

Dimensionality Reduction

Dr. Nimisha Roy
Georgia Tech

Outline

- Overview 
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

Motivating Example: Data Visualization

53 blood and urine samples
(features) from 65 people

- Matrix format (65x53)

	H-WBC	H-RBC	H-Hgb	H-Hct	H-MCV	H-MCH	H-MCHC
A1	8.0000	4.8200	14.1000	41.0000	85.0000	29.0000	34.0000
A2	7.3000	5.0200	14.7000	43.0000	86.0000	29.0000	34.0000
A3	4.3000	4.4800	14.1000	41.0000	91.0000	32.0000	35.0000
A4	7.5000	4.4700	14.9000	45.0000	101.0000	33.0000	33.0000
A5	7.3000	5.5200	15.4000	46.0000	84.0000	28.0000	33.0000
A6	6.9000	4.8600	16.0000	47.0000	97.0000	33.0000	34.0000
A7	7.8000	4.6800	14.7000	43.0000	92.0000	31.0000	34.0000
A8	8.6000	4.8200	15.8000	42.0000	88.0000	33.0000	37.0000
A9	5.1000	4.7100	14.0000	43.0000	92.0000	30.0000	32.0000

Difficult to see the correlations of different features

Motivating Example: Data Visualization

Is there a representation better than the coordinate axes?

Is it really necessary to show all the 53 dimensions?

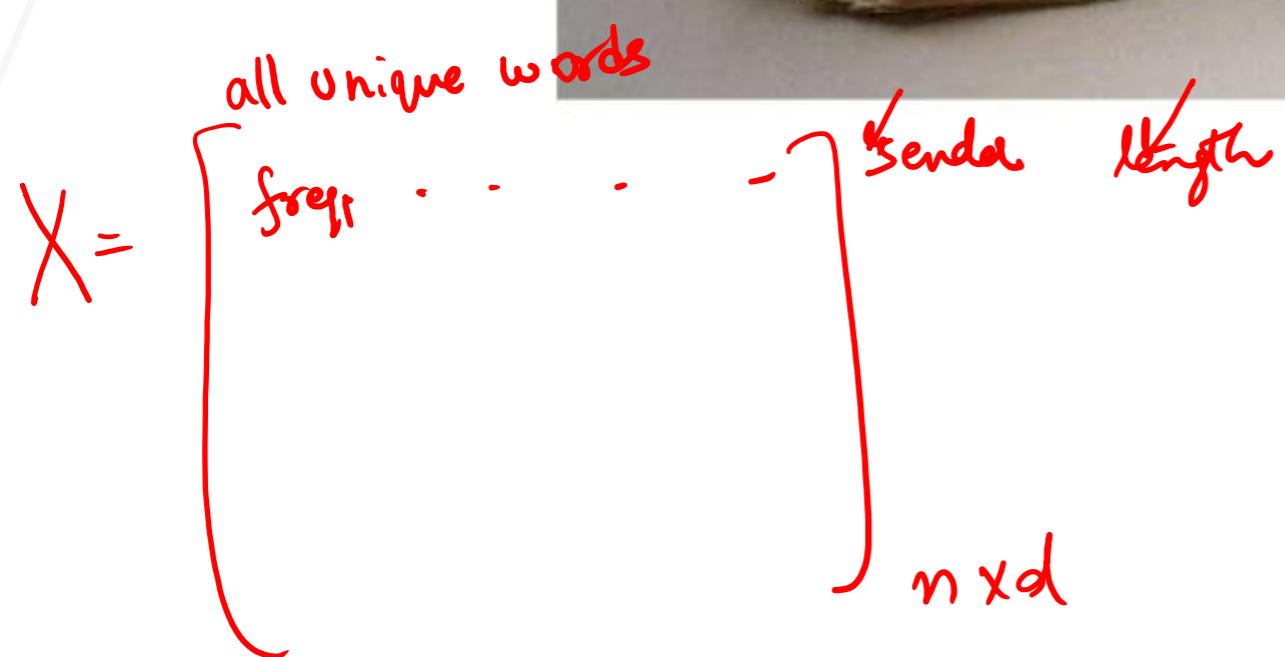
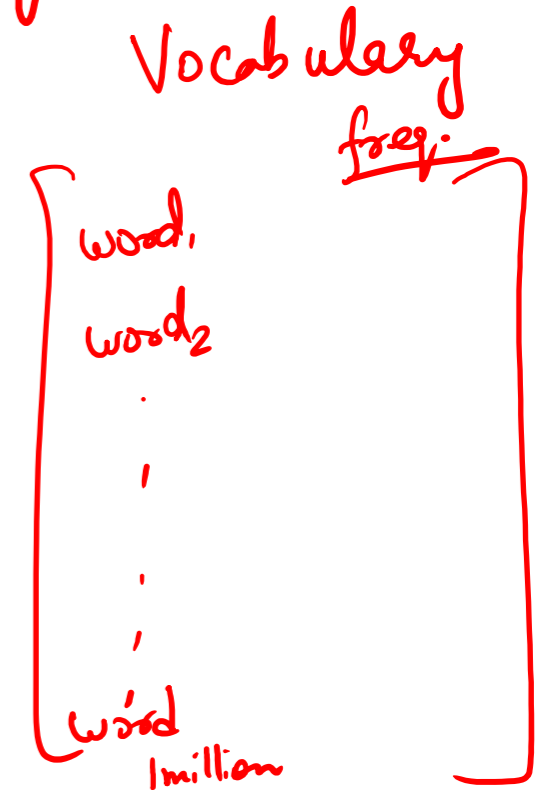
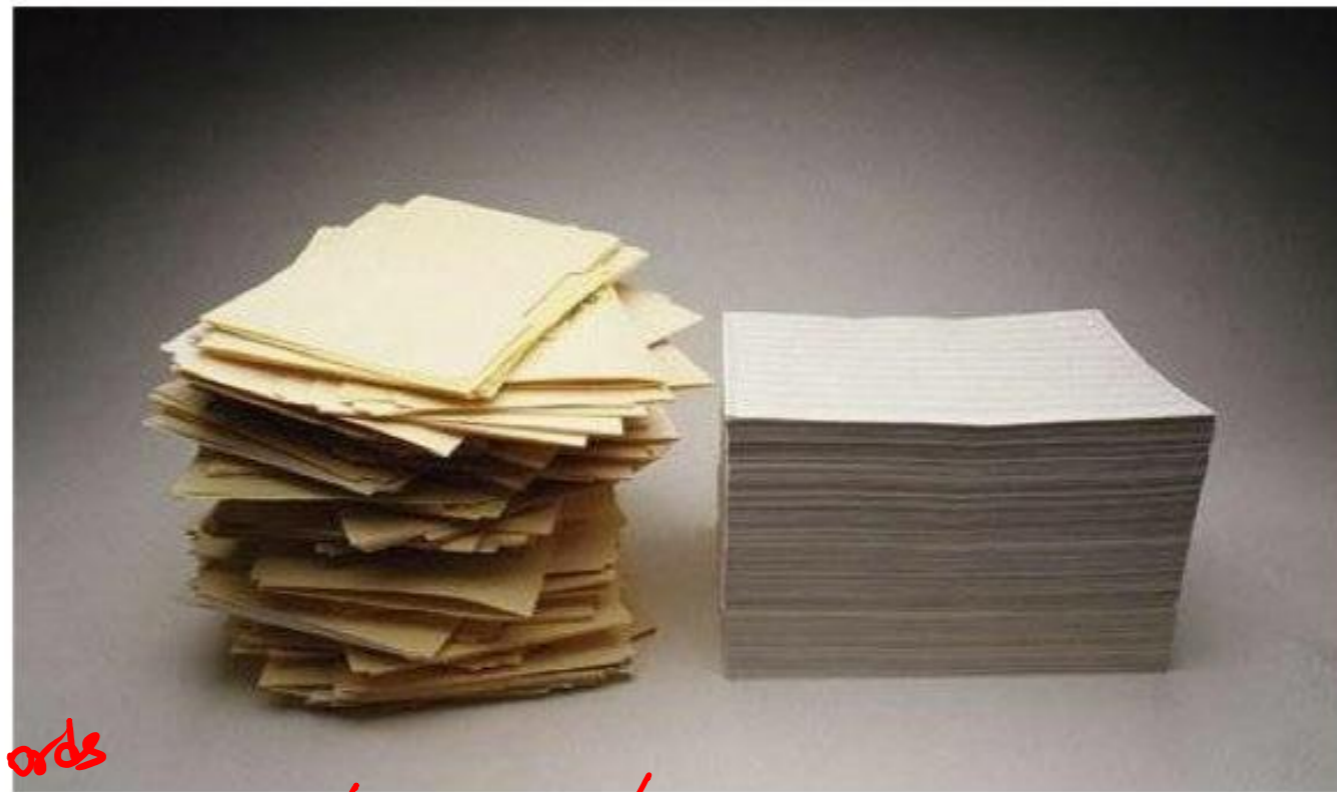
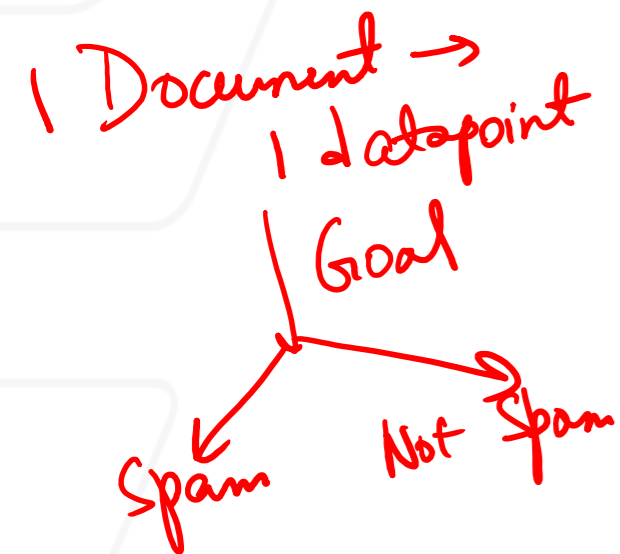
- ... what if there are strong correlations between the features?

How could we find the *smallest* subspace of the 53-D space that keeps the *most information* about the original data?

A Solution: Dimension Reduction

Another Example: Dimension Reduction for Text

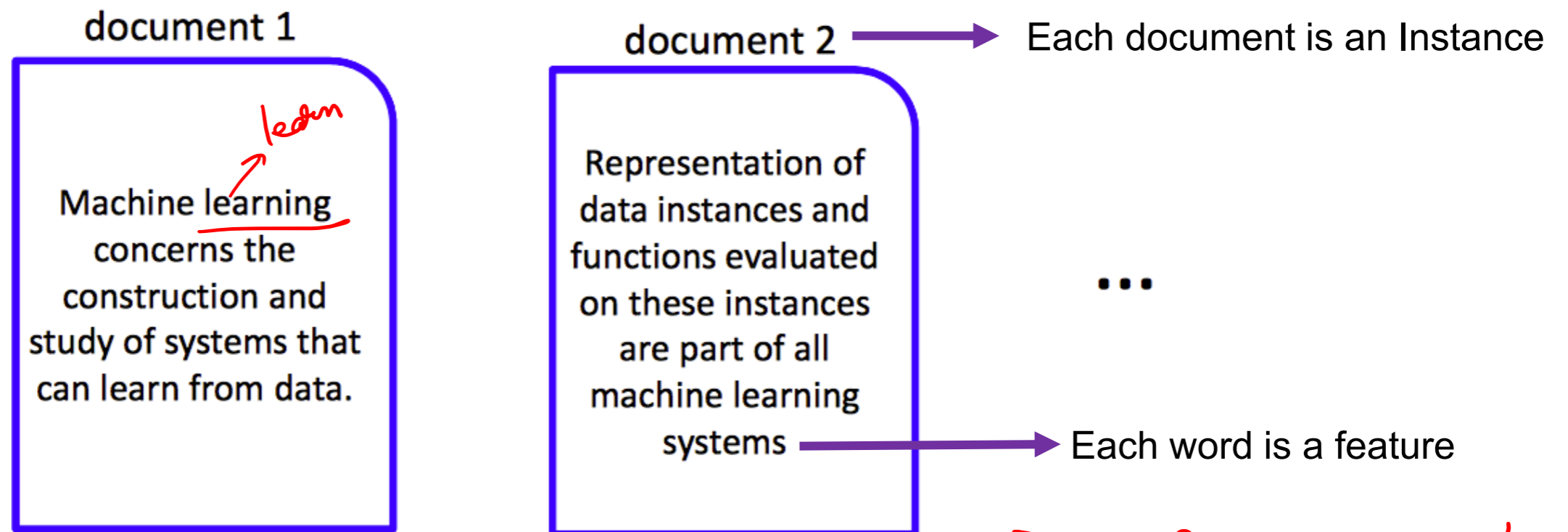
Preprocessing for textual data is dimensionality reduction



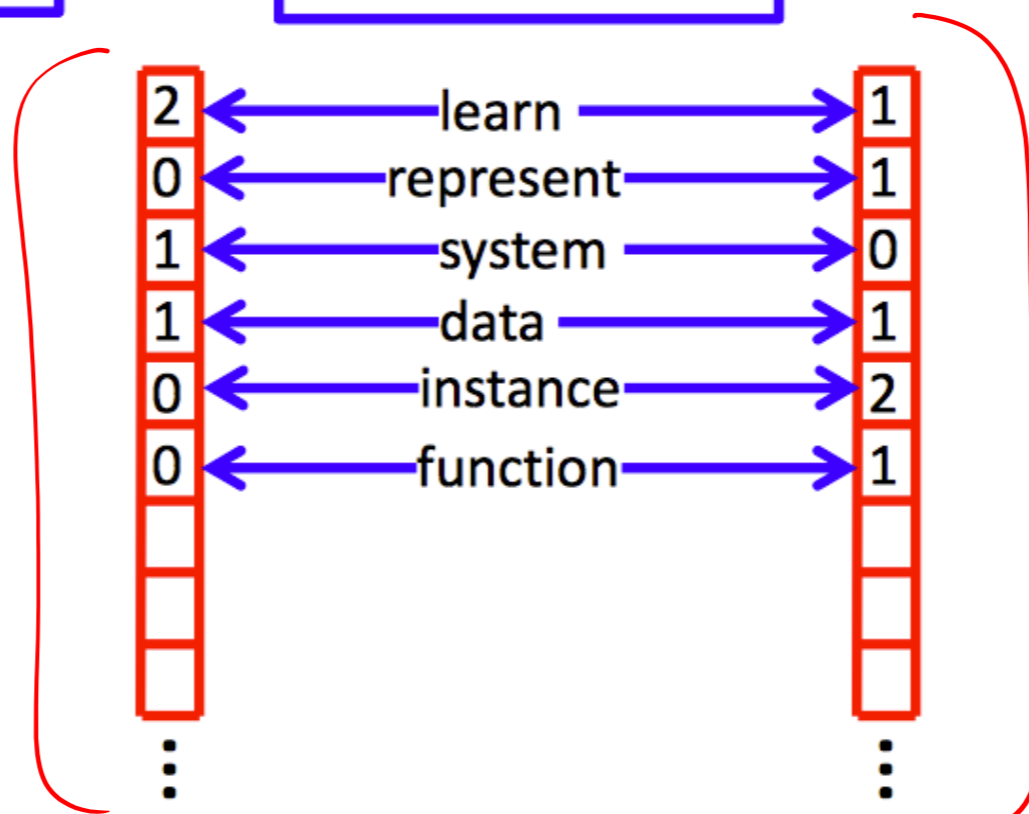
Bag-of-Words Representations

Remove stop words - a, an, the, ...

Lemmatization -



Vector in R^d



Bag of Words

Term-Document Data Matrix – Bag-of-words

database

	database	SQL	index	regression	likelihood	linear
d1	24	21	9	0	0	3
d2	32	10	5	0	3	0
d3	12	16	5	0	0	0
d4	6	7	2	0	0	0
d5	43	31	20	0	3	0
d6	2	0	0	18	7	16
d7	0	0	1	32	12	0
d8	3	0	0	22	4	2
d9	1	0	0	34	27	25
d10	6	0	0	17	4	23

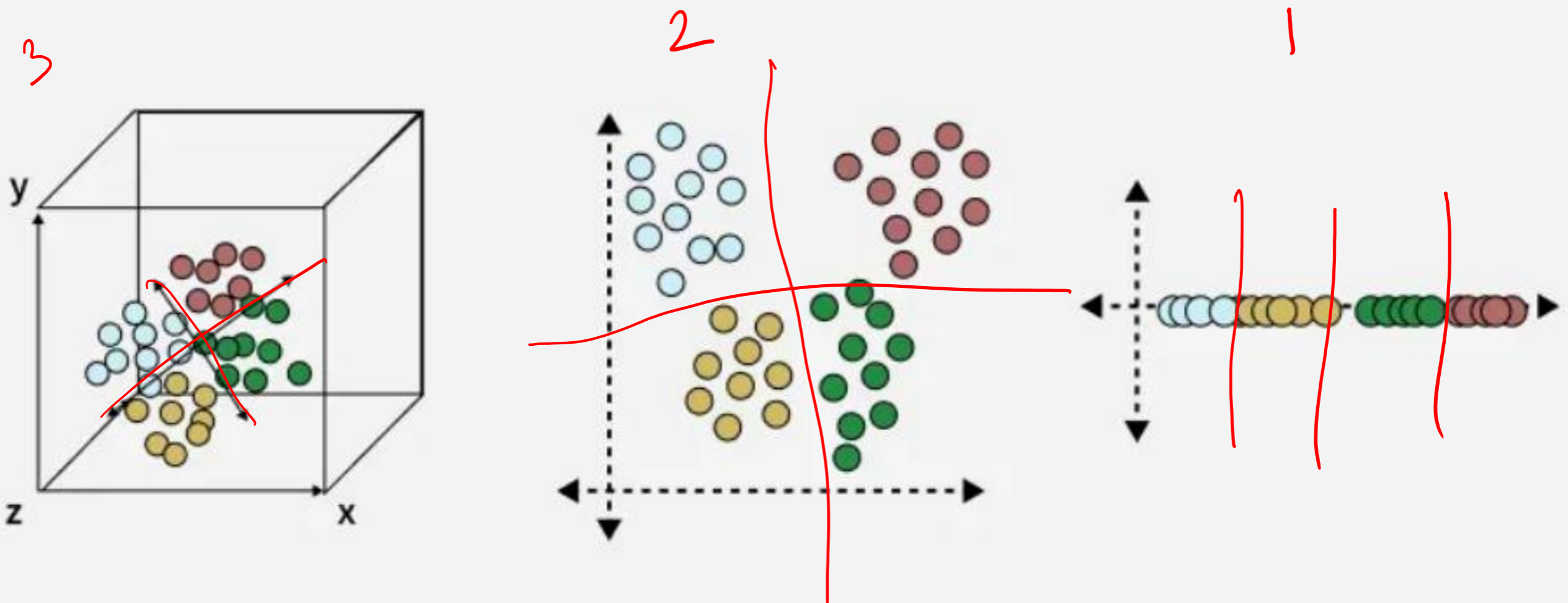
▪▪▪ Many more features

ML core algo

Solution:
**Dimension
Reduction**

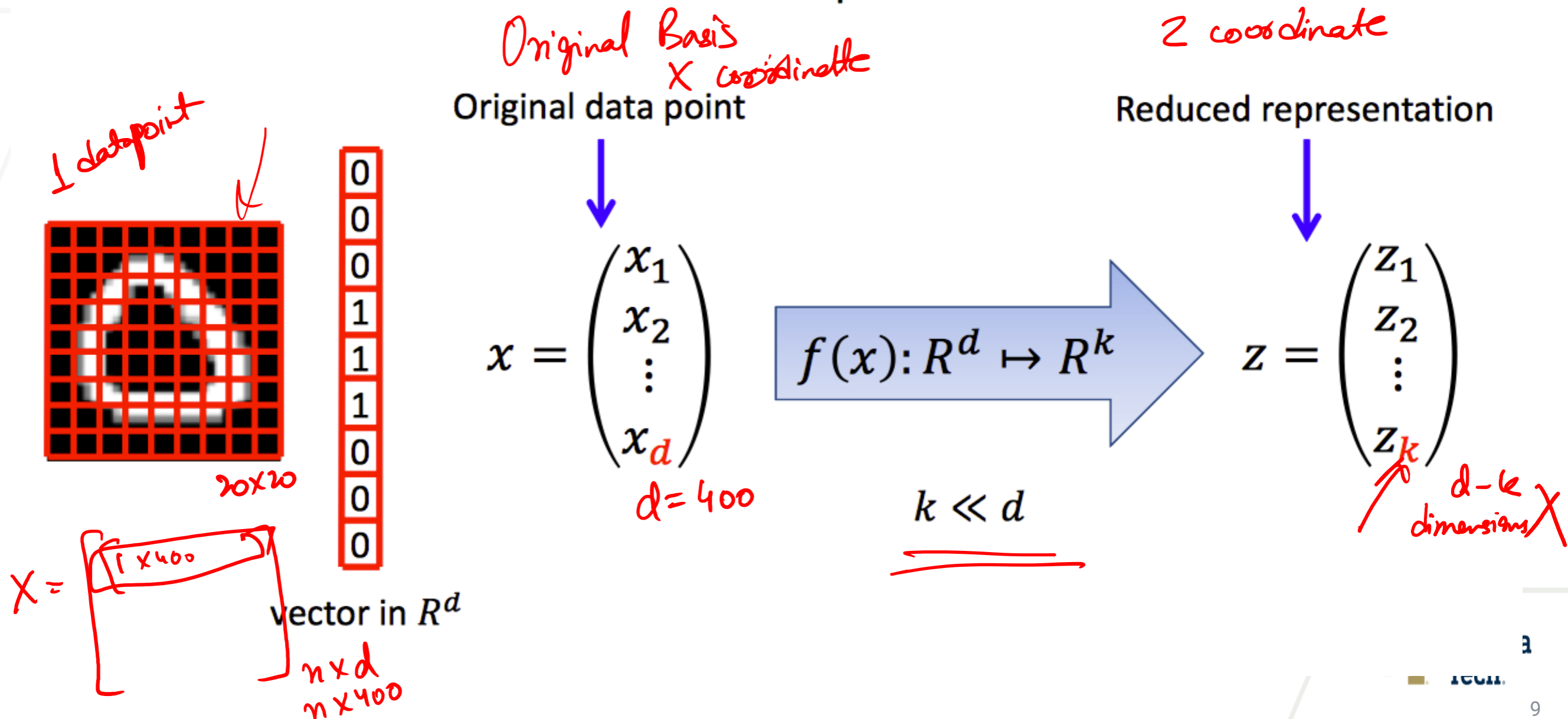
What is Dimensionality Reduction

Dimensionality Reduction is the process of reducing the number of input variables (features) in a dataset while preserving as much important information as possible.



What is Dimension Reduction?

- The process of reducing the number of random variables under consideration
 - One can combine, transform or select variables
 - One can use linear or nonlinear operations



Applications of Dimension Reduction

- The dimension-reduced data can be used for
 - Visualizing, exploring and understanding the data EDA
 - Aggregating weak signals in the data *features*
 - Cleaning the data
 - Speeding up subsequent learning task
 - Building simpler model later
- Key questions of a dimensionality reduction algorithm
 - What is the criterion for carrying out the reduction process?
 - What are the algorithm steps?

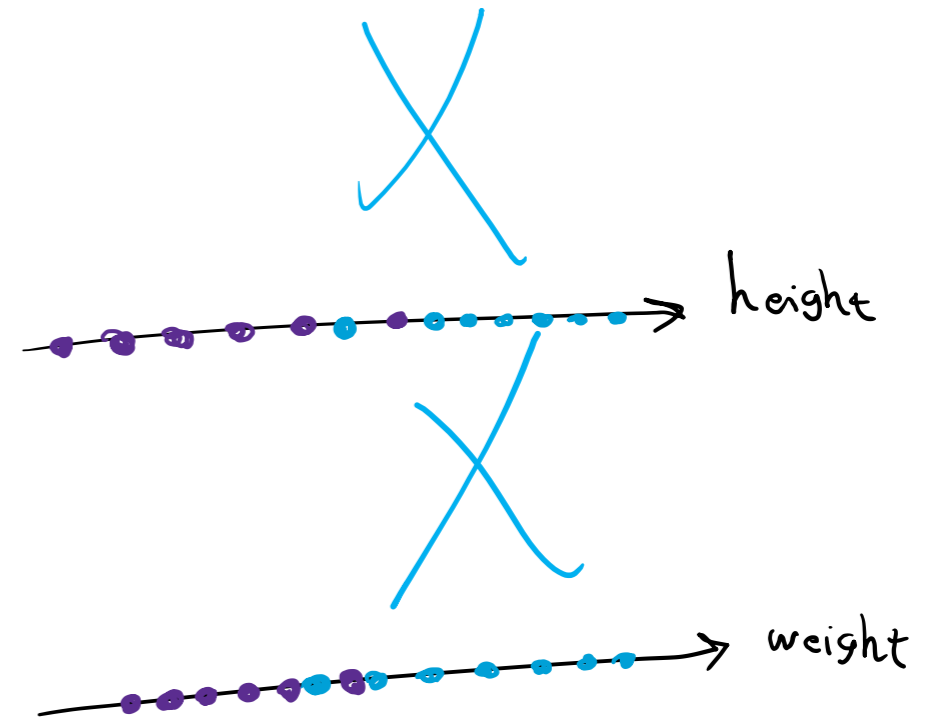
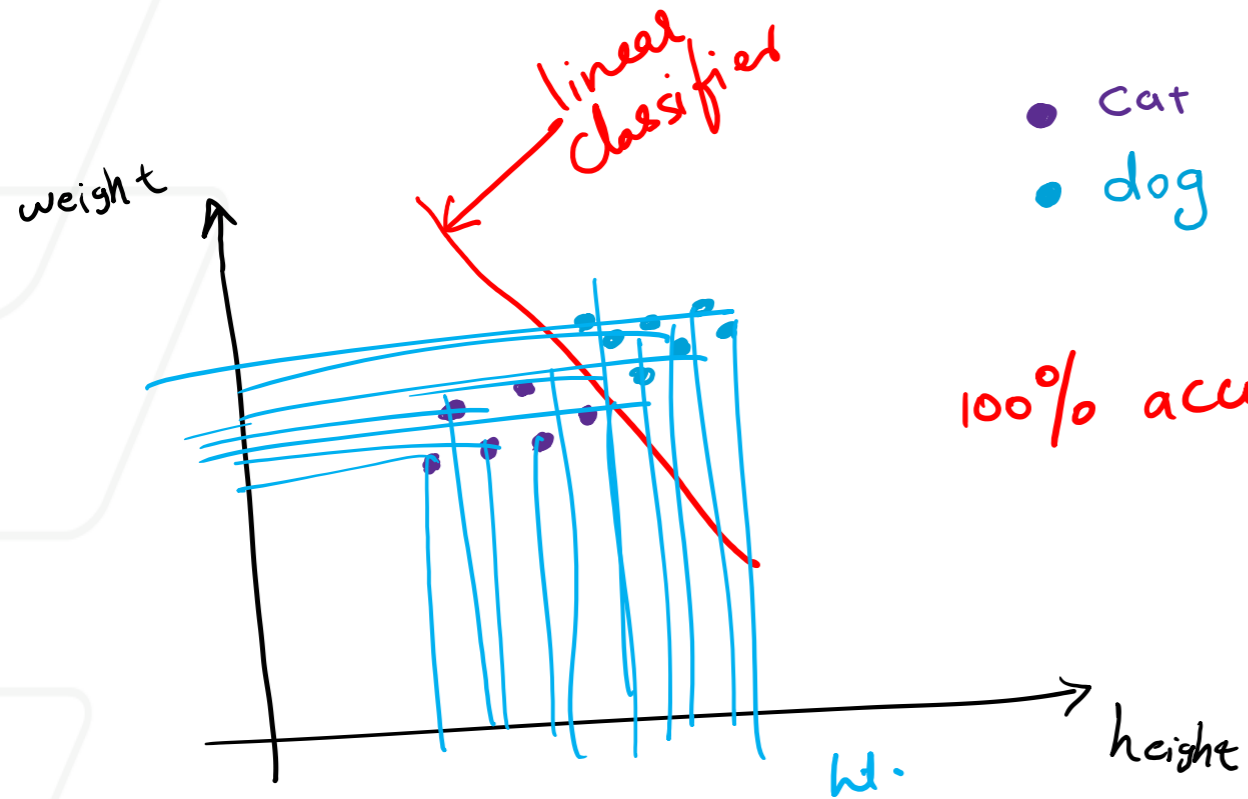
Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

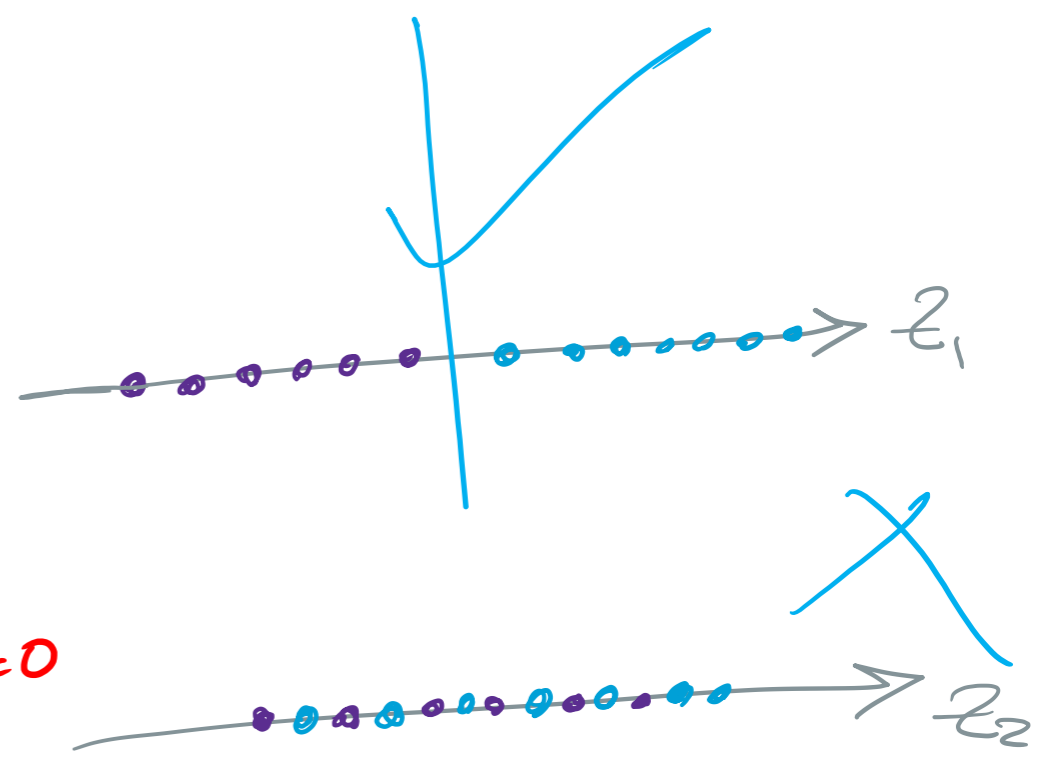
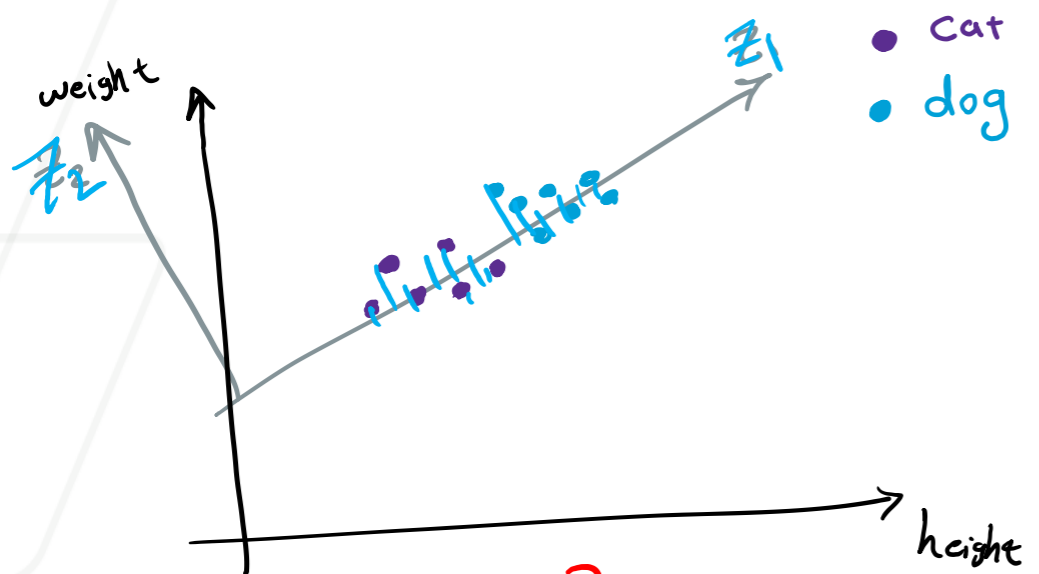


PCA: Dimension Reduction by Capturing Variation

- There are many criteria (geometric based, information theory based, etc.)
- One criterion: want to capture variation in data
 - variations are “signals” or information in the data
 - need to normalize each variables first
- In the process, also discover variables or dimensions highly **correlated**
 - represent highly related phenomena
 - combine them to form a stronger signal
 - lead to simpler presentation



Dropping / Truncating dimensions is not good



Normalization $\rightarrow [0, 1]$

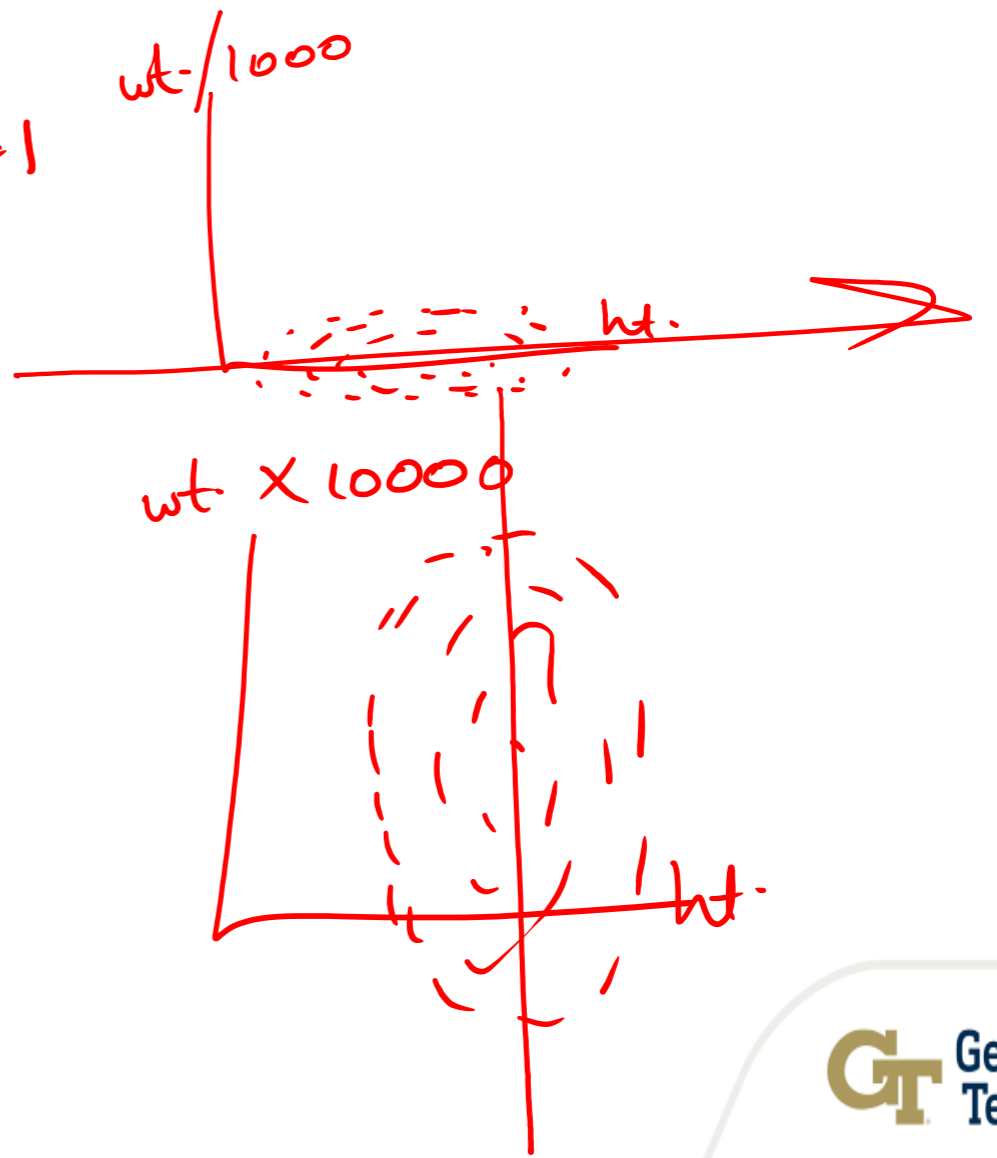
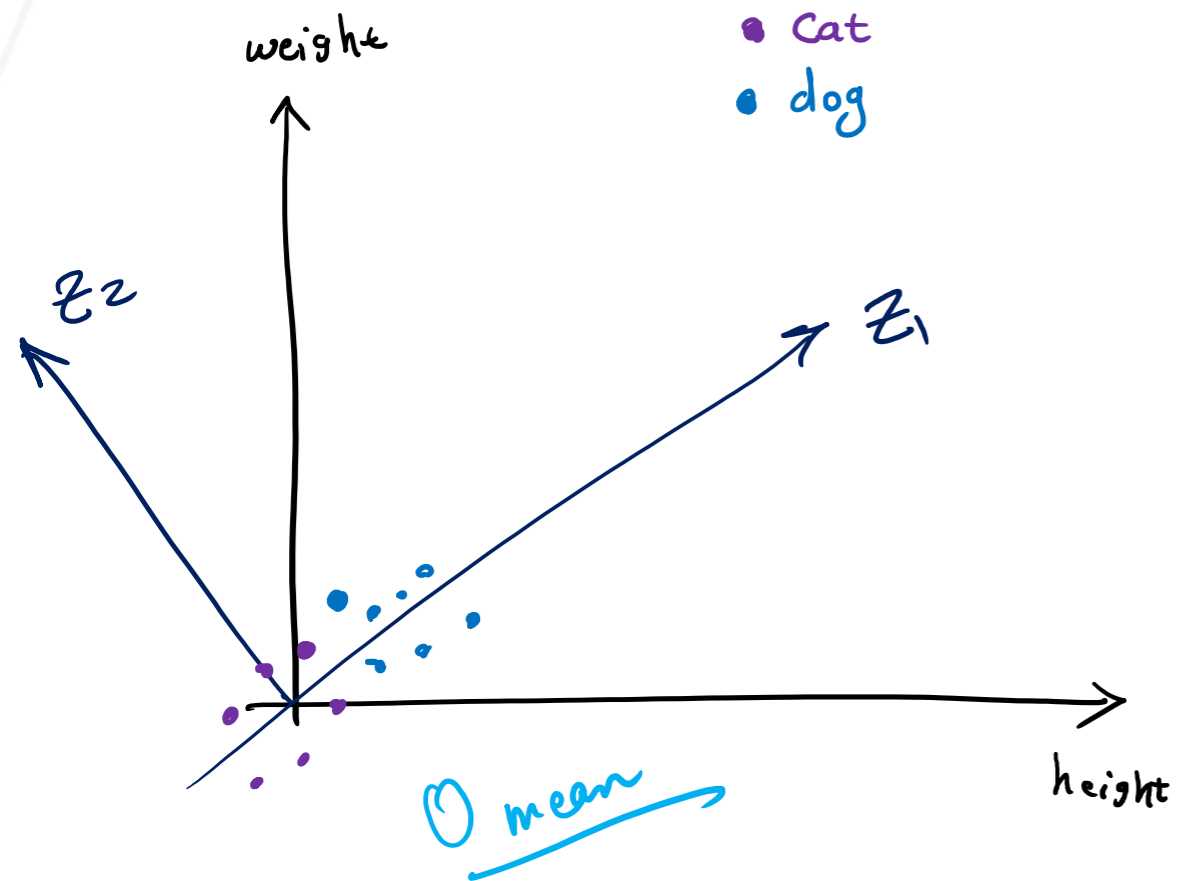
$$\frac{x - x_{min}}{x_{max} - x_{min}}$$

Standardization

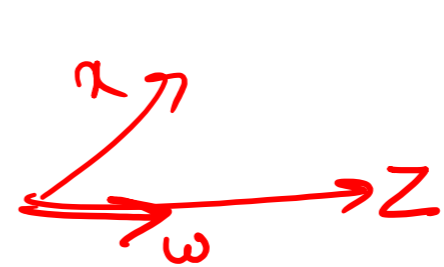
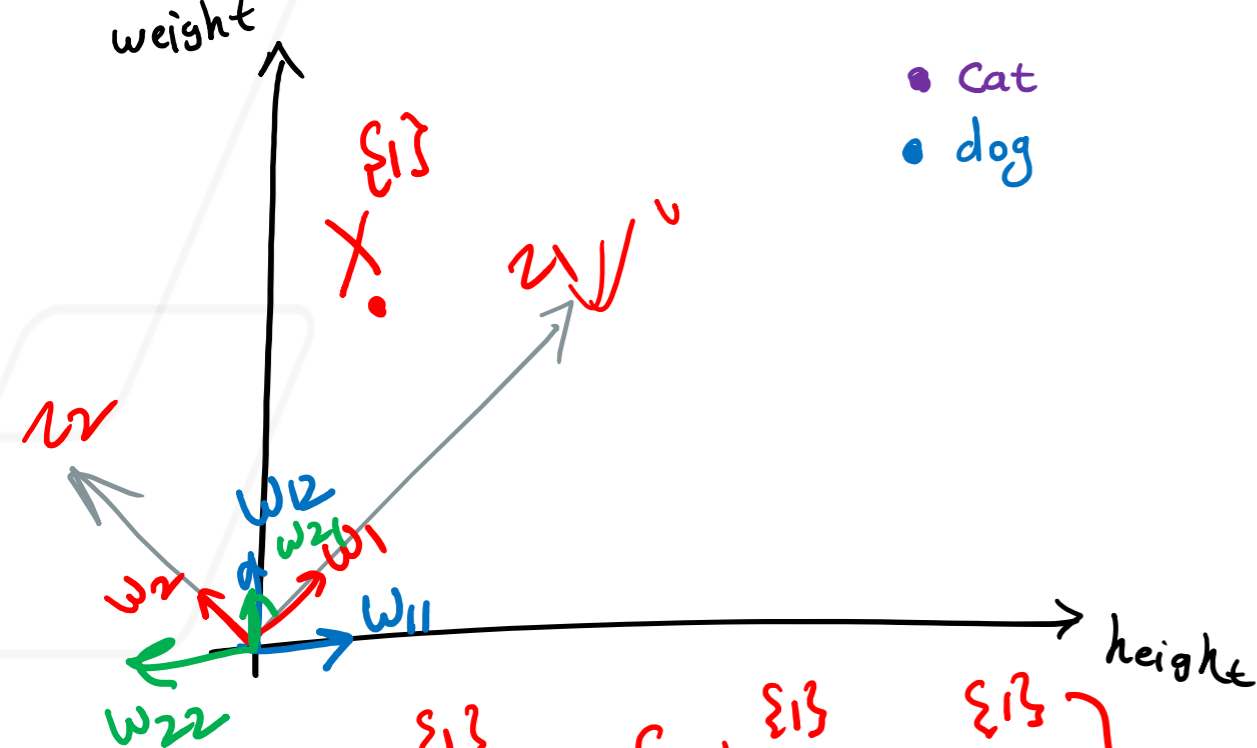
$$\frac{x - \mu}{\sigma}$$

mean = 0

Std. = 1



Projection



$$x \cdot w = x \cdot \frac{z}{\|z\|_2}$$

$x^{xi13} = \begin{bmatrix} h^{xi13} & e^{xi13} \end{bmatrix}$ → move to z space

unit vectors in z space is $\begin{bmatrix} w_1 & w_2 \end{bmatrix}$
 in x space = $\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$

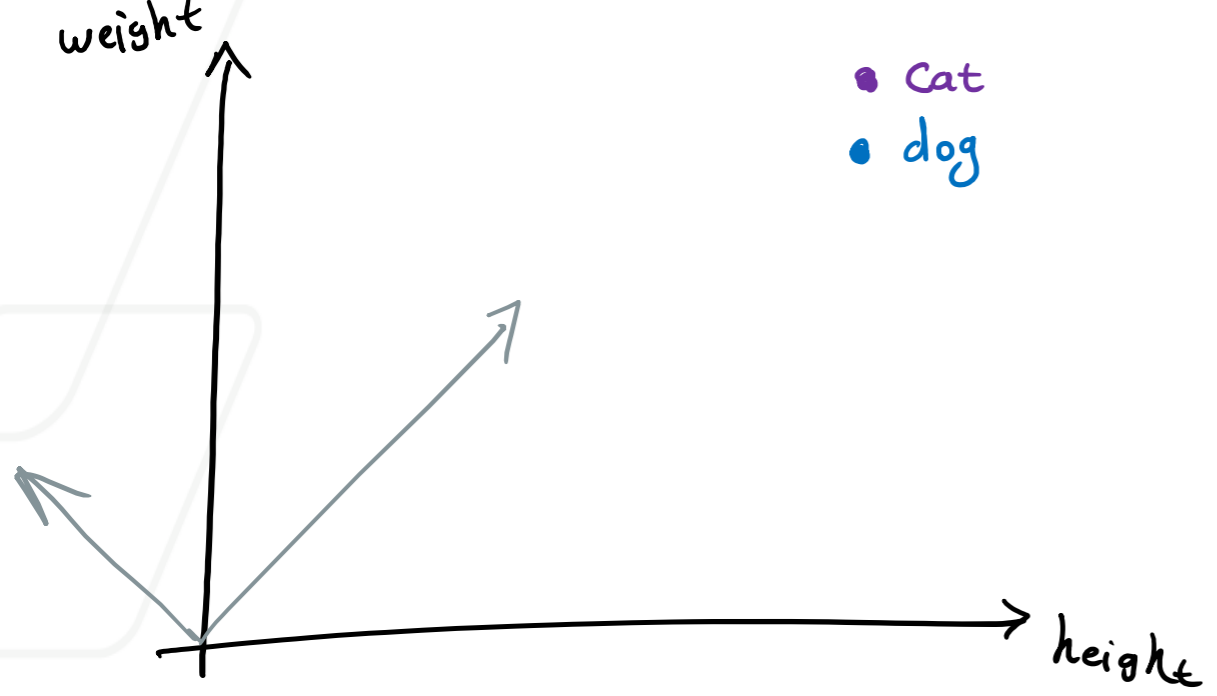
$$z_1^{xi13} = x^{xi13} \cdot w_1 = \begin{bmatrix} h^{xi13} \\ e^{xi13} \end{bmatrix} \begin{bmatrix} w_{11} & w_{12} \end{bmatrix} = h^{xi13} w_{11} + e^{xi13} w_{12}$$

$$z_2^{xi13} = x^{xi13} \cdot w_2 = \begin{bmatrix} h^{xi13} \\ e^{xi13} \end{bmatrix} \begin{bmatrix} w_{21} & w_{22} \end{bmatrix} = h^{xi13} w_{21} + e^{xi13} w_{22}$$

$$\begin{bmatrix} h^{xi13} & e^{xi13} \end{bmatrix} \xrightarrow{\text{transformed to z space}} \begin{bmatrix} h^{xi13} w_{11} + e^{xi13} w_{12} & h^{xi13} w_{21} + e^{xi13} w_{22} \end{bmatrix}$$

Standardization

- Cat
- dog



~~X~~

Standardize



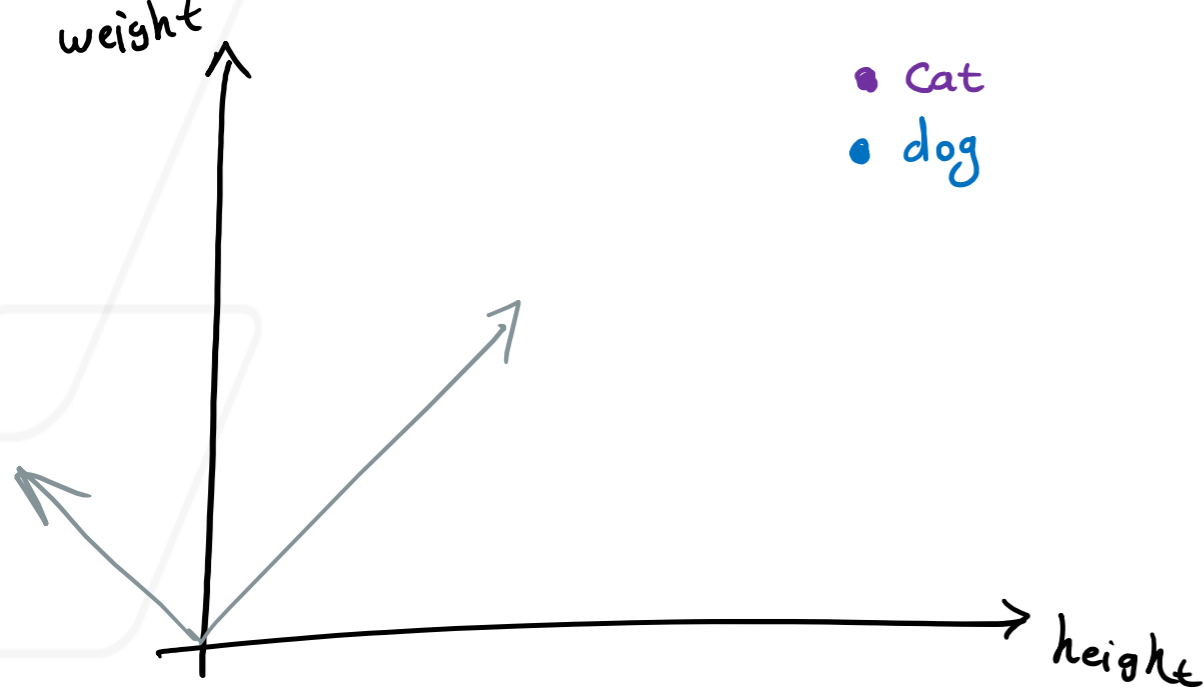
$$\frac{\overset{*}{X}}{\sigma} = \frac{X - \mu}{\sigma}$$

$$\frac{\overset{*}{X}}{\sigma} \cdot \overset{*}{\omega}$$

Truncate features

ω

Objective Function



$$\text{Var}(x) = \frac{1}{N} \sum_i (x_i - \mu)^2$$

Maximize variance in the
z space

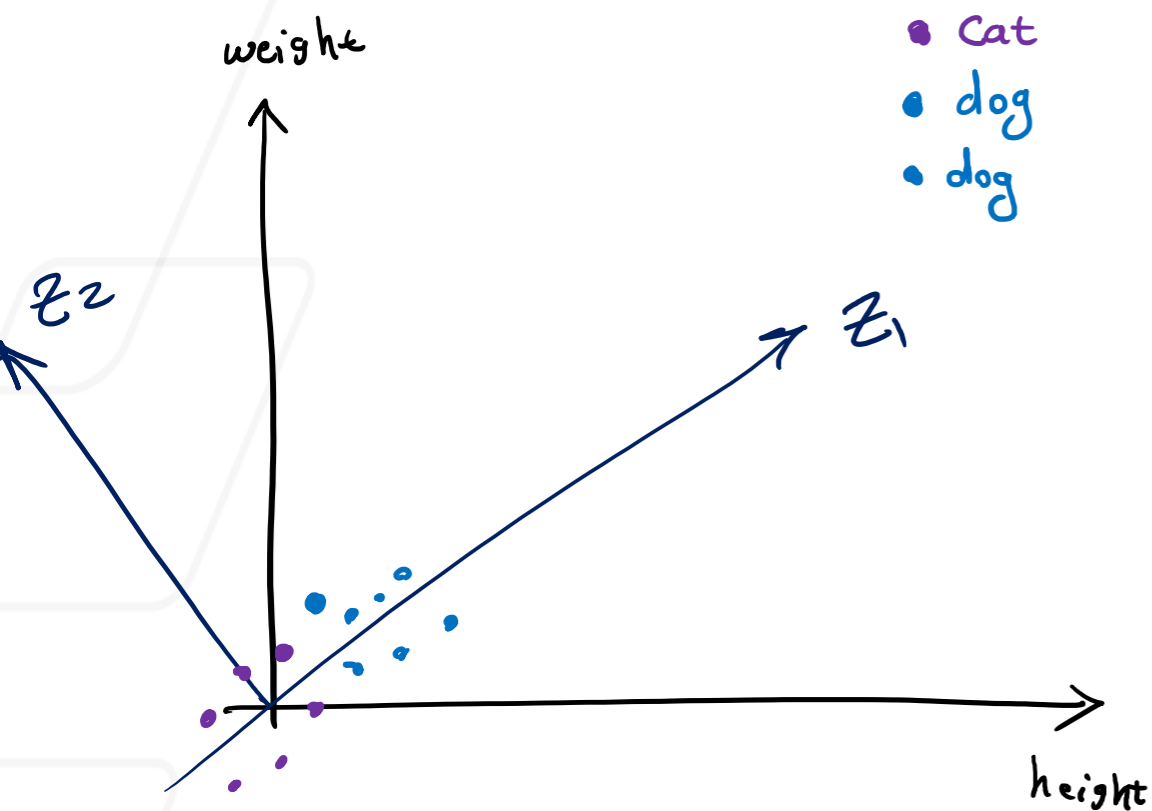
max.
$$\text{Var}(w) = \frac{1}{N} \sum_i (x_i \cdot w - \mu \cdot w)^2$$

s.t. $\|w\|_2 = 1$

Objective Function

$$L(w) = f(w) - \lambda g(w)$$

$\frac{\partial L}{\partial w} = 0 \rightarrow$ find that w



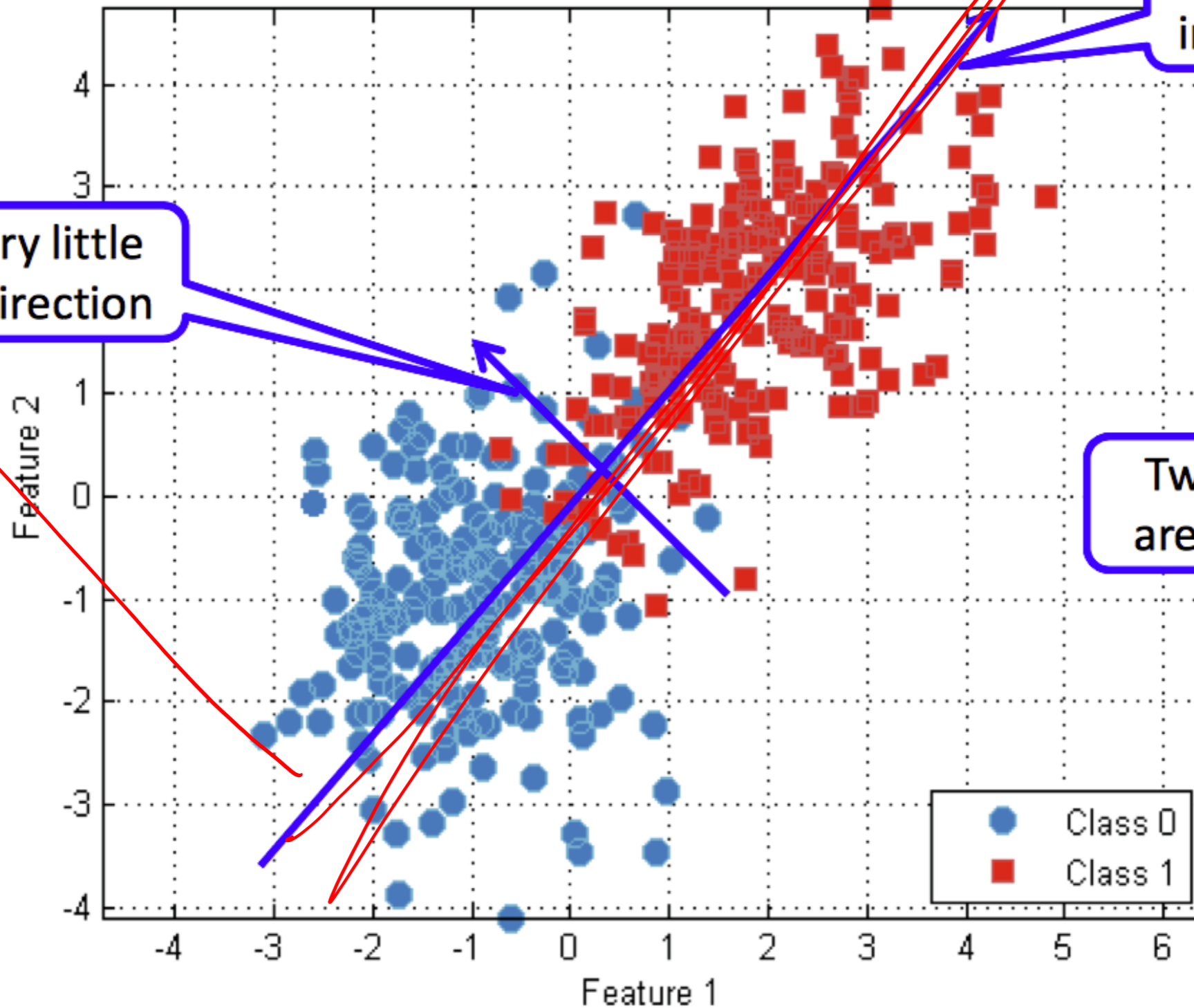
Yes

Were ht. and wt. features correlated?

So what does PCA do?

- **Combines correlated features** into a single axis that captures their shared information.
- The new direction z is aggregation of existing features. Aggregation of weak signals.
- **Reduces dimensionality** by keeping only the most informative directions — leading to simpler data representation while preserving most of the signal.

Capturing Variation in Data

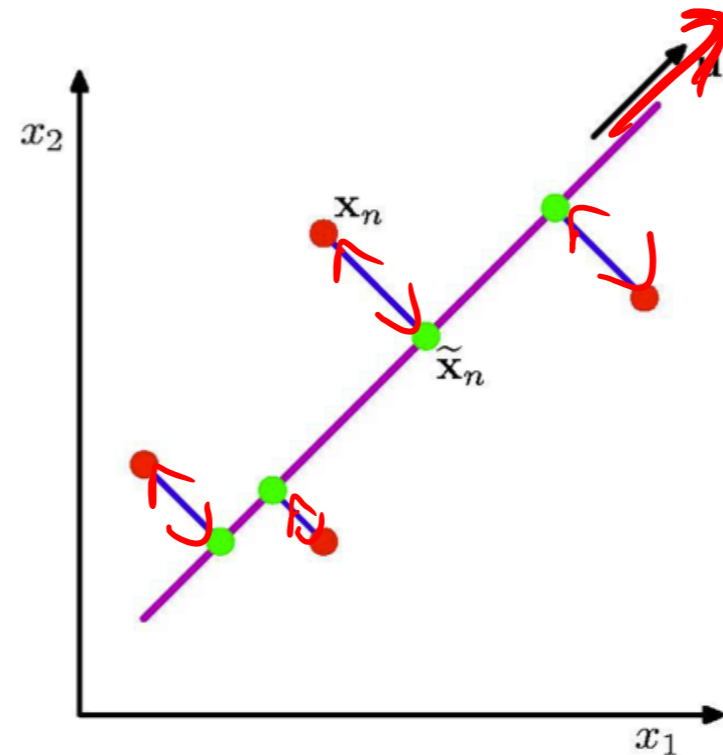


Data vary little
in this direction

Data vary a lot
in this direction

Two features
are correlated

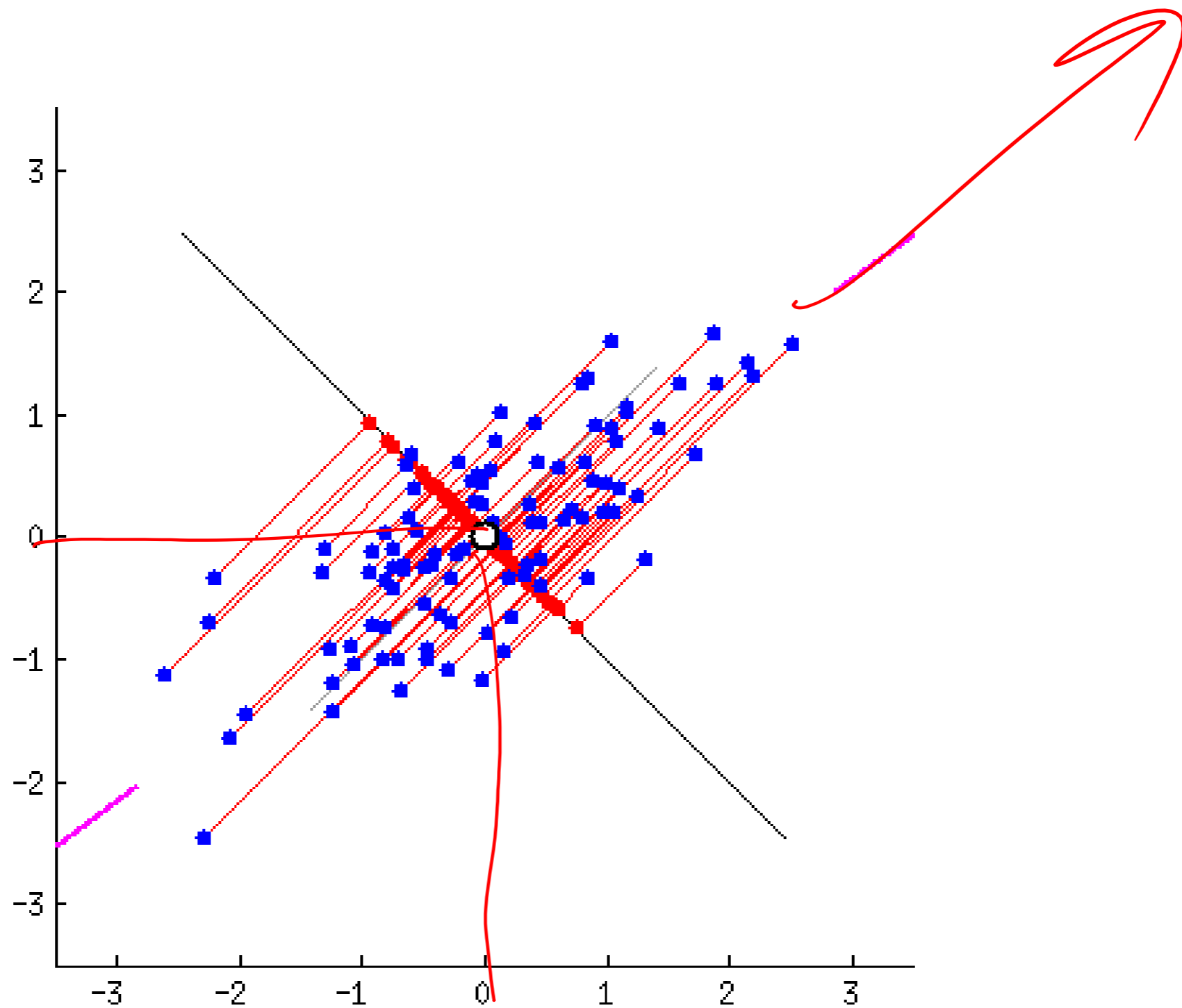
Two Equivalent Perspectives of PCA



PCA:

Orthogonal projection of the data onto a lower-dimension linear space that...

- ❑ maximizes variance of projected data (purple line)
- ❑ minimizes mean squared distance between
 - data point and
 - projections (sum of blue lines)



Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm ←
- PCA and SVD
- Summary

What is variance equation?

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Formulating the Problem

- Given n data points, $\{x_1, x_2, \dots, x_n\} \in R^d$ with their mean $\mu =$

$$\frac{1}{n} \sum_{i=1}^n x_i$$

- Find a direction $w \in R^d$ where

$$\|w\| = \sqrt{\sum_{j \in d} \omega_j^2} = 1$$

We constrain the norm of w to be equal to one to avoid having very large variance in each new dimension.

- Given n data points, $\{x_1, x_2, \dots, x_n\} \in R^d$ with their mean μ

$$\|w\| = \sqrt{\sum_{j \in d} \omega_j^2} = 1 \qquad \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

- Such that the variance (or variation) of the data along direction w is maximized

$$\max_{\|w\|=1} \frac{1}{n} \sum_{i=1}^n (x_i w - \mu w)^2$$

variance in new feature space

An Optimization Problem

- Manipulate the objective with linear algebra

$$\frac{1}{n} \sum_{i=1}^n (x_i w - \mu w)^2 = \frac{1}{n} \sum_{i=1}^n ((x_i - \mu)w)^2 =$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{((x_i - \mu)w)^T}_{A} \underbrace{((x_i - \mu)w)}_B = \frac{1}{n} \sum_{i=1}^n w^T (x_i - \mu)^T (x_i - \mu) w$$

$$\underbrace{(AB)^T}_{(A \ B)} = B^T A^T$$

$$w^T \left(\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^T (x_i - \mu) \right) w = w^T C w$$

Covariance matrix

Let's optimize it

$$f(w) = w^T C w$$

$$g(w) = \|w\|_2 - 1 = w^T w - 1$$

$$L(w, \lambda) = w^T C w - \lambda (w^T w - 1)$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow 2Cw - 2\lambda w = 0$$

$$\Rightarrow Cw = \lambda w$$

eigenvector of C
eigenvalues of C
eigenvector

Square
Symmetrical

$$X \underline{v} = \lambda \underline{v}$$

Data

eigenvalues

~~X~~ → find w or a different basis which captures max. variation in the data → find C → find eigenvectors of C

will

Equivalence to The Eigenvalue Problem

Objective function: $\max_{\|w\|=1} w^T C w$

- Form lagrangian function of the optimization problem

$$L(w, \lambda) = w^T C w + \lambda(1 - w^T w)$$

If w is a maximum of the original optimization problem, then there exist a λ , where (w, λ) is a **stationary point** of $L(w, \lambda)$

Therefore:

$$\frac{\partial L(w, \lambda)}{\partial w} = 0 = 2Cw - 2\lambda w \Rightarrow$$

$$Cw = \lambda w$$

Eigen-Value Problem

- Eigen-value problem

d : dimension

- Given a symmetric matrix $C \in R^{d \times d}$

C is also a positive semidefinite matrix

- Find a vector $w \in R^d$ and $\|w\| = 1$

- Such that

$$Cw = \lambda w$$

- There will be multiple solution of w_1, w_2, \dots, w_d for its corresponding $\lambda_1, \lambda_2, \dots, \lambda_d$

- They are ortho-normal: $w_i^T w_i = 1$ $w_i^T w_j = 0$

Principal Direction of the Data

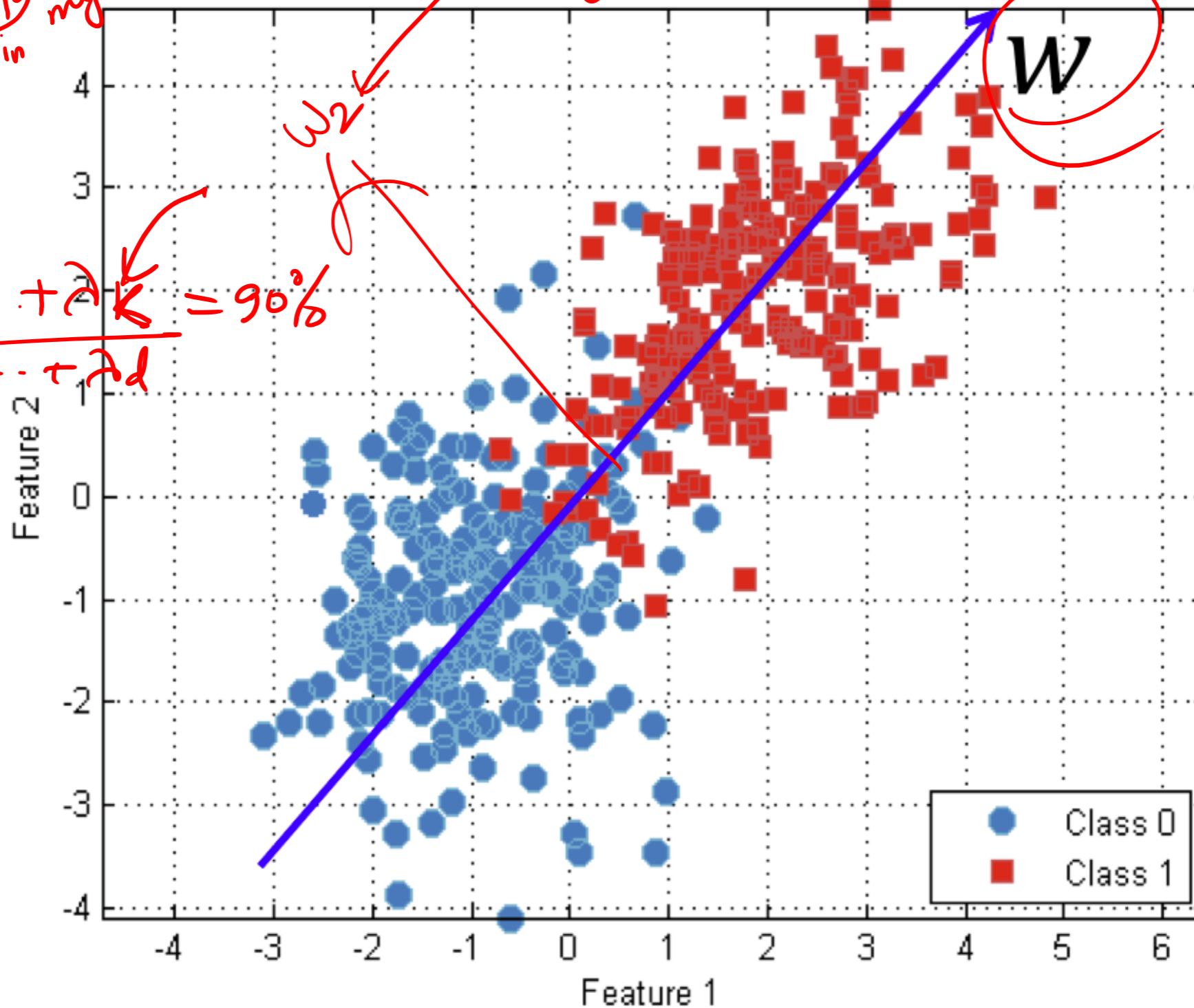
hyperparameter

I want to capture 90% variation in my data

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_d} = 90\%$$

principal direction for the rest of the subspace

principal direction



Variance in the Principal Direction

eq. from solving optimization

- Principal direction w satisfies

$$Cw = \lambda w = w\lambda$$

- Variance in principal direction is

Variance in 2 direction

$$w^T C w$$

*$\lambda_1 + \lambda_2 + \dots + \lambda_d$
all variation in data*

$$= w^T w \lambda$$

$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_d} = 90\%$

$$= \lambda$$

eigen-value

Multiple Principal Directions

- Directions w_1, w_2, \dots which has
 - the largest variances
 - but are **orthogonal** to each other
- Take the eigenvectors w_1, w_2, \dots of C corresponding to
 - the largest eigenvalue λ_1 ,
 - the second largest eigenvalue λ_2
 - ...

Relations Between Principal Components

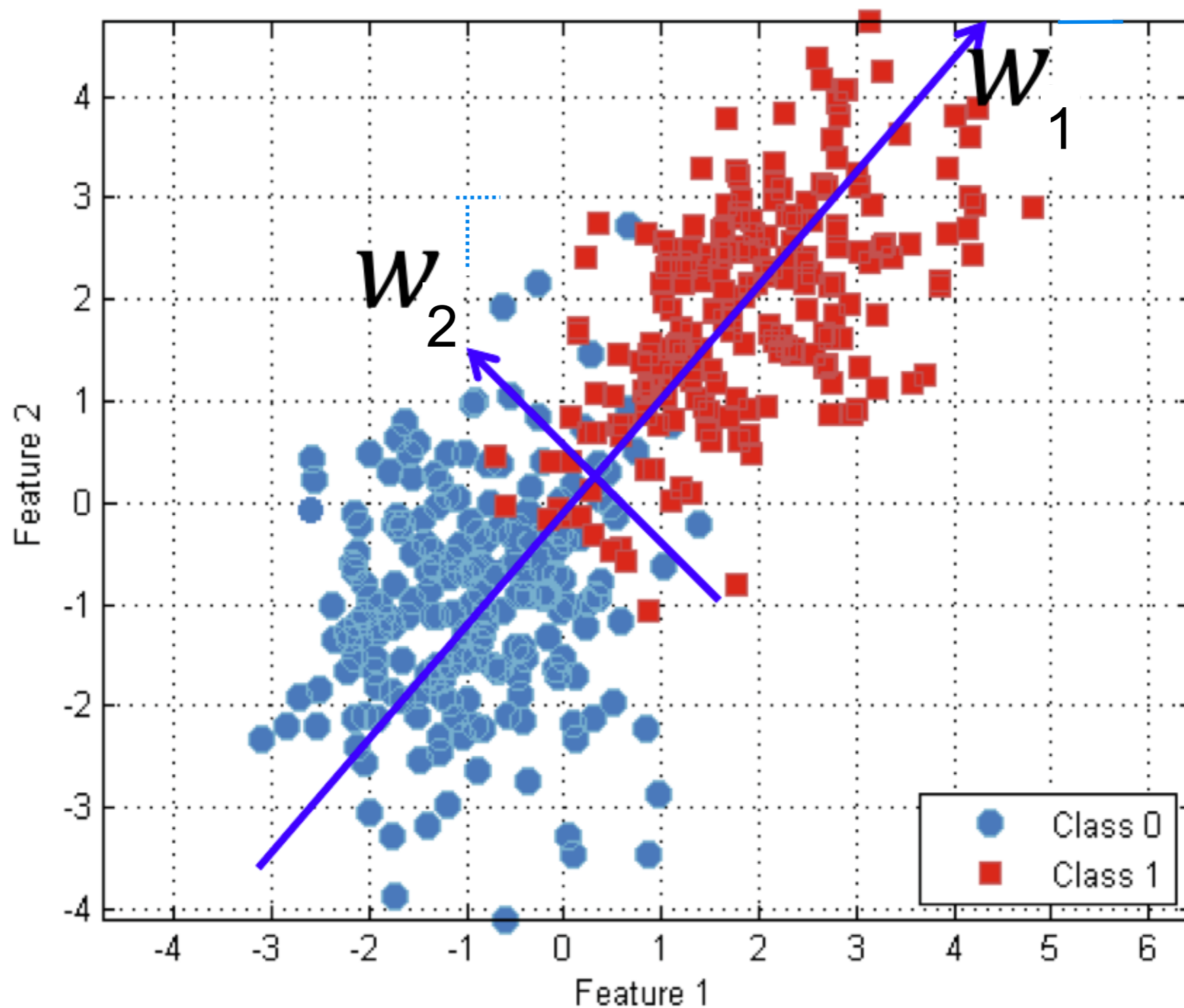
PC direction → eigenvector
direction of C
PC score → variation value

Principal component #1: points in the direction of the **largest variance**.

Each subsequent principal component

- is **orthogonal** to the previous ones, and
- points in the directions of the **largest variance of the residual subspace**

How to decide how many dimensions we want to reduce



The PCA Algorithm

- Given n data points, $\{x_1, x_2, \dots, x_n\} \in R^d$ with mean

- Step 1: Estimate the mean and covariance matrix from data

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad C = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^T (x_i - \mu)$$

Principal directions

- Step 2: Take the eigenvectors w_1, w_2, \dots of C corresponding to the largest eigenvalue λ_1 , the second largest eigenvalue $\lambda_2 \dots$

eigenvectors of largest eigenvalue

- Step 3: Compute reduced representation

$$z_i = \left(\frac{(x_i - \mu_1)}{\sigma_1} w_1 \quad \frac{(x_i - \mu_2)}{\sigma_2} w_2 \quad \dots \right) \quad z \Rightarrow \underline{n \times k}$$

$k \ll d$

Standardization

Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD ←
- Summary

Very computationally expensive

$$C = \frac{X^T X}{n}$$

eigen decomposition
→ square
→ orthogonal → symmetric

(np.svd(x)) → X can be any shape
→ numerically stable
→ Pretty fast

Singular Value Decomposition

$$\bar{X}_{n \times d}$$

n: instances

d: dimensions

X is a centered matrix

$$\bar{X} = U \Sigma V^T$$

$U_{n \times n} \rightarrow$ unitary matrix $\rightarrow U \times U^T = I$

$\Sigma_{n \times d} \rightarrow$ diagonal matrix

$V_{d \times d} \rightarrow$ unitary matrix $\rightarrow V \times V^T = I$

$$\begin{matrix}
 X = & \begin{bmatrix} u_{1 \times 1} & \dots & \dots & \dots & u_{1 \times n} \\ \vdots & \ddots & \dots & \dots & \vdots \\ \vdots & \vdots & \ddots & \dots & \vdots \\ \vdots & \vdots & \dots & \ddots & \vdots \\ u_{1 \times 1} & \dots & \dots & \dots & u_{n \times n} \end{bmatrix} & \times & \begin{bmatrix} \Sigma_{1 \times 1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma_{d \times d} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} & \times & \begin{bmatrix} v_{1 \times 1} & \dots & \dots & \dots & v_{1 \times d} \\ \vdots & \ddots & \dots & \dots & \vdots \\ \vdots & \vdots & \ddots & \dots & \vdots \\ \vdots & \vdots & \dots & \ddots & \vdots \\ v_{d \times 1} & \dots & \dots & \dots & v_{d \times d} \end{bmatrix} \\
 & U & & \Sigma & & V^T
 \end{matrix}$$

$n \times n$
 $n \times d$
 $d < n$
 $d \times d$

Principle direction \rightarrow (points to the first row of V)

Matrix compression: K dimensions out of d

According to PCA $\rightarrow Cw = \lambda w = w\lambda$

Centering X

$$\text{Covariance } C_{d \times d} = \frac{1}{n} \sum_{i=1}^n (x^i - \mu)^T \overbrace{(x^i - \mu)} = \frac{\bar{X}^T \bar{X}}{n}$$

$$\bar{X} = U\Sigma V^T$$

$$C = \frac{\bar{X}^T \bar{X}}{n}$$

$$C = \frac{V\Sigma^T U^T U \Sigma V^T}{n} = \frac{V\Sigma^2 V^T}{n}$$

$$C = \frac{V\Sigma^2 V^T}{n} = V \frac{\Sigma^2}{n} V^T$$

$$Cv = \lambda v$$

$$\frac{\Sigma^2}{n}$$

eigenvalues

$$CV = V \frac{\Sigma^2}{n} V^T V = V \frac{\Sigma^2}{n}$$

According to Eigen-decomposition definition $\rightarrow CV = V\Lambda$

V is the eigen vectors of covariance (Principal directions)

$\lambda_i = \frac{\sigma_i^2}{n} \rightarrow$ The eigenvalues of covariance matrix

Let's project the data (X) on principal directions:

$$\bar{X}V = U\Sigma V^T V = U\Sigma$$

$\bar{X}V$ is linear combination of the original data (x-space) features

Projection of one instance (x) on the first principal direction using k dimensions

$$p_1 = [u_{1 \times 1} \Sigma_{1 \times 1}, u_{1 \times 2} \Sigma_{2 \times 2}, \dots, u_{1 \times k} \Sigma_{k \times k}]$$

$$p_2 = [u_{2 \times 1} \Sigma_{1 \times 1}, u_{2 \times 2} \Sigma_{2 \times 2}, \dots, u_{2 \times k} \Sigma_{k \times k}]$$

$$U \Rightarrow n \times k$$

$$\Sigma \Rightarrow k \times k$$

Upper left corner

Eigen values $\lambda = \frac{\Sigma^2}{n}$

Eigenvectors (principal directions) V

$$\bar{X} = U \Sigma V^T$$

$n = 50$

$k = 10$

$U = 50 \times 50$
 $\Sigma = 50 \times d$

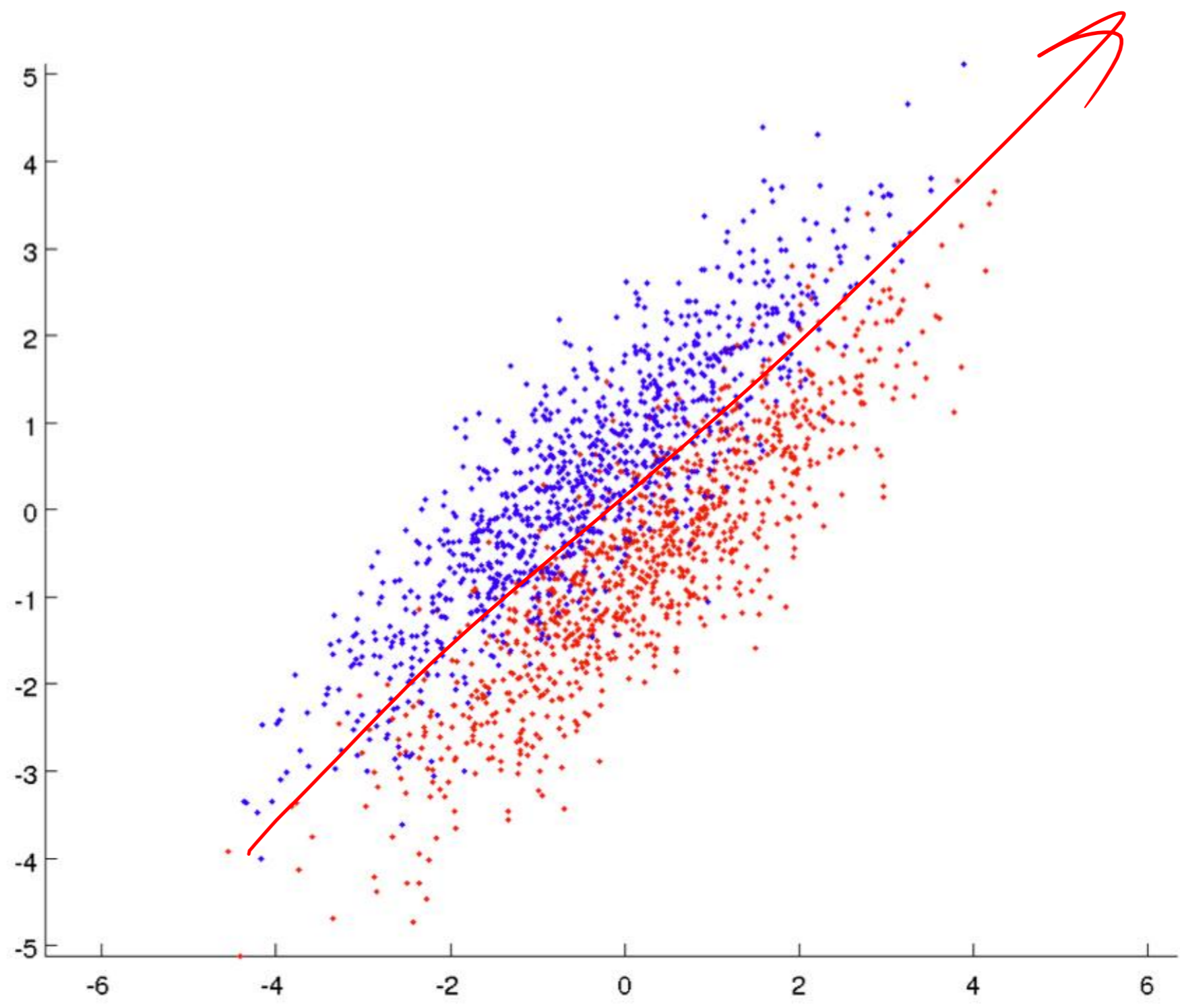
$U \begin{bmatrix} :10 : :10 \\ \end{bmatrix} 10 \times 10$
 $\Sigma = 10 \times d$

Principal components (Scores) or projections on principal directions

In fact, using the SVD to perform PCA makes much better sense numerically than forming the covariance matrix to begin with, since the formation of $X^T X$ can cause loss of precision.

Maybe

Are Principal Components Good for Classification?



Why PCA potentially works in classification?

The dimension with the largest variance corresponds to the dimension with the largest entropy and thus encodes the most information (Information Theory). The smallest eigenvectors will often simply represent noise components.

If a dimension has very little variance, the values are almost the same for all samples, so that dimension does not help distinguish between classes.

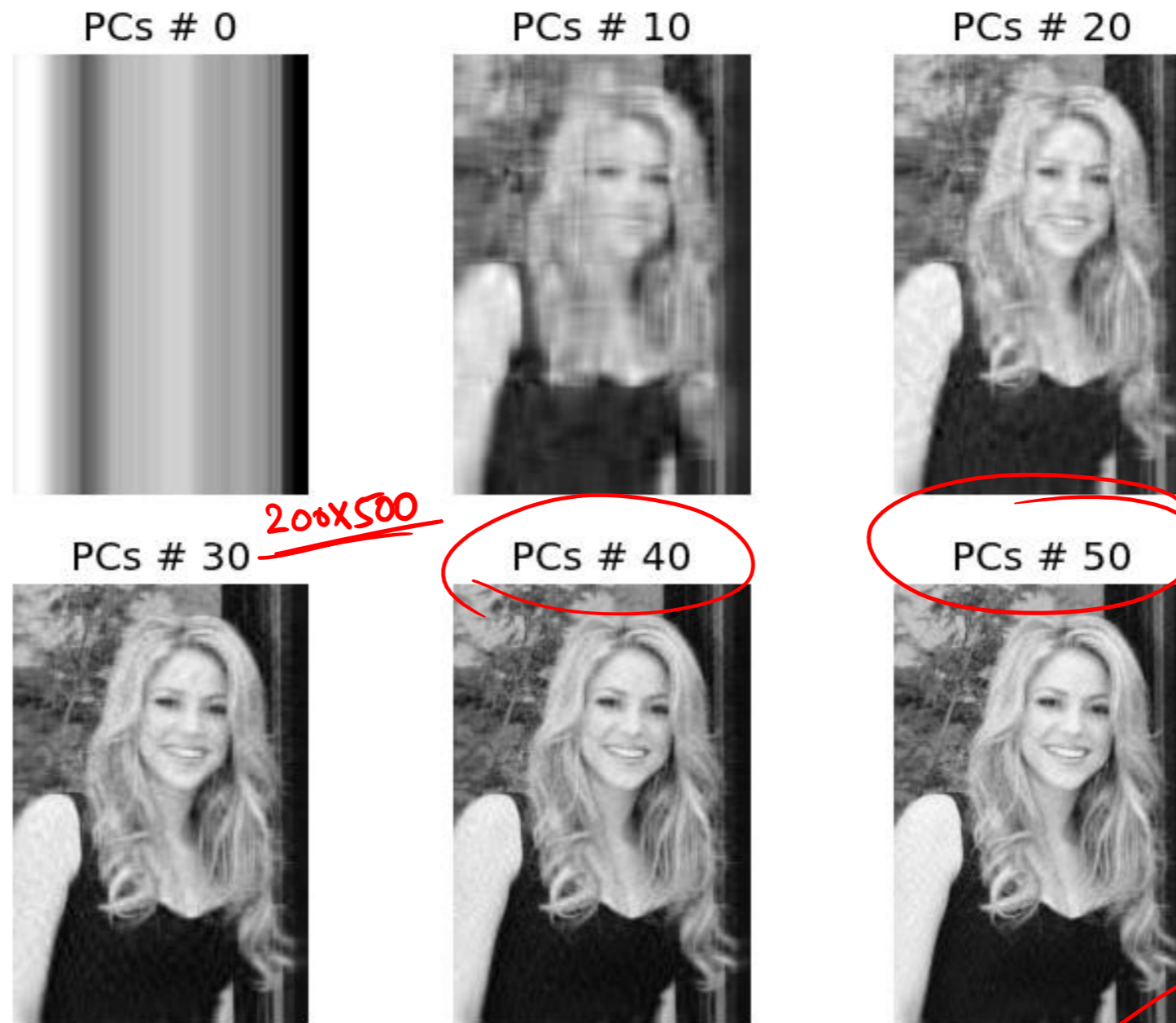
Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary 

Summary

- PCA
 - Finds orthonormal basis for data
 - Sorts dimensions in order of “importance”
 - Discard low significance dimensions
- Uses
 - Get concise low-dimensional representations
 - Remove noise
- Not magic
 - Doesn't know class labels
 - Can only capture linear variations

Image compression using PCA



When does PCA work and for what?

Use Case

1. Aggregating correlated dimensions (feature reduction)

2. Data visualization (2D or 3D)

3. Noise reduction / data compression

4. Preprocessing before classification or clustering

Why PCA Works

merges correlated features into a smaller number of *uncorrelated* components that explain most of the data's variance.

finds the axes of greatest variance, which are usually the most informative for viewing high-dimensional data.

keeps only the top (k) components with the largest eigenvalues, effectively filtering out low-variance (noise) directions.

PCA can remove redundancy and noise, making downstream algorithms (like k-NN, SVM, or k-means) faster and less sensitive to irrelevant features.

When to Use It

When features are linearly related (e.g., many sensors measuring similar signals).

When you want to plot or explore high-dimensional data (e.g., projecting 100D embeddings to 3D).

When you want to compress data while preserving the most structure — e.g., image compression, signal denoising.

When input features are highly correlated or too many for your classifier to handle efficiently.

When PCA struggles or fails?



Issue

Why it matters

Nonlinear relationships

PCA only captures *linear* patterns. If your data lies on a nonlinear manifold (like a spiral or circle), PCA will fail to uncover it. (Techniques like *kernel PCA* or *t-SNE* work better.)

Outliers

PCA is *variance-based*, so outliers can distort the direction of principal components.

Different scales

If you don't normalize your features, those with larger scales dominate the components.

Interpretability

Principal components are linear combinations of features — they may not have a clear real-world meaning.

When all features are uncorrelated

If features are already independent, PCA won't reduce dimensionality — it can't compress noise.

Quick Knowledge Check

- In PCA, which mathematical property defines the new axes (principal components)? A. Eigenvectors of the data. B. Eigenvectors of the covariance matrix. C. Feature correlation coefficients
- Why is it necessary to normalize features before applying PCA? A. To make features linearly independent. B. To ensure features with larger scales do not dominate the variance. C. To remove correlation between variables
- What does a large eigenvalue represent in PCA? A. A noisy or redundant direction. B. A direction with low variance. C. A direction capturing significant variation (signal)
- Why is SVD preferred over directly computing eigenvectors of $X^T X$? A. SVD can be used only for square matrices. B. $X^T X$ can lose precision and is computationally expensive. C. SVD requires manual normalization
- Which of the following is **not** a common application of PCA? A. Image compression. B. Noise reduction. C. Feature extraction for visualization. D. Handling categorical variables directly
- When does PCA tend to fail? A. When features are correlated. B. When data has nonlinear relationships. C. When dimensionality is high
- What do the smallest eigenvalues in PCA typically represent?. A. The most informative directions. B. Noise or unimportant variation in data. C. The directions of maximum entropy