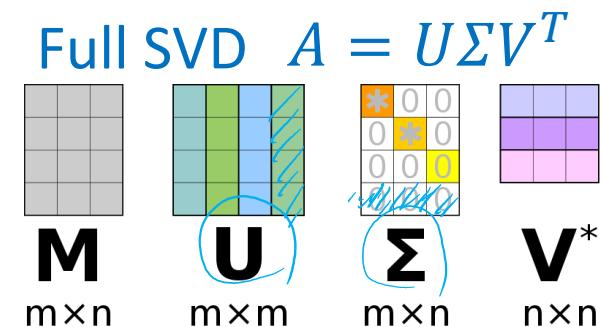
## ➤ Eigenvalue Decomposition:

- For an  $m \times n$  matrix, is  $O(n^3)$  if there are  $n$  features
- Impractical for large  $n$

$$A = PDP^{-1}$$

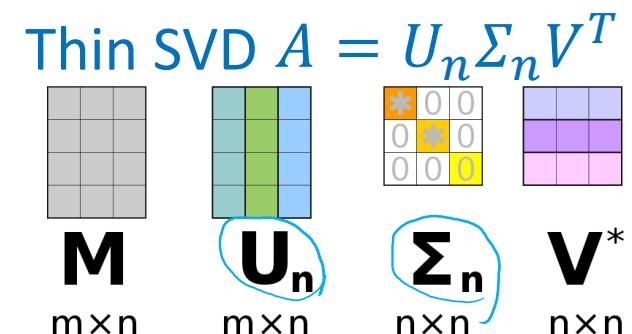
## ➤ Full SVD:

- For an  $m \times n$  matrix, it is  $O(mn \min(n, m))$
- More computationally tractable for large  $n$  or large  $m$



## ➤ Thin (economy-sized) SVD:

- For an  $m \times n$  matrix, it is  $O(\max(m, n) \cdot \min(m, n))$
- Even faster and more economical

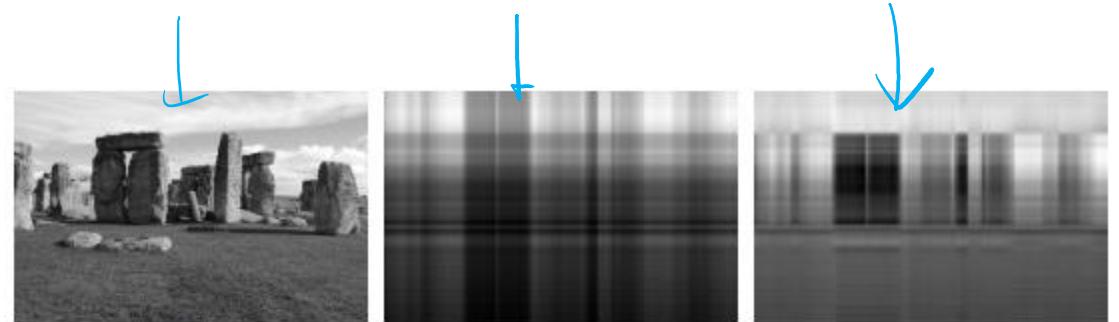


# Image Compression

- Analogous to decomposing face we can also decompose a single image for compression.

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top \in \mathbb{R}^{m \times n}$$

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = \underbrace{\sum_{i=1}^r \sigma_i \mathbf{A}_i}_{\text{rank } r}$$



(a) Original image  $\mathbf{A}$ . (b) Rank-1 approximation  $\hat{\mathbf{A}}(1)$ . (c) Rank-2 approximation  $\hat{\mathbf{A}}(2)$ .



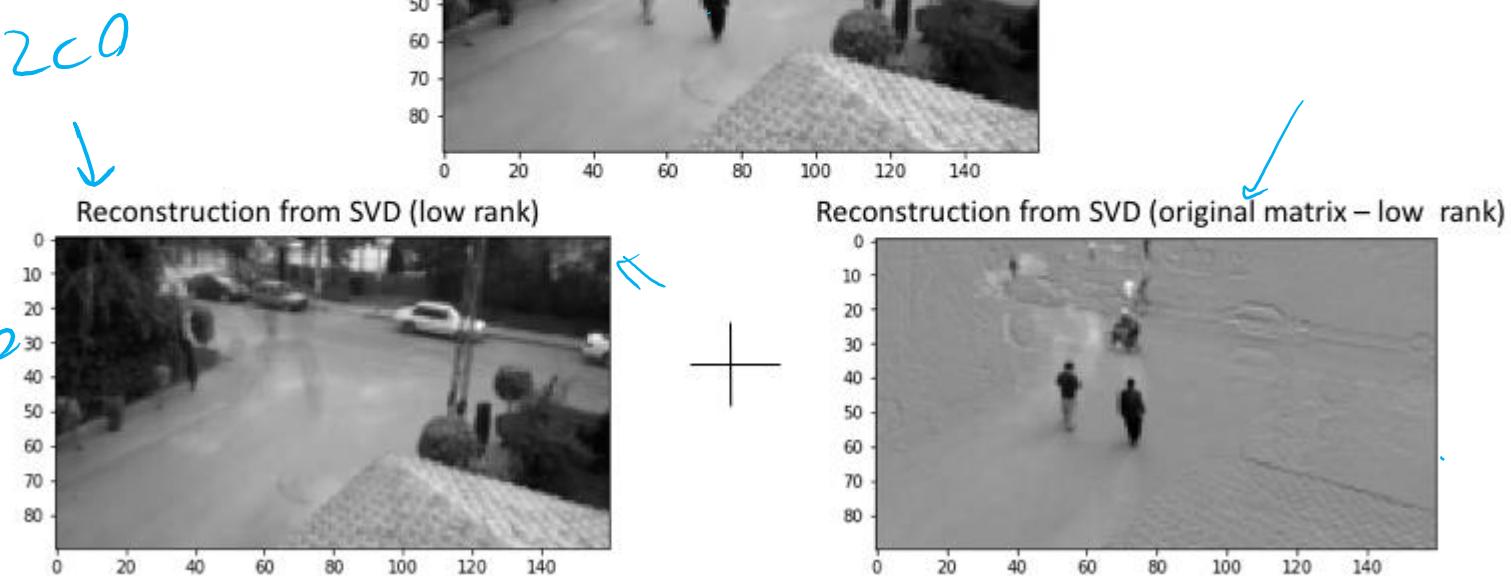
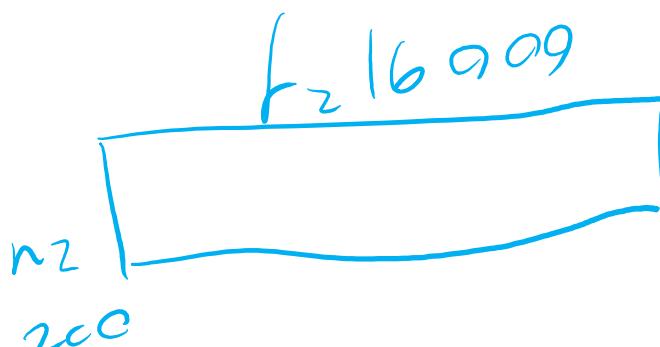
(d) Rank-3 approximation  $\hat{\mathbf{A}}(3)$ . (e) Rank-4 approximation  $\hat{\mathbf{A}}(4)$ . (f) Rank-5 approximation  $\hat{\mathbf{A}}(5)$ .

# Application 1 – Background Removal

- Given a video we can use SVD to remove the background

Each frame → Sample =  $200 \times 200 = 200^2 = 40000$

Each Pixel → Feature =  $16 \times 16 = 256$



$$100 \times 160 = 16000$$

$$10 \times 20 = 200$$

# Example Google Colab Code

# Application 2 – Recommender Systems

➤ SVD can also be used to find structure in data for making recommendations.

➤ SVD is strongly related to recommender systems applications:

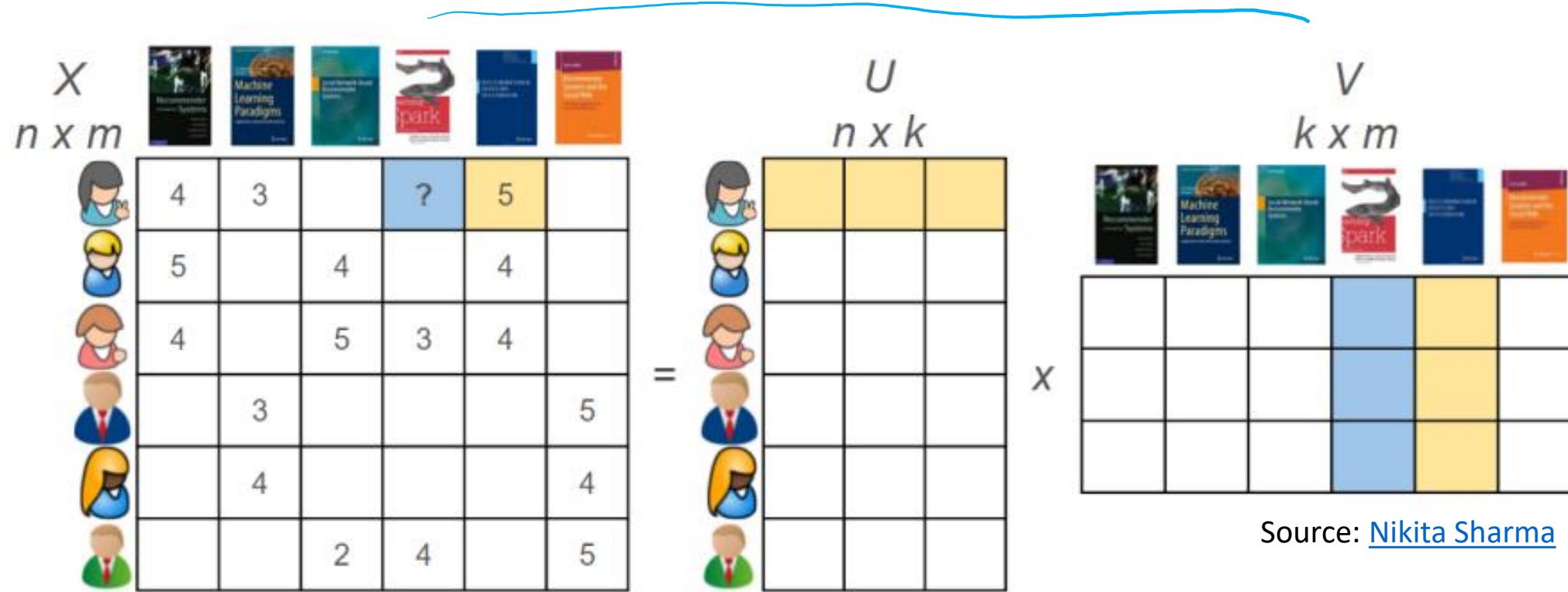
- Movie recommendations
- Product recommendations
- Restaurant recommendations
- Website recommendations
- ...

SVD

The diagram illustrates the Singular Value Decomposition (SVD) of a movie rating matrix. The matrix is shown as a grid of numbers, with rows representing movies and columns representing users. The rows are labeled with movie titles: Star Wars, Blade Runner, Amelie, and Delicatessen. The columns are labeled with user names: Ali, Beatrix, Chandra, and D. The matrix entries range from 0 to 5. Handwritten blue arrows point from the matrix to the decomposition components. A curved arrow points from the top-left of the matrix to the first column of the matrix  $V$ . Another curved arrow points from the bottom-right of the matrix to the last row of the matrix  $U$ . A third curved arrow points from the middle of the matrix to the diagonal elements of the matrix  $\Sigma$ .

$$= \begin{bmatrix} -0.6710 & 0.0236 & 0.4647 & -0.5774 \\ -0.7197 & 0.2054 & -0.4759 & 0.4619 \\ -0.0939 & -0.7705 & -0.5268 & -0.3464 \\ -0.1515 & -0.6030 & 0.5293 & -0.5774 \end{bmatrix}$$
$$= \begin{bmatrix} 9.6438 & 0 & 0 \\ 0 & 6.3639 & 0 \\ 0 & 0 & 0.7056 \\ 0 & 0 & 0 \end{bmatrix}$$
$$= \begin{bmatrix} -0.7367 & -0.6515 & -0.1811 \\ 0.0852 & 0.1762 & -0.9807 \\ 0.6708 & -0.7379 & -0.0743 \end{bmatrix}$$

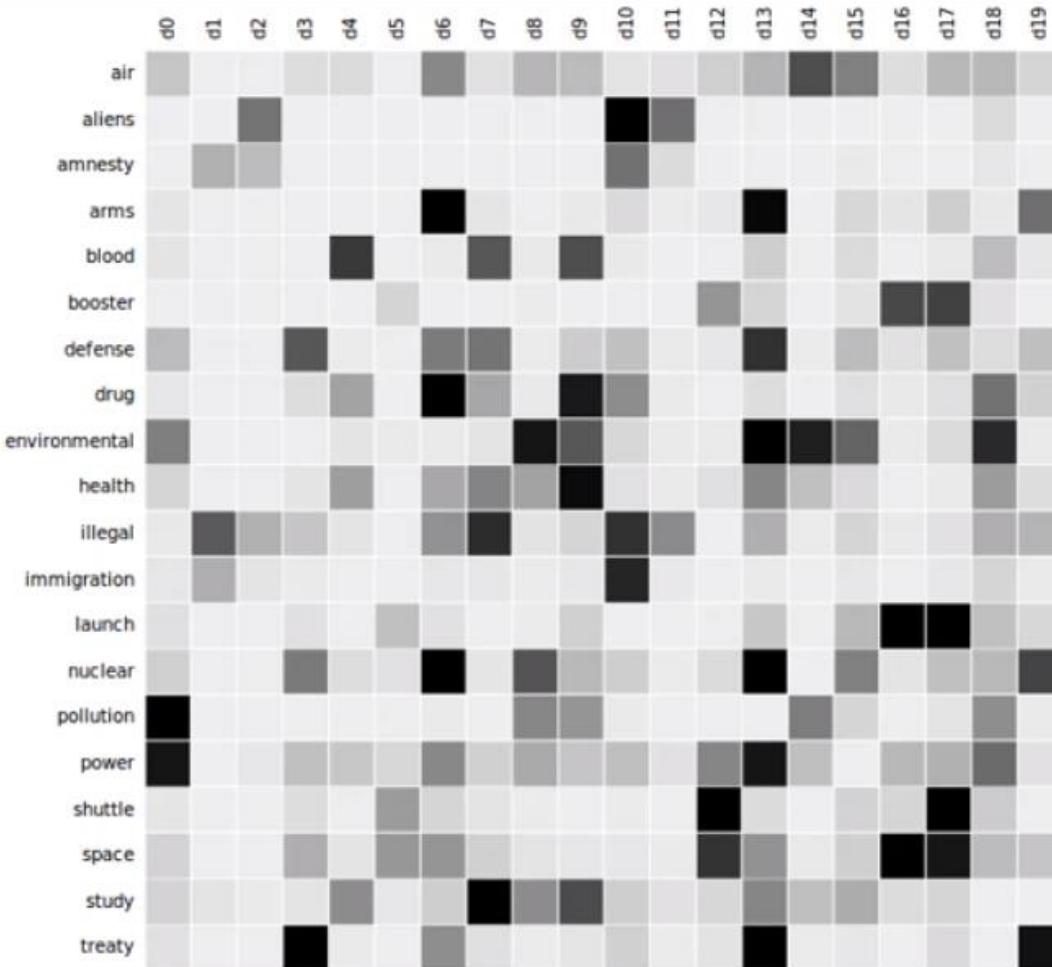
# Example: Collaborative Filtering



Matrix decomposition using SVD does not work well with missing data.  
Gradient descent can be used to learn  $U$  and  $V$  matrices to make movie recommendations.

# Example Google Colab Code

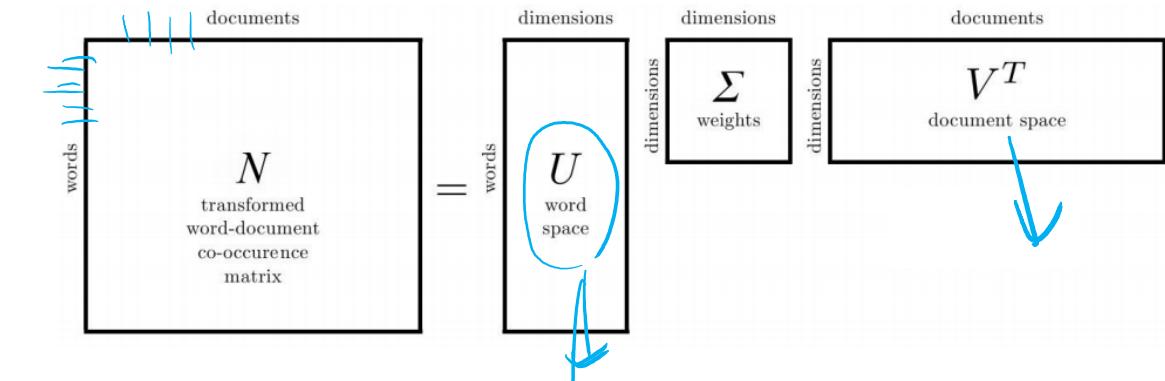
# Application 3 – Text Embeddings



Source: Wikipedia

- Word relationships from documents
- Latent semantic representation

10 three



# Example Google Colab Code

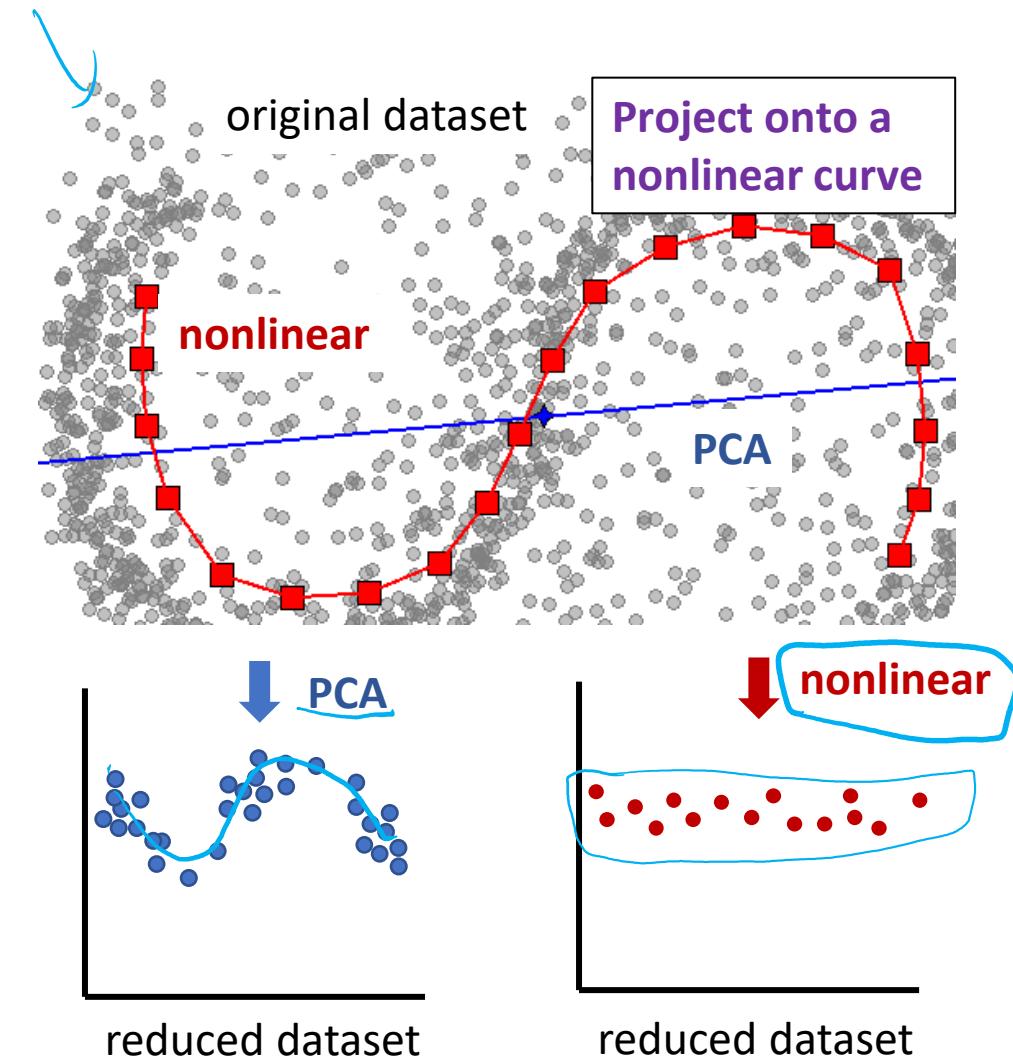
# Summary

- Eigendecomposition and Singular Value Decomposition offer ways to factor matrices, much like we can factor numbers into primary numbers.
- This can lead to dimensionality reduction, data insights and algorithm speed ups.
- Eigendecomposition is limited to square matrices, so we use **SVD which is guaranteed to work in all circumstances**.
- Both apply only for linear transformations/mappings

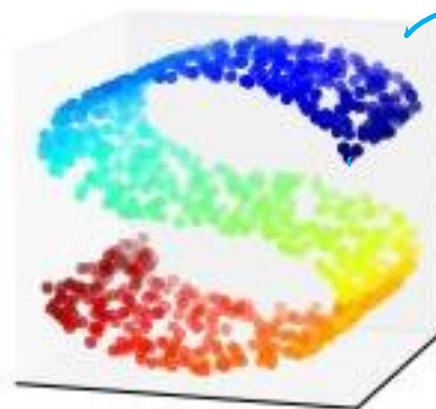
10 : 25

# Nonlinear Dimensionality Reduction

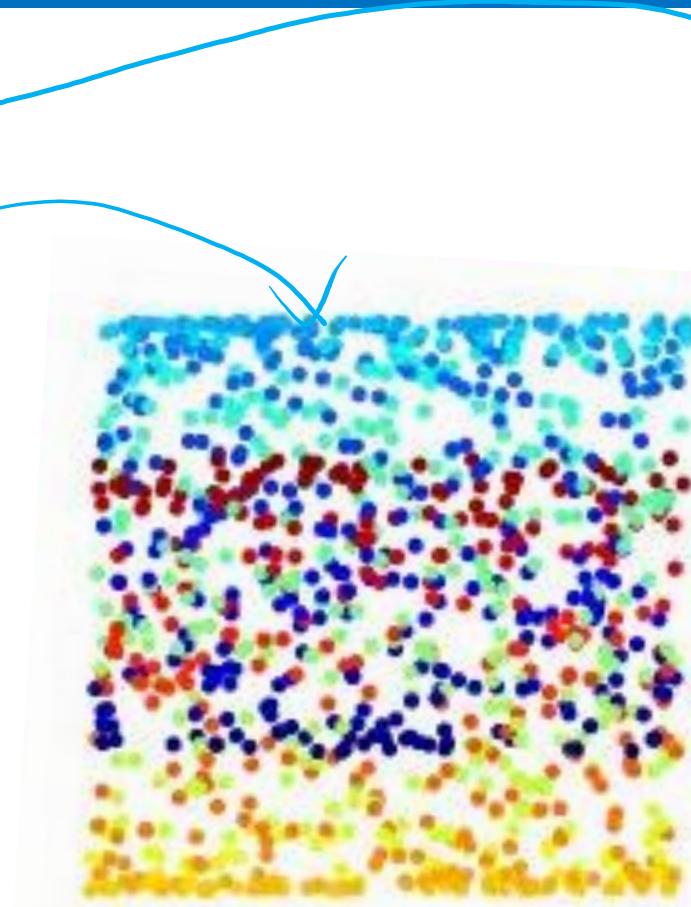
- Data can have nonlinear relationships for which our decomposition techniques would be ineffective.
- Data Visualization/Reduction
  - t-SNE
  - Isomap
  - LLE
  - Kernel PCA
  - Deep Autoencoders  
(type of neural network)



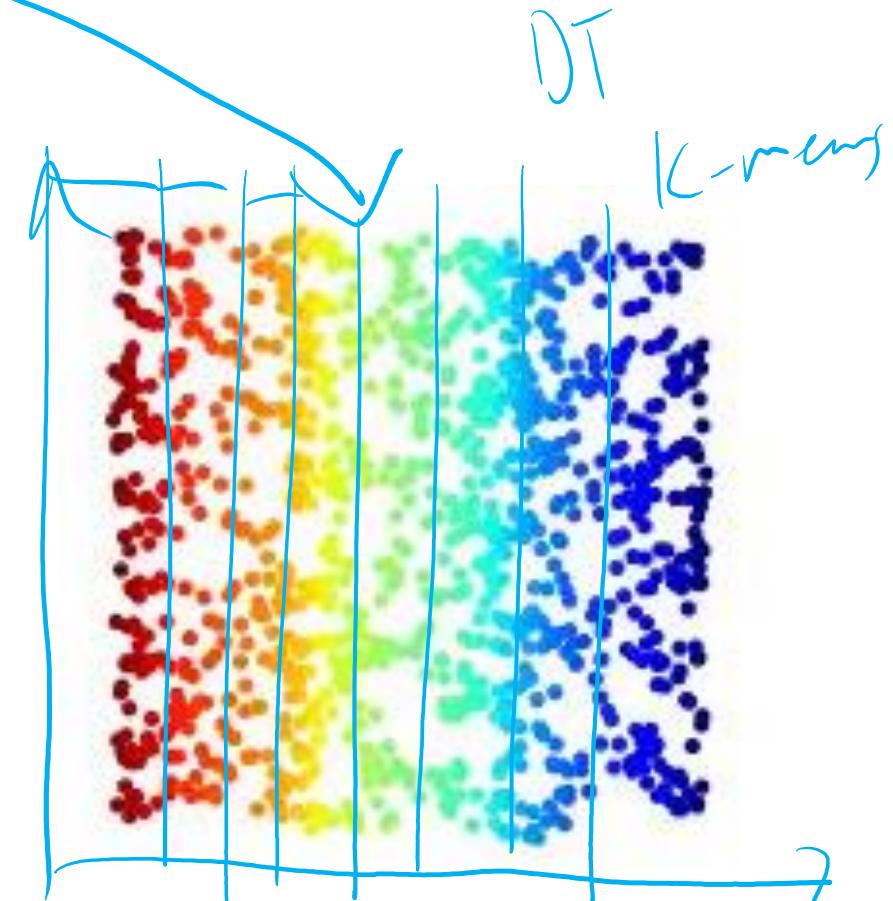
# Example: Nonlinear Projections



Nonlinear Data



PCA Projection



Isomap Projection

# Vector Calculus

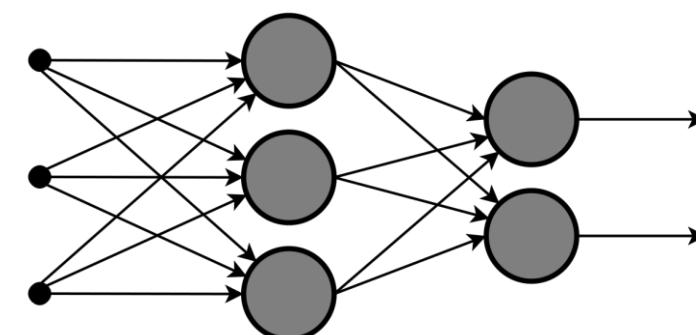
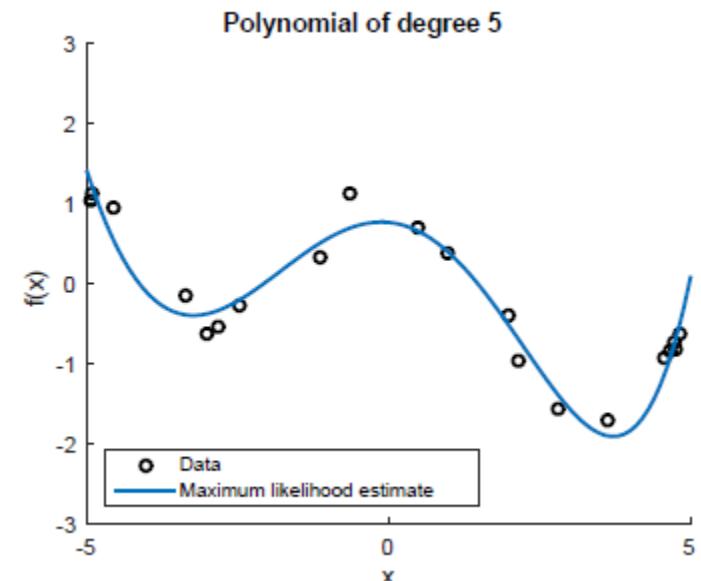


**Readings:**

- **MML Chapter 5.1-4**

# Why Vector Calculus?

- Vector calculus is used extensively in optimization
  - We will see it the next couple of lectures when we discuss model-based learning using **linear regression** and **neural network models**
- Most of the python frameworks such as Tensorflow, PyTorch, etc. handle these calculations using numerical methods



# Types of Differentiation

1. Scalar differentiation:  $f : \mathbb{R} \rightarrow \mathbb{R}$

$y \in \mathbb{R}$  w.r.t.  $x \in \mathbb{R}$

2. Multivariate case:  $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$y \in \mathbb{R}$  w.r.t. vector  $x \in \mathbb{R}^N$

gradient  $\left[ \frac{\partial f}{\partial x_1}, \dots \right]$

3. Vector fields:  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

vector  $y \in \mathbb{R}^M$  w.r.t. vector  $x \in \mathbb{R}^N$

Jacobian  
 $M \times N$

4. General derivatives:  $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$

matrix  $y \in \mathbb{R}^{P \times Q}$  w.r.t. matrix  $X \in \mathbb{R}^{M \times N}$

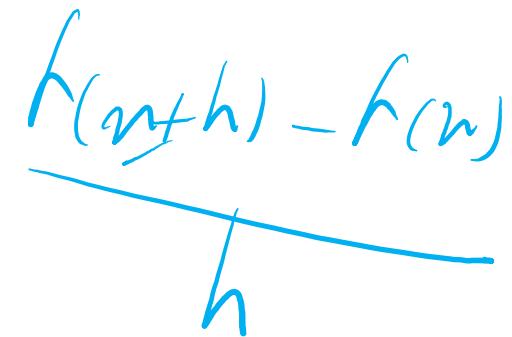
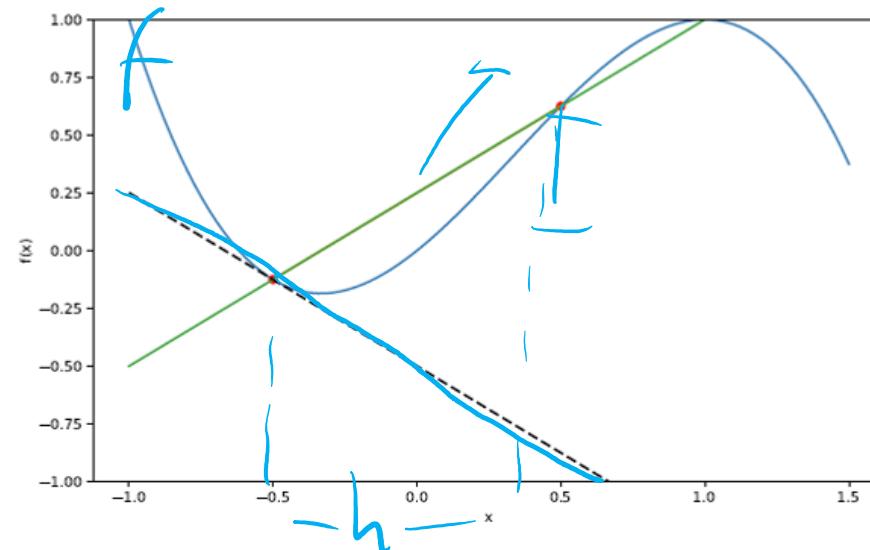
Tensor

# 1. Scalar Differentiation

- Derivative defined as the limit of the difference quotient:

$$f'(x) = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- Slope of a secant line through  $f(x)$  and  $f(x+h)$



# Some Examples

$$\curvearrowright f(x) = x^n$$

$$f(x) = \sin(x)$$

$$f(x) = \tanh(x)$$

$$f(x) = \exp(x)$$

$$f(x) = \log(x)$$

$$\curvearrowright f'(x) = nx^{n-1}$$

$$f'(x) = \cos(x)$$

$$f'(x) = 1 - \tanh^2(x)$$

$$f'(x) = \exp(x)$$

$$f'(x) = \frac{1}{x}$$

# Rules

- ▶ Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

- ▶ Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x) = \underbrace{\frac{df(x)}{dx}g(x)}_{\text{derivative of } f(x)} + f(x)\underbrace{\frac{dg(x)}{dx}}_{\text{derivative of } g(x)}$$

- ▶ Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg(f(x))}{df} \frac{df(x)}{dx}$$

- ▶ Quotient Rule

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f(x)'g(x) - f(x)g(x)'}{(g(x))^2} = \frac{\frac{df}{dx}g(x) - f(x)\frac{dg}{dx}}{(g(x))^2}$$

# Example: Scalar Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = \cancel{g}'(\cancel{f(x)}) \cancel{f'(x)} = \cancel{\frac{dg}{df}} \frac{df}{dx}$$

$$\begin{aligned} g(z) &= 6z + 3 \\ z &= f(x) = -2x^2 + 5 \end{aligned}$$

$$(g \circ f)'(x) = \underbrace{(6)}_{dg/df} \underbrace{(-2)}_{df/dx}$$

$$= -12$$

## 2. Multivariate Differentiation $f: \mathbb{R}^N$ to $\mathbb{R}$

$$y = \underline{f(\underline{x})},$$

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

➤ Partial derivative (change one coordinate at a time)

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_N) - f(x)}{h}$$

➤ The gradient collects all partial derivatives:

$$\frac{df}{dx} = \left[ \frac{\partial f}{\partial x_1} \quad \cdots \quad \frac{\partial f}{\partial x_N} \right] \in \mathbb{R}^{1 \times N} \quad \text{results in a row vector}$$

# Example: Multivariate Differentiation

➤ Given:

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$$

$\begin{matrix} 1 \\ \times \\ 2 \end{matrix}$

➤ Partial derivative

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 x_2 + x_2^3$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = x_1^2 + 3x_1 x_2^2$$

➤ Gradient

$$\frac{df}{dx} = \left[ \frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right] \in \mathbb{R}^{1 \times 2}$$

$$= \boxed{\begin{bmatrix} 2x_1 x_2 + x_2^3 & x_1^2 + 3x_1 x_2^2 \end{bmatrix}}$$

### 3. Vector Field Differentiation $f: \mathbb{R}^N$ to $\mathbb{R}^M$

$$\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

➤ Jacobian matrix (collection of all partial derivatives)

$$\begin{bmatrix} \frac{dy_1}{d\mathbf{x}} \\ \vdots \\ \frac{dy_M}{d\mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \dots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

# Example: Vector Field Differentiation

➤ Given:

$$f(\boldsymbol{x}) = \mathbf{A}\boldsymbol{x}, \quad f(\boldsymbol{x}) \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \boldsymbol{x} \in \mathbb{R}^N$$

| | : | \zeta

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\boldsymbol{x}) \\ \vdots \\ f_M(\boldsymbol{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

➤ Compute the Jacobian

$$f_i(\boldsymbol{x}) = \sum_{j=1}^N A_{ij}x_j$$

$$\frac{\partial f_i}{\partial x_j} = A_{ij}$$

$$\frac{df}{d\boldsymbol{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = \mathbf{A} \in \mathbb{R}^{M \times N}$$

# Dimensionality of the Jacobian

- In general: A function  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$  has a gradient that is an  $M \times N$  matrix with:

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Jacobian dimension: #target dimensions x #input dimensions

# Example: Chain Rule

Given  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $x: \mathbb{R} \rightarrow \mathbb{R}^2$ ,  $t \in \mathbb{R}$

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \times \frac{dx}{dt}$$

What are the dimensions of  $\frac{df}{dx}$  and  $\frac{dx}{dt}$ ?

$$1 \times 2 \text{ and } 2 \times 1$$

Compute the gradient  $\frac{df}{dt}$  using the chain rule:

$$\begin{aligned} \frac{df}{dt} &= \frac{df}{dx} \frac{dx}{dt} = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right] \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} = \begin{bmatrix} 2 \sin t & 2 \end{bmatrix} \begin{bmatrix} \cos t \\ -\sin t \end{bmatrix} \\ &= 2 \sin t \cos t - 2 \sin t = 2 \sin t(\cos t - 1) \end{aligned}$$

$$\begin{aligned} f(x) &= f(x_1, x_2) = x_1^2 + 2x_2, \\ x(t) &= \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix} \end{aligned}$$

$$\frac{\partial f}{\partial x_1} = 2x_1, \quad \frac{\partial f}{\partial x_2} = 2$$

$$\begin{bmatrix} 2x_1 & 2 \end{bmatrix} \begin{bmatrix} \cos t \\ -\sin t \end{bmatrix}$$

$$\begin{bmatrix} 2 \sin t & 2 \end{bmatrix} \begin{bmatrix} \cos t \\ -\sin t \end{bmatrix}$$

# 4. Derivatives with Respect to Matrices

- Recall: A function  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$  has a gradient that is an  $M \times N$  matrix with:

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: #target dimensions  $\times$  #input dimensions

- This generalizes to when the inputs (**N**) or targets (**M**) are **matrices**.

- Function  $f: \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$  has a gradient that is a  $(P \times Q) \times (M \times N)$  **tensor**:

$$\frac{df}{dX} \in \mathbb{R}^{(P \times Q) \times (M \times N)}, \quad df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

# More Examples:

# Example:

Given:  $f(x) = Ax$ ,  $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ ,  $x = [x_1 \ x_2]^T$ , find  $\frac{df}{dx}$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}_{2 \times 2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} x_1 + 2x_2 \\ 3x_1 + 4x_2 \end{bmatrix}^T_{2 \times 1} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}^T_{2 \times 1}$$
$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$g(n) \approx f_n$$
$$g' \approx f$$

# Example:

Given:  $f(x) = \underline{x^T A x}$ ,  $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ ,  $x = [x_1 \ x_2]^T$ , find  $\underline{df(x)/dx}$

$$f(x) = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \underbrace{\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}}_{2 \times 2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 \end{bmatrix}_{1 \times 2} \begin{bmatrix} x_1 + 2x_2 \\ 3x_1 + 4x_2 \end{bmatrix}_{2 \times 1} = x_1^2 + \underbrace{2x_1x_2 + 3x_1x_2}_{5x_1x_2} + 4x_2^2$$

$$\frac{df}{dx} = \begin{bmatrix} 2x_1 + 5x_2 & 5x_1 + 8x_2 \end{bmatrix}_{1 \times 2}$$

# PCA by derivation

- Requires some familiarity with vector calculus for matrix differentiation
- Suggested reading: [Deep Learning Textbook \(pgs 45 – 50\)](#) by Ian Goodfellow, Yoshua Bengio and Aaron Courville.

# Next Time

- Week 8 Q/A Support Session: Project 3 Support ←
- Project 3 is due on 13 March at 23:00 (extended deadline) ✎
- Guest Lecture on 15 March at 10:00 ←
  - Dr. Sophie Lohmann: “Limits of measurement - who are we measuring?”
  - Zoom link <https://utoronto.zoom.us/j/86722516215>
- Week 9 Lecture – Linear Regression → March 16 wed
  - Monte Carlo Simulation
  - Empirical Risk Minimization
  - Maximum Likelihood
  - Probabilistic Modelling and Inference