

M5: Dimensionality Reduction via Principal Component Analysis (**PCA**) - [another popular (unsupervised ML) appn. of linear algebra]

Manikandan Narayanan

Week 9 (Sep 22-, 2025)

PRML Jul-Nov 2025 (Grads Section)

Acknowledgment of Sources

- Slides based on content from related
 - Courses:
 - IITM – Profs. Arun/Harish/Chandra’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited respectively as [AR], [HR]/[HG], [CC], [BR] in the bottom right of a slide.
 - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
 - Books:
 - PRML by Bishop. (content, figures, slides, etc.) – cited as **[CMB]**
 - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
 - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]
 - Information Theory, Inference and Learning Algorithms by David JC MacKay – [DJM]

Outline for Module M5

- M5. Dimensionality Reduction
 - **M5.0 Introduction**
 - M5.1 PCA
 - Intuition for two formulations
 - Maximizing variance formulation
 - Minimizing error formulation
 - M5.2 PCA applications/extensions/variants (very brief)

Context: Two unsupervised ML problems

- Unsupervised ML: Recognize patterns in the dataset (set of n data points in \mathbb{R}^d) without any labels on the data points.
- Dimensionality reduction:
 - Transforms data points from high to low dimensions without much loss of “information”, assuming the data points lie “effectively in/close to” a low-dim. manifold of the original space.
 - Many approaches possible: **PCA**, t-SNE, UMAP, **Laplacian Eigenmaps (spectral)**, etc.
- Clustering:
 - Grouping n objects into k clusters based on their similarity
 - Again, many approaches possible: k-means, hierarchical, **spectral**, etc.

Dimensionality Reduction

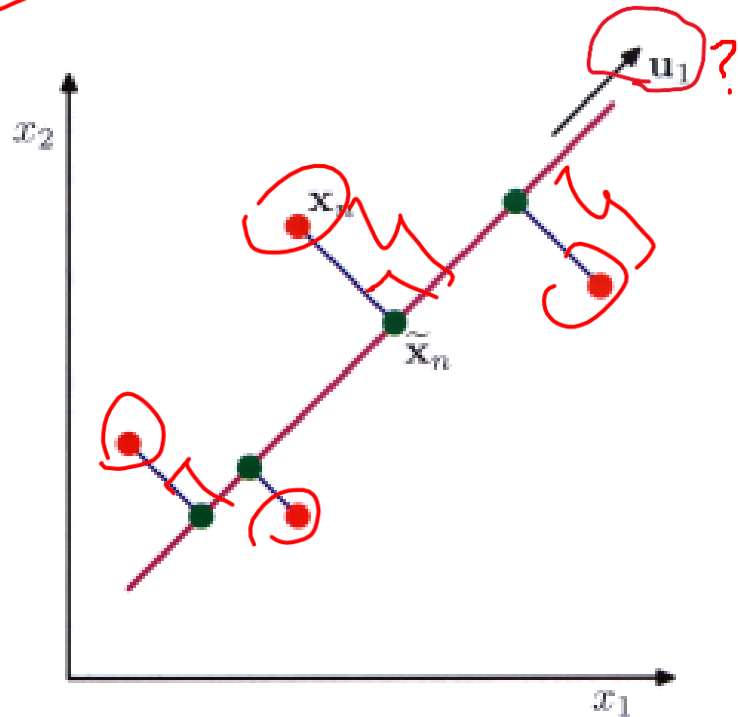
- Transforms data points from high to low dimensions without much loss of “information”, assuming the data points lie “effectively in/close to” a low-dim. manifold of the original space.
- A popular unsupervised PR/ML task like clustering with many apps.:
 - Visualization of data (in different domains)
 - Feature extraction/engineering (preprocessing step for classification/regression)
 - Lossy compression/Denoising
 - Preprocessing step for clustering (think spectral)
 - Identifying unknown confounding factors
- Our approach in this lecture: **Mostly [CBM]**, and (deterministic/hard) PCA (but just be aware that probabilistic/soft PCA (based on linear Gaussian) also exists).

Principal Component Analysis or PCA

- A widely-used dimensionality reduction method, aka Karhunen-Loeve transform
 - PCA represents D -dimensional data points using M -dimensional vectors with $M < D$, while maximizing variance (Hotelling, 1933) or equivalently minimizing error (Pearson, 1901), using the spectrum of the data matrix.
- Why is it called **linear** PCA?
 - Orthogonal projection of data onto a low-dimensional **linear** subspace, known as the principal subspace, s.t. the variance of the projected data is maximized or equivalently the mean squared distance between data pts and their projections is minimized.
 - Can be viewed as a continuous latent variable model (LVM) with **linear** Gaussian assumptions (with links to EM algo, which may be seen later during mixture density estimation if time permits)!

PCA in pictures and notations

$D=2 \rightarrow M=1$



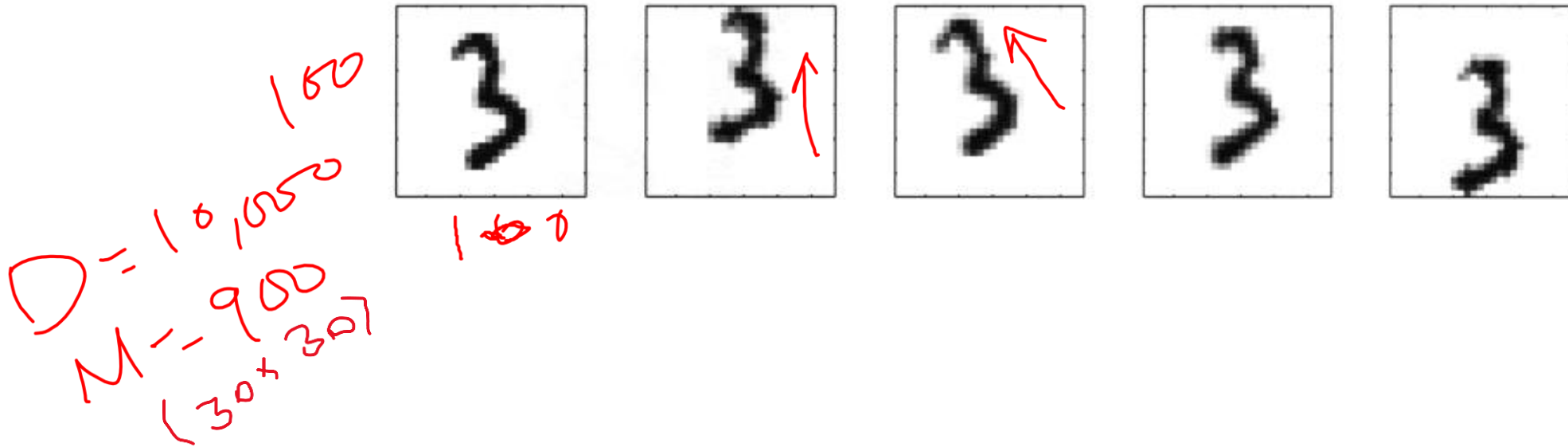
I/P: $x_1, \dots, x_N \in \mathbb{R}^D$

O/P: $\tilde{x}_1, \dots, \tilde{x}_N \in \mathbb{R}^D$

s.t:

- 1) \tilde{x}_n "approx." x_n
- 2) \tilde{x}_n can be repr. using $M < D$ dimensions.

Why dimension reduction works?: Latent space;
latent variables/features and their extraction



Face recogn.: Eigenfaces (latent space)



EV #1



EV #2



EV #3



EV #4



EV #5



EV #6



EV #7



EV #8



EV #9



EV #10



EV #50



EV #100



EV #150

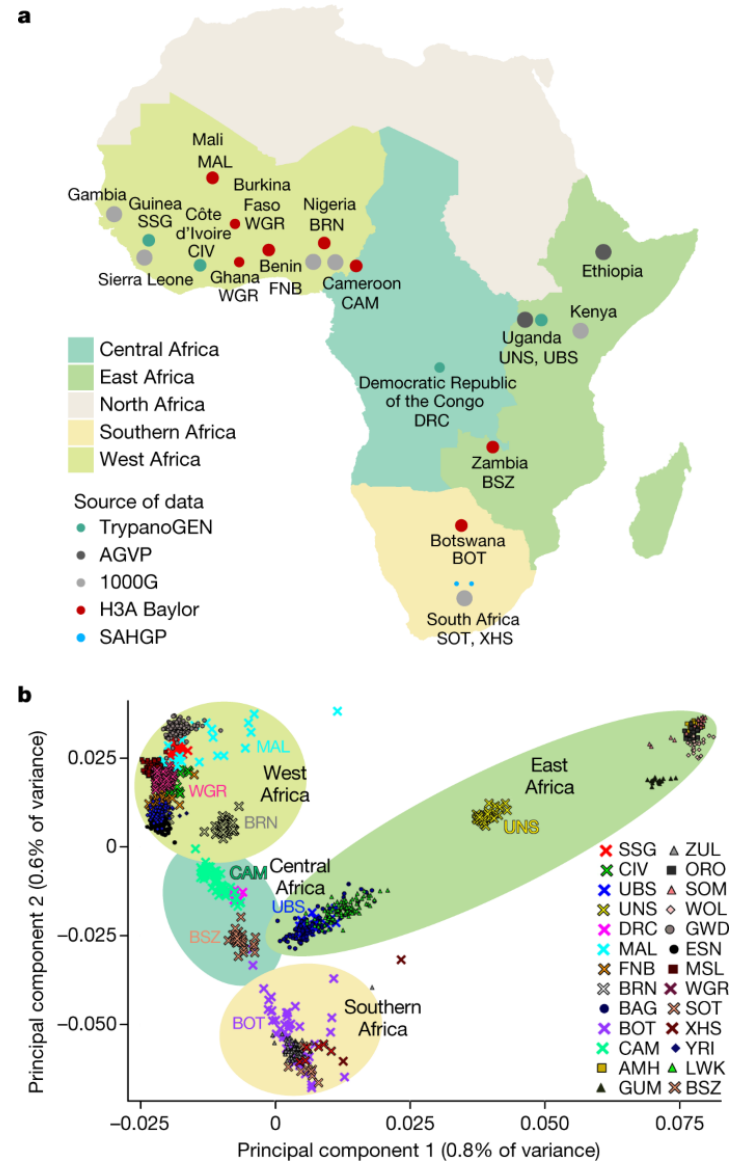


EV #200



EV #500

An example application: bioinf. visualization



[From Chowdhury et al. Nature 2020]

Single-cell revolution in biology – visualized in reduced dimensions

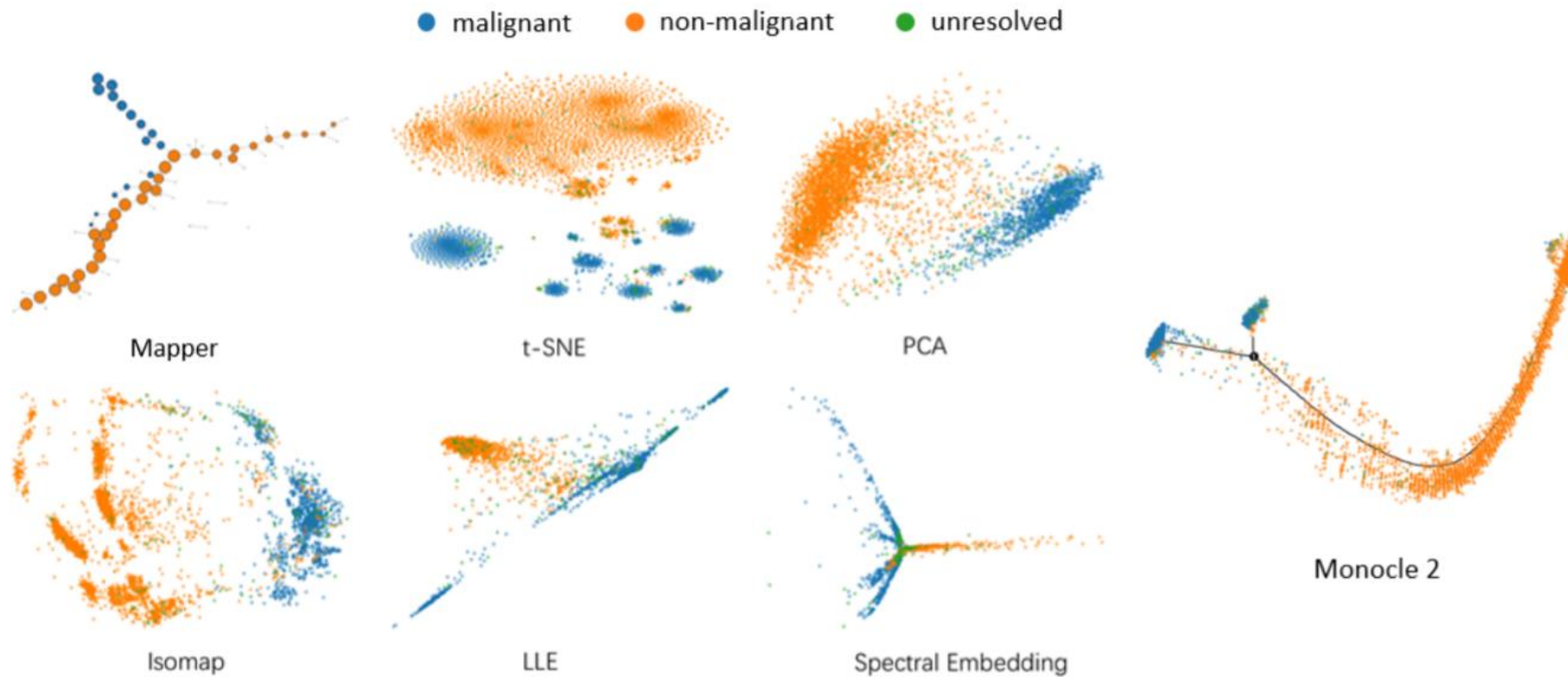


Figure 2. Visualization of melanoma cells.

Outline for Module M5

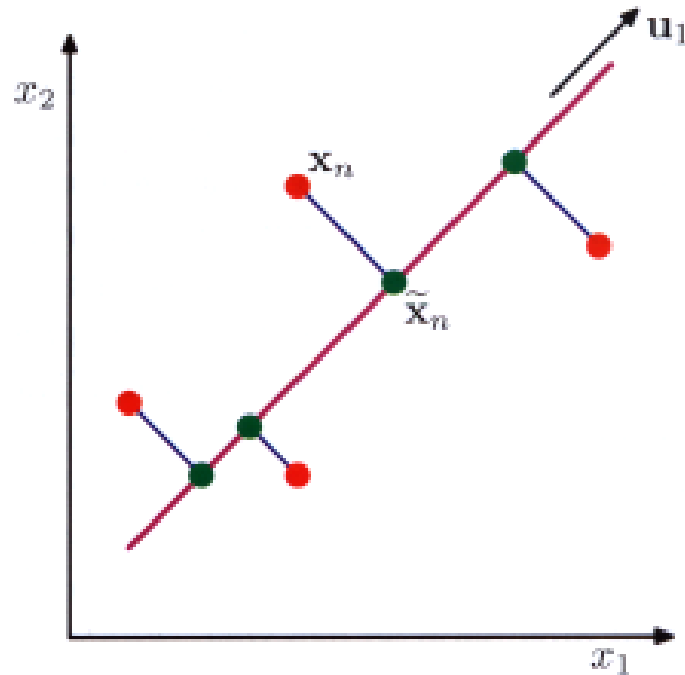
- M5. Dimensionality Reduction
 - M5.0 Introduction
 - **M5.1 PCA**
 - **Intuition for two formulations**
 - Maximizing variance formulation
 - Minimizing error formulation
 - M5.2 PCA applications/extensions/variants (very brief)

1) The M -dimensional subspace has to pass thru' the center of the data cloud (average \bar{x})!

- What if $M=0$?
- What if $M=1$?
- Illustrate in board (that u_1 passes thru' \bar{x})

2) PCA in pictures, and notation ($M=1$)

(Given that u_1 passes thru' \bar{x} , what is the optimal direction (rotation angle) of u_1 ?)

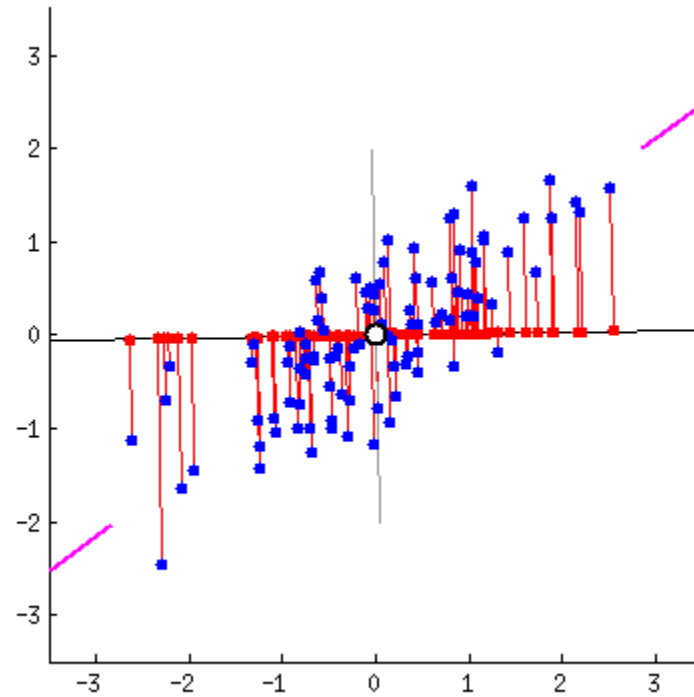


1/p: $x_1, \dots, x_N \in \mathbb{R}^D$
 o/p: Find $u_1 \in \mathbb{R}^D$ s.t.

Var $(x_1^T u_1, \dots, x_N^T u_1)$ is
 maximized (s.t. $\|u_1\|=1$)

$$\Leftrightarrow \min \frac{1}{N} \sum_{n=1}^N \| \tilde{x}_n - x_n \|^2$$

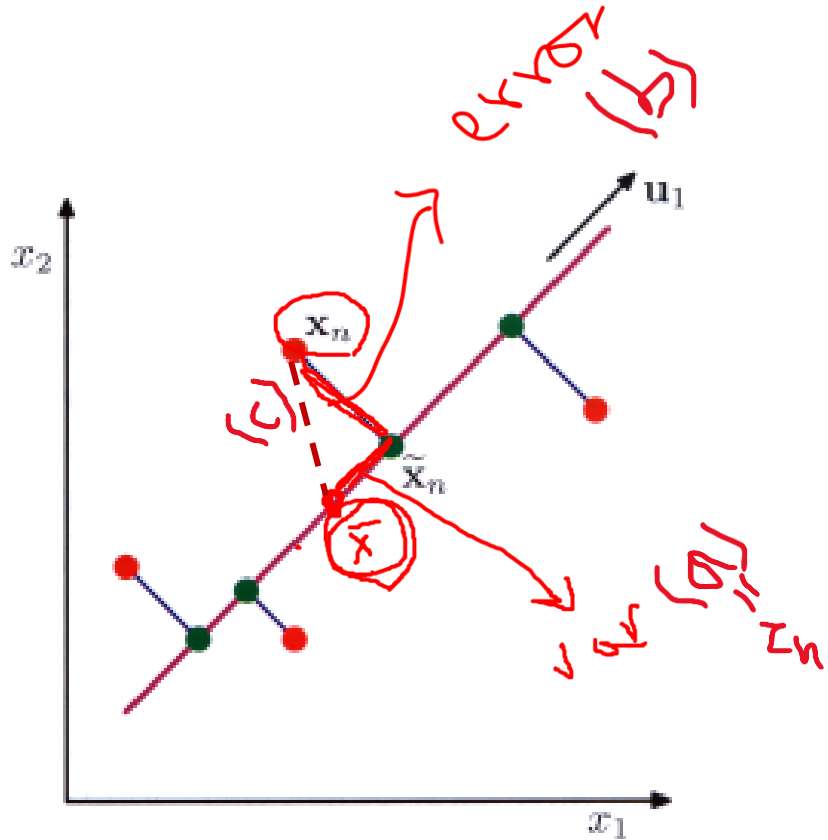
3) Animation of the two equiv. objectives!



<https://stats.stackexchange.com/a/140579>

(Also, show book example for 3D \rightarrow 2D reduction!)

4) Formal statement of the two objectives ($M=1$)



$$\begin{aligned}
 \text{I/P: } & X_1, \dots, X_N \in \mathbb{R}^D \\
 \text{O/P: Find } & u_1 \in \mathbb{R}^D \text{ s.t.} \\
 & \text{Var}(X_1^T u_1, \dots, X_N^T u_1) \text{ is} \\
 & \text{maximized (s.t. } \|u_1\| = 1)
 \end{aligned}$$

$$\max \text{Var}(\{z_n\}) = \text{Var}(\{X_n^T u_1\}) \iff \min \frac{1}{N} \sum_{n=1}^N \left\| \underbrace{\sum_{n=1}^N \bar{x}}_{\bar{x} + \sum_{n=1}^N z_n u_1} - x_n \right\|^2$$

“Avg. of original data pts = Avg. of projected data pts” for the opt. u_1 vector. So only vectors for which $\bar{x} = \bar{\tilde{x}}$ need to be considered. So, c is a constant, & minimizing a^2 or maximizing b^2 are the same by Pythagoras thm. ($c^2 = a^2 + b^2$). [CBM]

Outline for Module M5

- M5. Dimensionality Reduction
 - M5.0 Introduction
 - **M5.1 PCA**
 - Intuition for two formulations
 - **Maximizing variance formulation**
 - Minimizing error formulation
 - M5.2 PCA applications/extensions/variants (very brief)

What is the variance of projected data when $M=1$?

$$\text{Var}(\{\mathbf{x}_n^T \mathbf{u}_1\}_{n=1, \dots, N}) = \frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{x}_n^T \mathbf{u}_1 - \bar{\mathbf{x}}^T \mathbf{u}_1 \right\}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

What \mathbf{u}_1 maximizes the above variance?

$$\begin{aligned} &= \frac{1}{N} \sum_n \left((\mathbf{x}_n^T - \bar{\mathbf{x}}^T) \mathbf{u}_1 \right)^2 \\ &= \frac{1}{N} \sum_n \mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_1 \\ &= \mathbf{u}_1^T \left[\frac{1}{N} \sum_n (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \right] \mathbf{u}_1 \end{aligned}$$

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

$$\begin{aligned} (AB)^T &= B^T A^T \\ (\hat{\mathbf{u}}^T \mathbf{v})^2 &= \|\mathbf{u}_v\|^2 = (\mathbf{u}^T \mathbf{v})^2 \\ &= (\mathbf{u}^T \mathbf{v})^2 \end{aligned}$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T.$$

How do you maximize variance, subject to unit length constraint? Lagrangian multiplier

Maximize $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$.

yields: $\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$

with maximum value being: $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$

$$\nabla_x (x^T A x) = (A + A^T) x$$
$$\nabla_x (x^T x) = 2x$$

$$\begin{aligned} & \mathbf{u}_1^T \lambda_1 \mathbf{u}_1 \\ &= \lambda_1 \mathbf{u}_1^T \mathbf{u}_1 = \lambda_1 \end{aligned}$$

Ex: Prove
 S is real, symm.
psd.

$M > 1$, Problem statement for PCA (max. var.)

$\mu, x_1, \dots, x_n \in \mathbb{R}^D$
Find: $u_1, \dots, u_M \in \mathbb{R}^D$
orthonorm.

$$u_i^T u_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

s.t. $J = \sum_{i=1}^M \text{Var}(x_i^T u_i) = \sum_{i=1}^M u_i^T S u_i$ is maximized.

Global max of J achieved when $\{u_i\}$ are top M eigen vecs. of S with max value being $\sum_{i=1}^M \lambda_i$.

$M > 1$ (alternate) proof sketch

One proof very similar to that of spectral clustering M score vectors' proof.

Alternate proof sketch: It can also be shown that

- Q: What is the direction u_2 that maximizes variance of projected data along u_2 , under the constraint that u_2 is orthogonal to the top eigenvector u_1 ?
- A: Eigenvector corresp. to the 2nd largest eigenvalue.
- Extending the argument by induction shows that “top” M eigen vectors of S are the “top” M PCs, and projection onto them maximizes the sum of variances of projected data along these directions.

M > 1 proof (similar to spectral clustering M score vectors' proof)

Maximize $\tilde{J} = \text{Tr} \left\{ \hat{\mathbf{U}}^T \mathbf{S} \hat{\mathbf{U}} \right\} + \text{Tr} \left\{ \mathbf{H} (\mathbf{I} - \hat{\mathbf{U}}^T \hat{\mathbf{U}}) \right\}$

$$\tilde{J} = \sum_{i=1}^M \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i + \sum_{i,j=1}^M H_{ij} (\delta_{ij} - \mathbf{u}_i^T \mathbf{u}_j)$$

$$\hat{\mathbf{U}}_{N \times M} = \begin{bmatrix} | & & | \\ \mathbf{u}_1 & & \mathbf{u}_M \\ | & & | \end{bmatrix}$$

yields: $\mathbf{S} \hat{\mathbf{U}} = \hat{\mathbf{U}} \mathbf{H}$

One soln: $\mathbf{S} \mathbf{u}_i = \lambda_i \mathbf{u}_i$ for $i = 1, \dots, M$,

with maximum value being: $\tilde{J} = \sum_{i=1}^M \lambda_i$.

(enough to consider this solution alone, as the symmetric \mathbf{H} can be assumed to be diagonal ($\mathbf{H}_{ii} := \lambda_i$) wlog)

Outline for Module M5

- M5. Dimensionality Reduction
 - M5.0 Introduction
 - **M5.1 PCA**
 - Intuition for two formulations
 - Maximizing variance formulation
 - **Minimizing error formulation**
 - M5.2 PCA applications/extensions/variants (very brief)

Representing a datapoint (and its approxmn.) via a set of D ($M < D$) orthonormal vectors in \mathbb{R}^D

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}.$$

(orthonormality requirement) u_1, \dots, u_M

$$\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i \quad \mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i.$$

(perfect/loss-free reconstruction using D dimensions, i.e., along D u_i s)

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

(approximation using $M < D$ dimensions;
Note: b_i does **not** depend on n)

x_1, \dots, x_N ($N \times D$) \mapsto $\tilde{x}_1, \dots, \tilde{x}_N$ ($N \times D$) $\{u_i\}^D$ scalars $\{z_{ni}\}^{N \times M}$ $\{b_i\}^{D-M}$ scalars

Problem statement for PCA (min. error)

1/p: $\{x_n\}_{n=1}^N \in \mathbb{R}^D$

o/p: Find $\{z_{ni}\}$, $\{b_i\}$ s.t. $\{u_i\}$ orthonorm.

s.t. $\frac{1}{N} \sum_n \|x_n - \tilde{x}_n\|^2$ is minimized;

where

$$\tilde{x}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

Objective #1: choosing z_{ni} and b_i to minimize error (given an orthonormal set $\{u_i\}$)

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2.$$

$$= \frac{1}{N} \sum_n \left[\sum_{i=1}^M (\alpha_{ni} - z_{ni})^2 + \sum_{i=M+1}^D (\alpha_{ni} - b_i)^2 \right]$$

minimized
when

$$z_{ni} = \mathbf{x}_n^T \mathbf{u}_i \quad b_i = \bar{\mathbf{x}}^T \mathbf{u}_i$$

Minimum value:

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i.$$

Objective #2: choosing $\{u_i\}$ to minimize error
(when $M=1$ & $D=2$ again)

Minimize $\tilde{J} = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2 + \lambda_2 (1 - \mathbf{u}_2^T \mathbf{u}_2)$

yields: $\mathbf{S} \mathbf{u}_2 = \lambda_2 \mathbf{u}_2$

with minimum value being: $\mathbf{u}_2^T \mathbf{S} \mathbf{u}_2 = \lambda_2$

General D,M (i.e., when $(D-M) > 1$, proof similar to spectral clustering multiple score vectors' proof)

Minimize $\tilde{J} = \text{Tr} \{ \hat{\mathbf{U}}^T \mathbf{S} \hat{\mathbf{U}} \} + \text{Tr} \{ \mathbf{H}(\mathbf{I} - \hat{\mathbf{U}}^T \hat{\mathbf{U}}) \}$

$$\tilde{J} = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i + \sum_{i,j=M+1}^D H_{ij} (\delta_{ij} - \mathbf{u}_i^T \mathbf{u}_j)$$

$H_{(i-M)(j-M)}$

$\hat{\mathbf{U}} \in \mathbb{R}^{D \times (D-M)}$

$$\hat{\mathbf{U}} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{u}_{M+1} & \dots & \mathbf{u}_D \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

yields: $\mathbf{S} \hat{\mathbf{U}} = \hat{\mathbf{U}} \mathbf{H}$

One soln: $\mathbf{S} \mathbf{u}_i = \lambda_i \mathbf{u}_i$ for $i = M + 1, \dots, D$,

with minimum value being: $\tilde{J} = \sum_{i=M+1}^D \lambda_i$.

(enough to consider this solution alone, as the symmetric \mathbf{H} can be assumed to be diagonal ($\mathbf{H}_{ii} := \lambda_i$) wlog)

Outline for Module M5

- M5. Dimensionality Reduction

- M5.0 Introduction

- M5.1 PCA

- Intuition for two formulations
 - Maximizing variance formulation
 - Minimizing error formulation

- **M5.2 PCA applications/extensions/variants (very brief)**

$\{x_n\}$ $X \in \mathbb{R}^{N \times D}$ (mean-centered) $\rightarrow S = \frac{X^T X}{N}$ $\xrightarrow{\text{EVD}}$ Top M EV (PCs) $\rightarrow \{\tilde{x}_n\}$ ($N \times M$)

$(N \times M, D^2, D-M)$

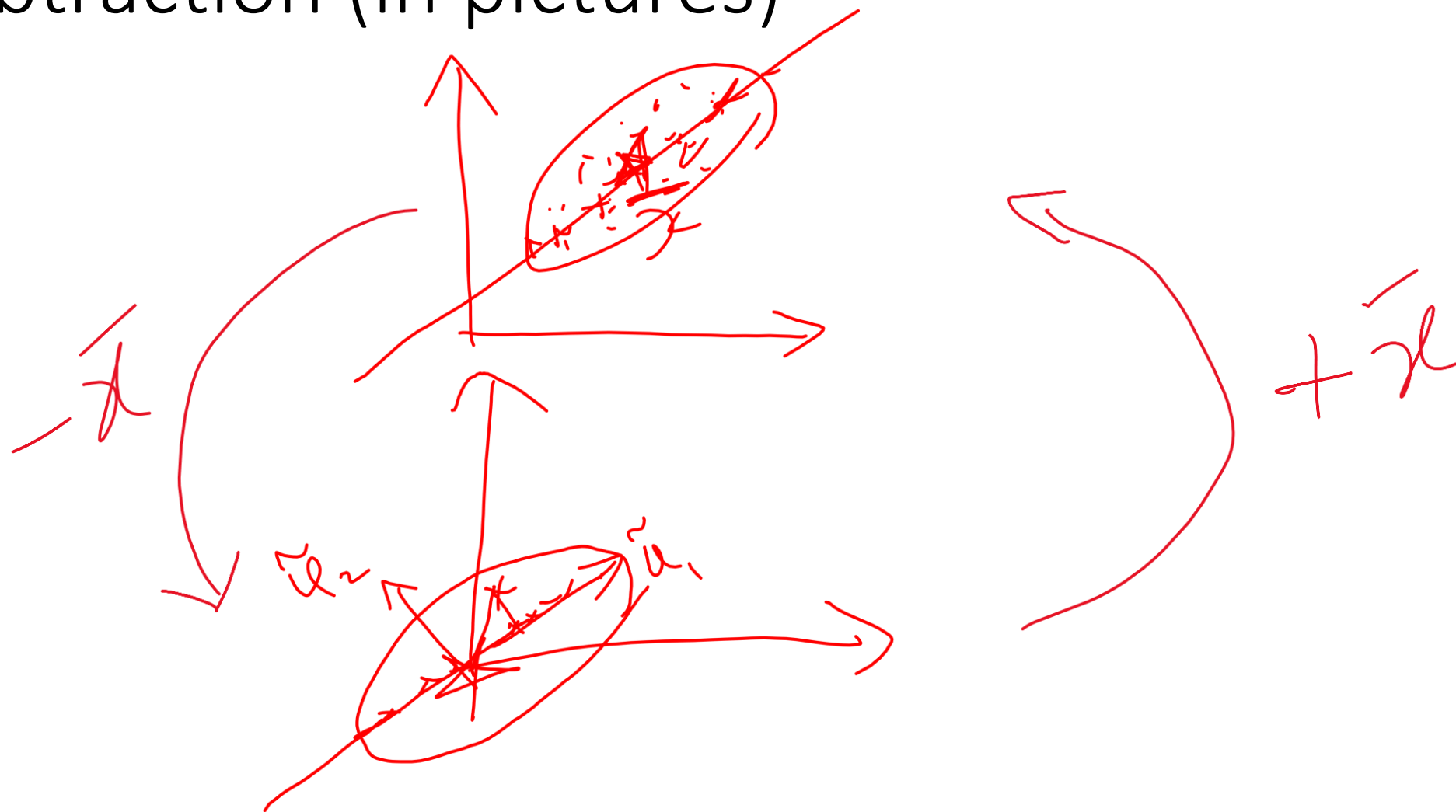
$$S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T.$$

Data compression: but first, what is all the talk about mean-subtracting?

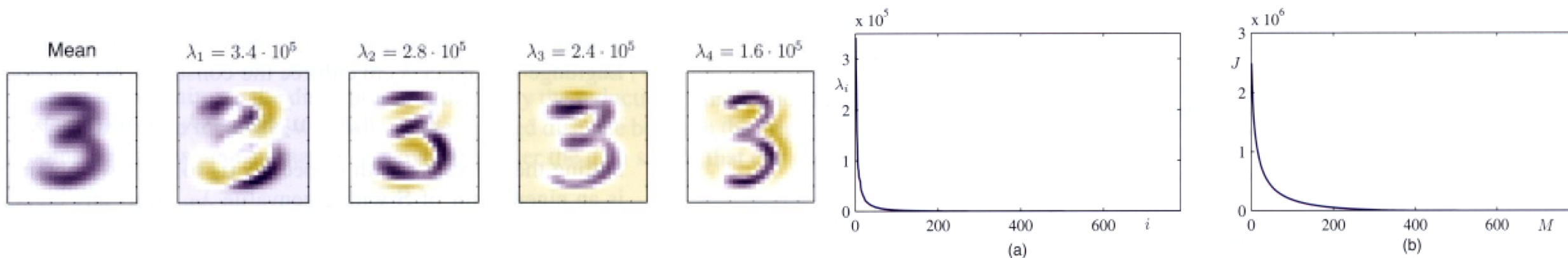
$$\begin{aligned}\tilde{\mathbf{x}}_n &= \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \quad + \bar{\mathbf{x}} - \bar{\mathbf{x}} \\ &= \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i = \bar{\mathbf{x}} + \sum_{i=1}^M ((\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i) \mathbf{u}_i\end{aligned}$$

$$\bar{\mathbf{x}} = \sum_{i=1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$$

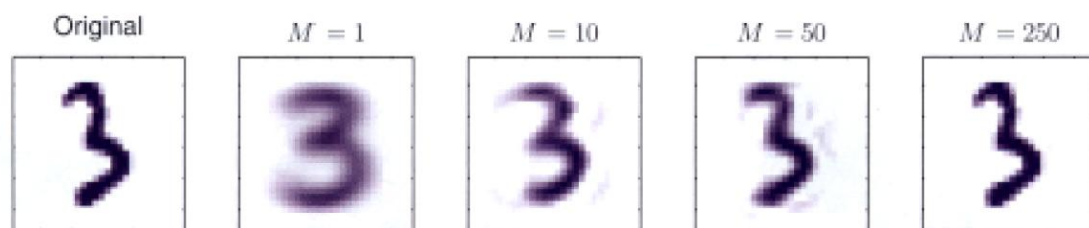
Mean-subtraction (in pictures)



Data compression (on mean-subtracted data):



ND
 $x_n \in \mathbb{R}^{9 \times 10}$




MM

$\sum_{i=1}^M \lambda_i^2$

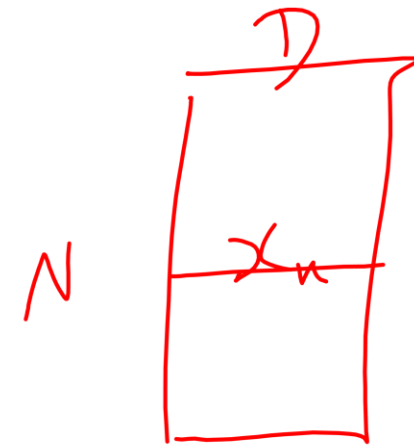
Other PCA extensions/variants

- PCA for High-dimensional Data
 - (transpose, compute Eigen vectors in less time, convert to Eigen vectors in original space – has connections to SVD) – see Assignment question!
- From Deterministic PCA to Probabilistic PCA to Bayesian PCA / Mixture of Probabilistic PCA / Factor analysis / etc.
- Kernel PCA



A handwritten diagram showing a rectangular box representing a matrix. To the left of the box is the letter 'N'. Above the box is the letter 'D'. Inside the box, the text 'hi-dim.' is written. Below the box, the expression $O(N^3)$ is written. To the right of $O(N^3)$ is the equation $\tilde{S} = \frac{1}{N} X X^T$.

$$O(N^3) \quad \tilde{S} = \frac{1}{N} X X^T$$



A handwritten diagram showing a rectangular box representing a matrix. To the left of the box is the letter 'N'. Above the box is the letter 'D'. Inside the box, the text 'x_n' is written. To the right of the box is the letter 'S'. Below the box is the letter 'X'. To the right of 'S' is the equation $S = \frac{1}{N} X X^T$. Below 'X' is the text '= mean. sub.'. To the right of this text is the expression $O(D^3)$.

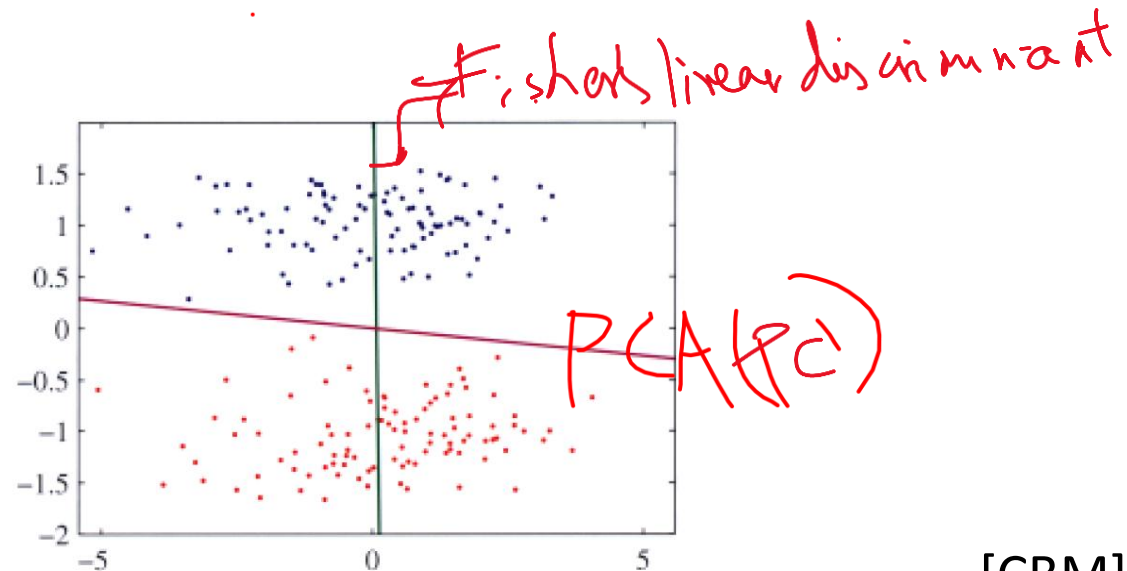
$$S = \frac{1}{N} X X^T$$

$X = \text{mean. sub. } O(D^3)$

Summary of dimensionality reduction (PCA)

- PCA is a method to transform data points from a high-dimensional to a linear low-dimensional space, by maximizing variance or equiv. minimizing error.
 - Linear because low-dimensional subspace is a vector subspace (provided, the data is mean-centered i.e., the mean passes through zero making the subspace to include zero/origin as well; even otherwise, PCA holds by viewing the subspace as an affine subspace instead of a vector subspace)
- One of the many available unsupervised methods for dimension reduction.

- Supervised dim. redn. methods also exist (not covered here) and may be better at finding relevant dimensions.



Thank you!