

# Principal Component Analysis

Chetan Arora

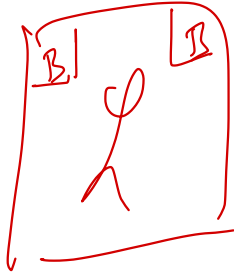
Disclaimer: The contents of these slides are taken from various publicly available resources such as research papers, talks and lectures. To be used for the purpose of classroom teaching, and academic dissemination only.



# Covariance

- Variance and Covariance are a measure of the “spread” of a set of points around their center of mass (mean)
- **Variance** is a measure of the deviation from the mean for points in one dimension e.g. variance in the measuring the length of the same object by different people.
- **Covariance** is a measure of how much each of the dimensions vary from the mean with respect to each other.

$x_i$        $x_j$



A



B



C

$x_1$

$x_2$

$$z_1 = \cos \theta x_1 + \sin \theta x_2$$

$$X = \{x_1, \dots, x_n\}$$



# Covariance

- Covariance is measured between 2 dimensions to see if there is a relationship between the 2 dimensions e.g. number of hours studied & marks obtained.
- The covariance between one dimension and itself is the variance

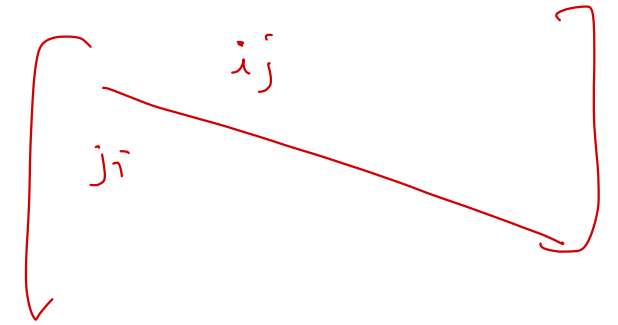


# Covariance

$$\text{Covariance}(X, Y) = \frac{\sum_{i=1}^n (\overset{x_1}{\underset{i}{X}} - \bar{X}) (\overset{x_2}{\underset{i}{Y}} - \bar{Y})}{n}$$

$$\sum_{k=1}^n (X_i^{(k)} - \bar{X}) (X_j^{(k)} - \bar{X})$$

- For a 3-dimensional data set  $(x, y, z)$ , one can measure the covariance between:
  - $x$  and  $y$  dimensions,
  - $y$  and  $z$  dimensions, and
  - $x$  and  $z$  dimensions.



- Covariance is symmetrical:  $\text{Cov}(X, Y)$  =  $\text{Cov}(Y, X)$



# Covariance Matrix

- Covariance between various dimensions is typically represented as a 2D **covariance matrix ( $C$ )**. For a 3-dimensional data  $(X, Y, Z)$ :

$$C = \begin{bmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) & \text{Cov}(X, Z) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) & \text{Cov}(Y, Z) \\ \text{Cov}(Z, X) & \text{Cov}(Z, Y) & \text{Cov}(Z, Z) \end{bmatrix}$$

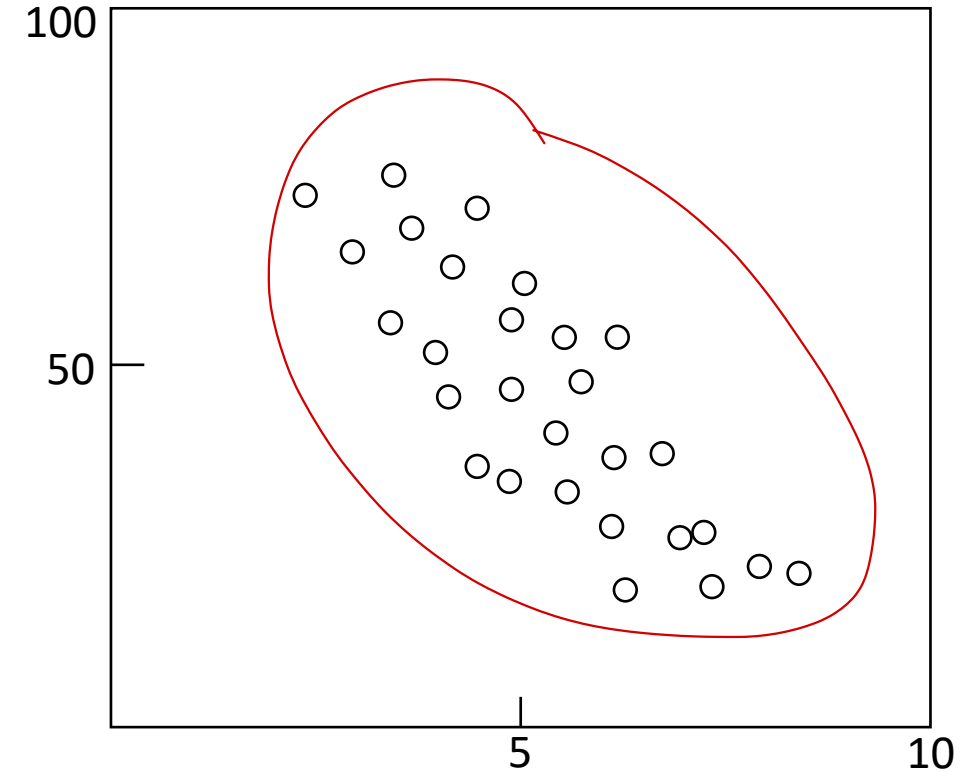
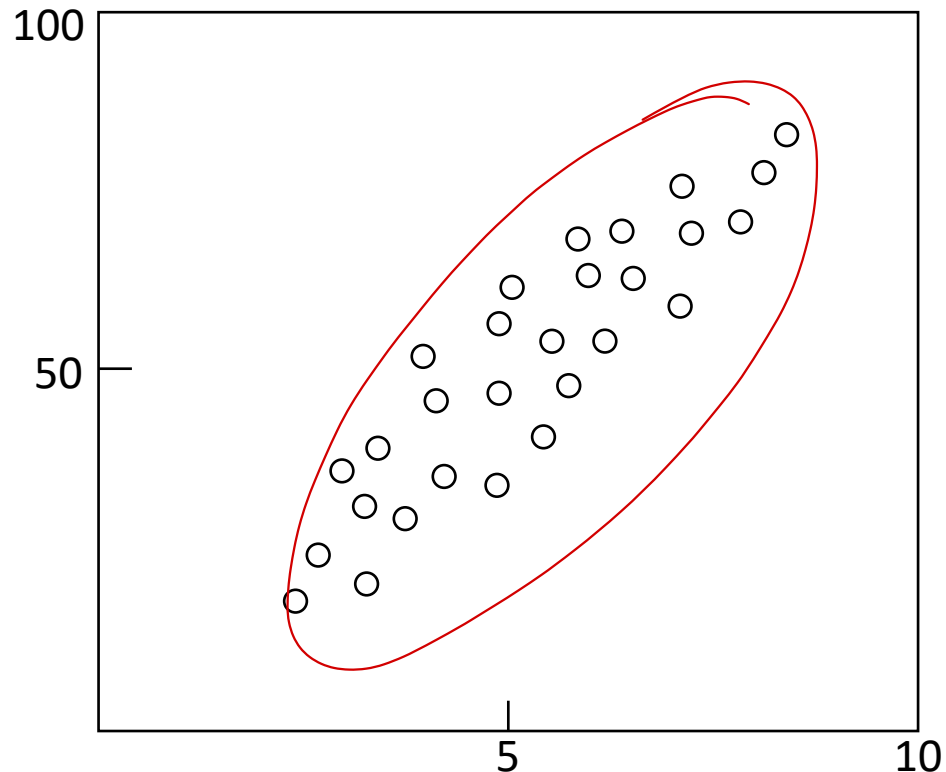
$d \times d$

- Diagonal elements are the variances of  $X$ ,  $Y$  and  $Z$ .
- The covariance matrix  $C$  is symmetrical about the diagonal



# Interpreting Covariance

- **Example:** A 2-dimensional data set.
  - $x$ : number of hours studied for a subject
  - $y$ : marks obtained in that subject





# Interpreting Covariance

- A positive value of covariance indicates both dimensions increase or decrease together e.g. as the number of hours studied increases, the marks in that subject increase.
- A negative value indicates while one increases the other decreases, or vice-versa e.g. number of hours on facebook vs performance in CS dept.
- If covariance is zero: the two dimensions are independent of each other e.g. heights of students vs the marks obtained in a subject





# Interpreting Covariance

Q. Why bother with calculating covariance when we could just plot the 2 values to see their relationship?



# Interpreting Covariance

Q. Why bother with calculating covariance when we could just plot the 2 values to see their relationship?

A. Covariance calculations are used to find relationships between dimensions in high dimensional data sets (usually greater than 3) where visualization is difficult.



# PCA

- Principal components analysis (PCA) is a technique that can be used to simplify a dataset
- It is a linear transformation that chooses a new coordinate system for the data set such that greatest variance by any projection of the data set comes to lie on the first axis (then called the first principal component), the second greatest variance on the second axis, and so on.
- PCA can be used for reducing dimensionality by eliminating the later principal components.



# PCA: Toy Example

- Consider following 3D points:

$x_1$	1	2	4	3	5	6
$x_2$	2	4	8	6	10	12
$x_3$	3	6	12	9	15	18

- Number of integers to be stored: 18



# PCA: Toy Example

- How about the following representation?

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 1 * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

*Handwritten notes: Red circles around the first vector and the scalar 1. Red arrows point to the vectors. Red text 'i' and 'n' are written next to the vectors.*

$$\begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix} = 2 * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

*Handwritten notes: Red circle around the scalar 2. Red arrow points to the vector.*

$$\begin{bmatrix} 4 \\ 8 \\ 12 \end{bmatrix} = 4 * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

*Handwritten notes: Red circle around the scalar 4. Red arrow points to the vector.*

$$\begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix} = 3 * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$\begin{bmatrix} 5 \\ 10 \\ 15 \end{bmatrix} = 5 * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$\begin{bmatrix} 6 \\ 12 \\ 18 \end{bmatrix} = 6 * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

- Number of integers to be stored: 9

*Handwritten: d =*

*Handwritten: k subspace*

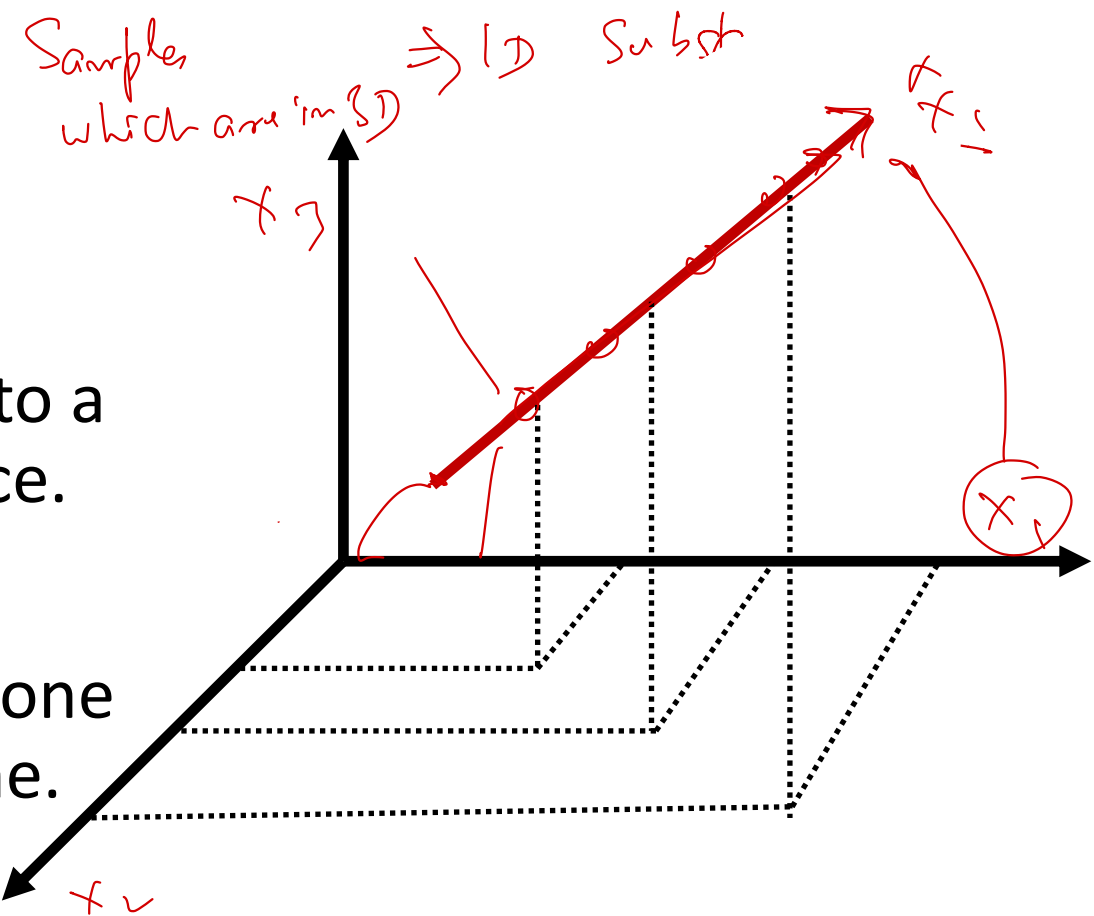
*Handwritten: k < d < c d*

*Handwritten: Basis Vector*



# Geometrical Interpretation

- View each point in 3D space. But in this example, all the points happen to belong to a line: a 1D subspace of the original 3D space.
- Consider a new coordinate system where one of the axes is along the direction of the line.
- In this coordinate system, every point has only one non-zero coordinate.
- We only need to store the direction of the line (3 integers) and the non-zero coordinate for each of the 6 points (6 integers).





# Principal Component Analysis

- Given a set of points, how do we know if they can be compressed like in the previous example?
- The answer is to look into the correlation between the points
- The tool for doing this is called PCA



# Principal Component Analysis

- By finding the eigenvalues and eigenvectors of the covariance matrix, we find that the eigenvectors with the largest eigenvalues correspond to the dimensions that have the strongest correlation in the dataset.
- This is the principal component.
- PCA is a useful statistical technique that has found applications in:
  - fields such as face recognition and image compression
  - finding patterns in data of high dimension.





*Vector* *Column vector*

# PCA Theorem

- Let  $x_1, x_2, \dots, x_n$  be a set of  $n$ ,  $d \times 1$  vectors and let  $\bar{x}$  be their average:

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}$$

- Let  $X$  be the  $d \times n$  matrix with columns:  $[x_1 - \bar{x} \quad x_2 - \bar{x} \quad \dots \quad x_n - \bar{x}]$

$$X = \begin{bmatrix} x_{11} & \dots & x_{n1} \\ x_{12} & \dots & x_{n2} \\ \vdots & \ddots & \vdots \\ x_{1d} & \dots & x_{nd} \end{bmatrix}$$



# Covariance Matrix

- Let  $Q = XX^T$  be the  $d \times d$  covariance matrix:

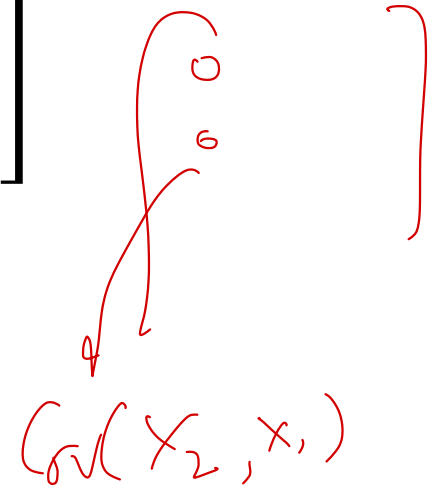
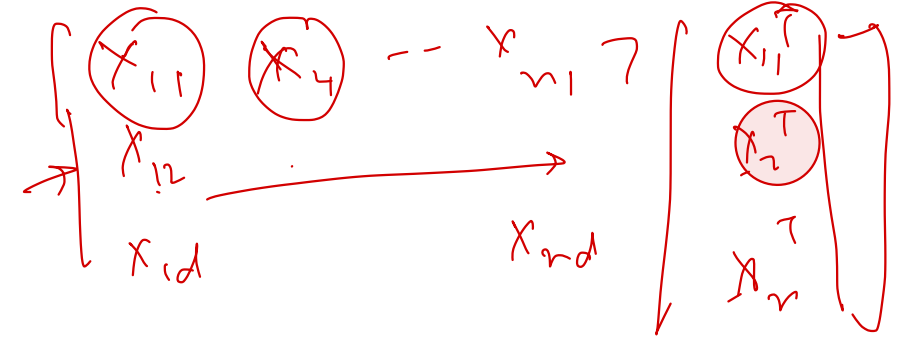
$$Q = [x_1 - \bar{x} \quad x_2 - \bar{x} \quad \dots \quad x_n - \bar{x}] \begin{bmatrix} (x_1 - \bar{x})^T \\ (x_2 - \bar{x})^T \\ \vdots \\ (x_n - \bar{x})^T \end{bmatrix}$$

Handwritten notes: The first matrix is a  $n \times d$  matrix where each row is a data point  $x_i$  minus the mean  $\bar{x}$ . The second matrix is a  $d \times n$  matrix where each column is a data point  $x_i$  minus the mean  $\bar{x}$ . The product  $Q$  is a  $d \times d$  matrix. A handwritten note  $\text{Cov}(x_2, x_1)$  points to the element at row 2, column 1 of  $Q$ .

- Also called **Scatter Matrix** in the context of PCA.

- $Q$  is square as well as symmetric

- $Q$  can be very large (in vision,  $d$  is often the number of pixels in an image!)



$$\begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \\ \vdots \\ x_{1d} \end{bmatrix} \quad \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ x_{id} \end{bmatrix} \quad \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & \vdots & \vdots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{id} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$$



32x32

1024

$d \times d$   
 $d$  eigen vectors  
 $\downarrow$   
 $d$

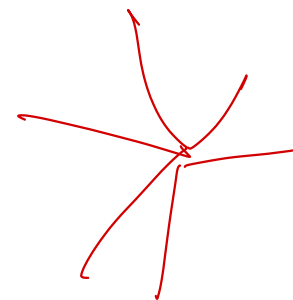
# PCA Theorem

## Theorem:

Each  $x_j$  can be written as:  $x_j = \bar{x} + \sum_{i=1}^d g_{ji} e_i$ , where  $e_i$  are the  $d$  eigenvectors of  $Q$  with non-zero eigenvalues.

- The eigenvectors  $e_1, e_2, \dots, e_d$  span an **eigenspace**.

- $e_1, e_2, \dots, e_d$  are  $d \times 1$  orthonormal vectors.



$e_i$

1  $\hat{a}$   
2  $\hat{b}$   
3  $\hat{c}$

- The scalars  $g_{ji}$  are the coordinates of  $x_j$  in the eigenspace:

$$g_{ji} = (x_j - \bar{x}) \cdot e_i$$

$e_i$

$\lambda_1 > \lambda_2 > \dots > \lambda_d$



# Using PCA to Compress Data

- Expressing  $x$  in terms of  $e_1, \dots, e_d$  has not changed the size of the data
- If the points are highly correlated many of the coordinates of  $x$  will be zero or close to zero (if they indeed lie in a lower-dimensional linear subspace)
- Sort the eigenvectors  $e_i$  according to their eigenvalue:  $\lambda_1 \geq \lambda_2 \geq \dots \lambda_d$

- Assuming  $\lambda \approx 0$ , if  $i > k$ :  $x_j \approx \bar{x} + \sum_{i=1}^k g_{ji} e_i$
- Handwritten notes for the above equation:
- A red circle around  $x_j$  with an arrow pointing to the  $x_j$  in the original equation.
  - A red circle around  $\bar{x}$  with an arrow pointing to the  $\bar{x}$  in the original equation.
  - A red circle around the summation term  $\sum_{i=1}^k g_{ji} e_i$  with an arrow pointing to the summation in the original equation.
  - A red circle around  $x_j$  with a superscript 2, followed by an equals sign and the summation term  $\bar{x} + \sum$ .

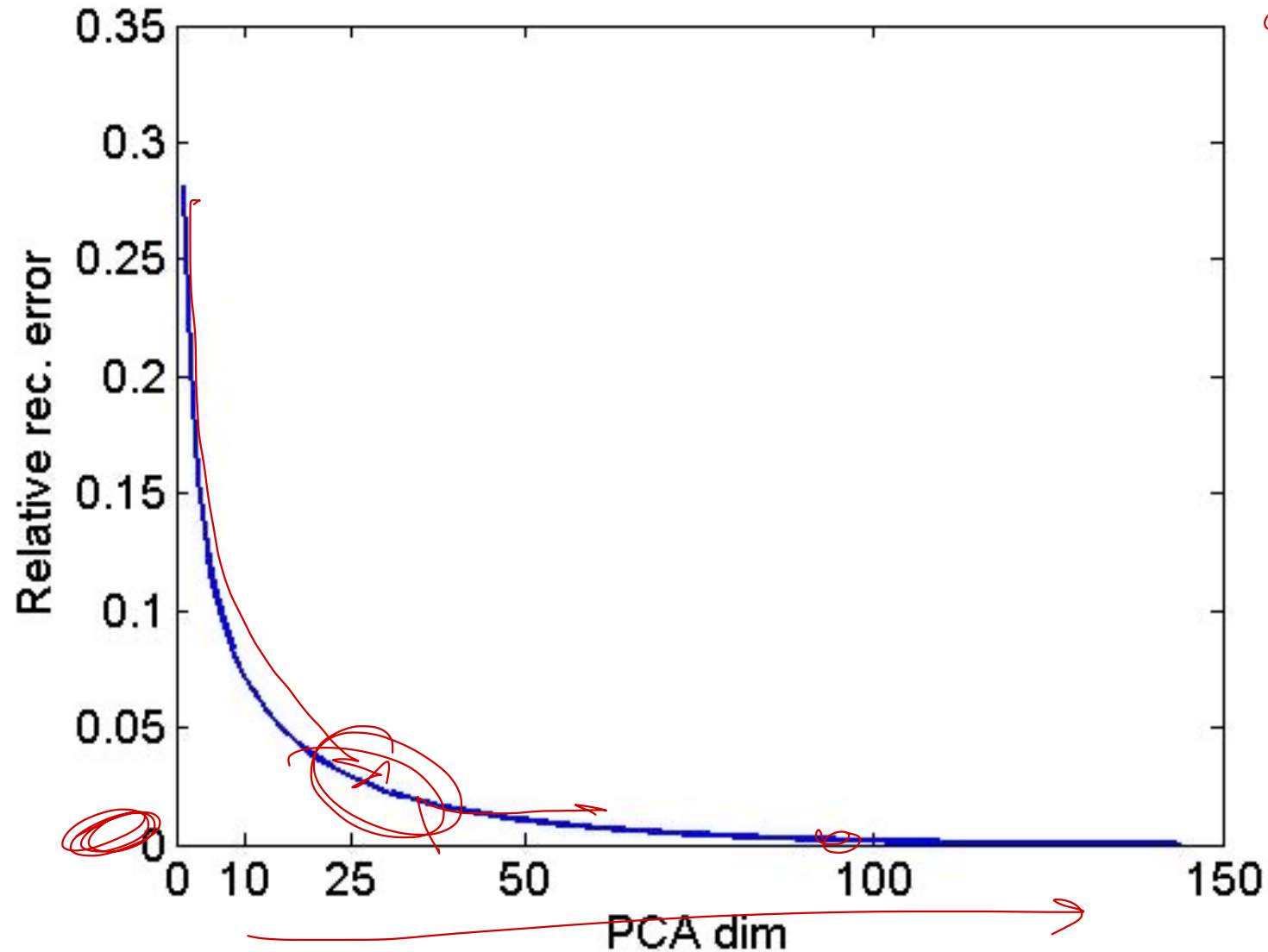
Handwritten notes:

$$\|x_j - \tilde{x}_j\|_2$$

Reconstruction error

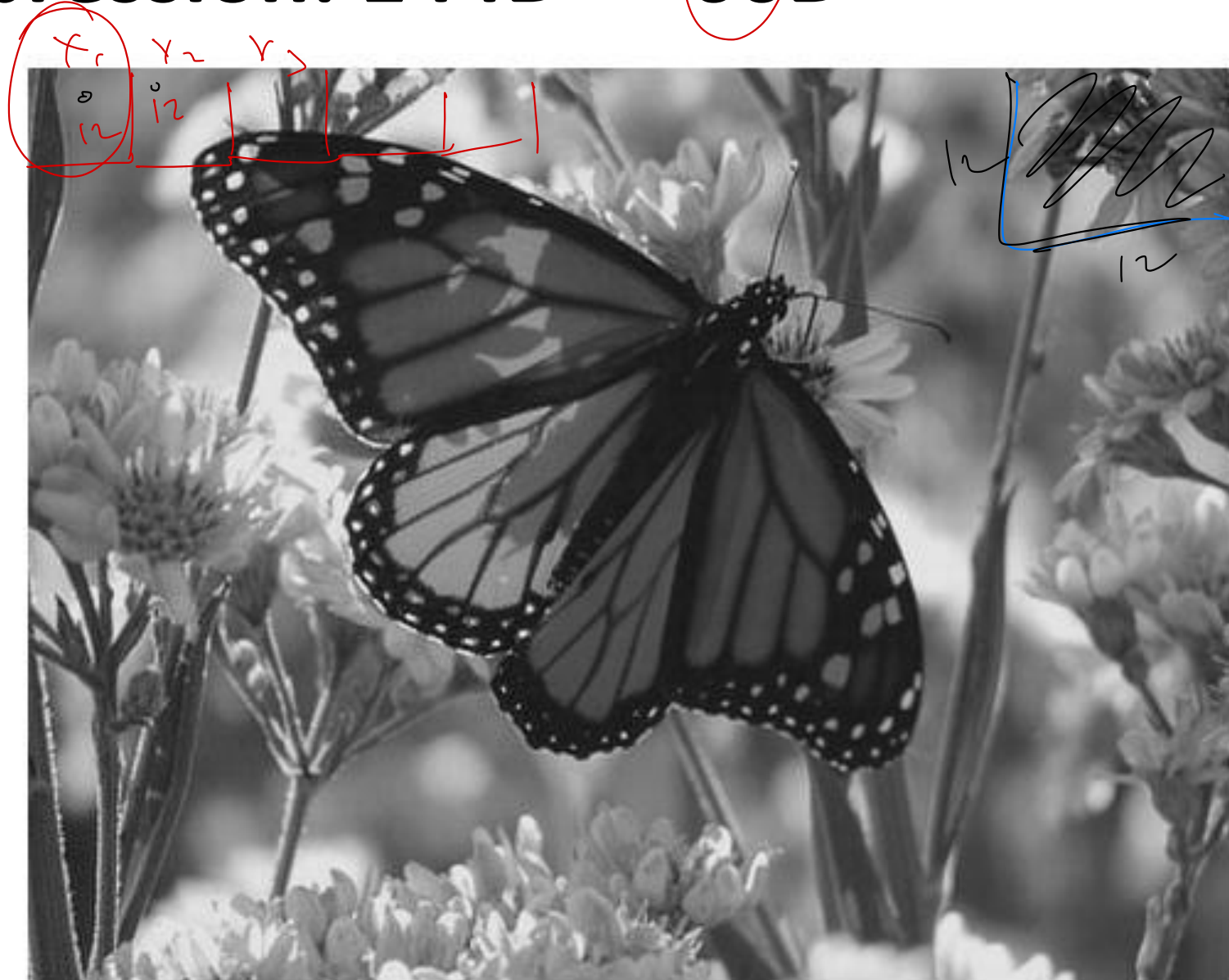


# $L_2$ error and PCA dim





# PCA compression: 144D $\rightarrow$ 60D





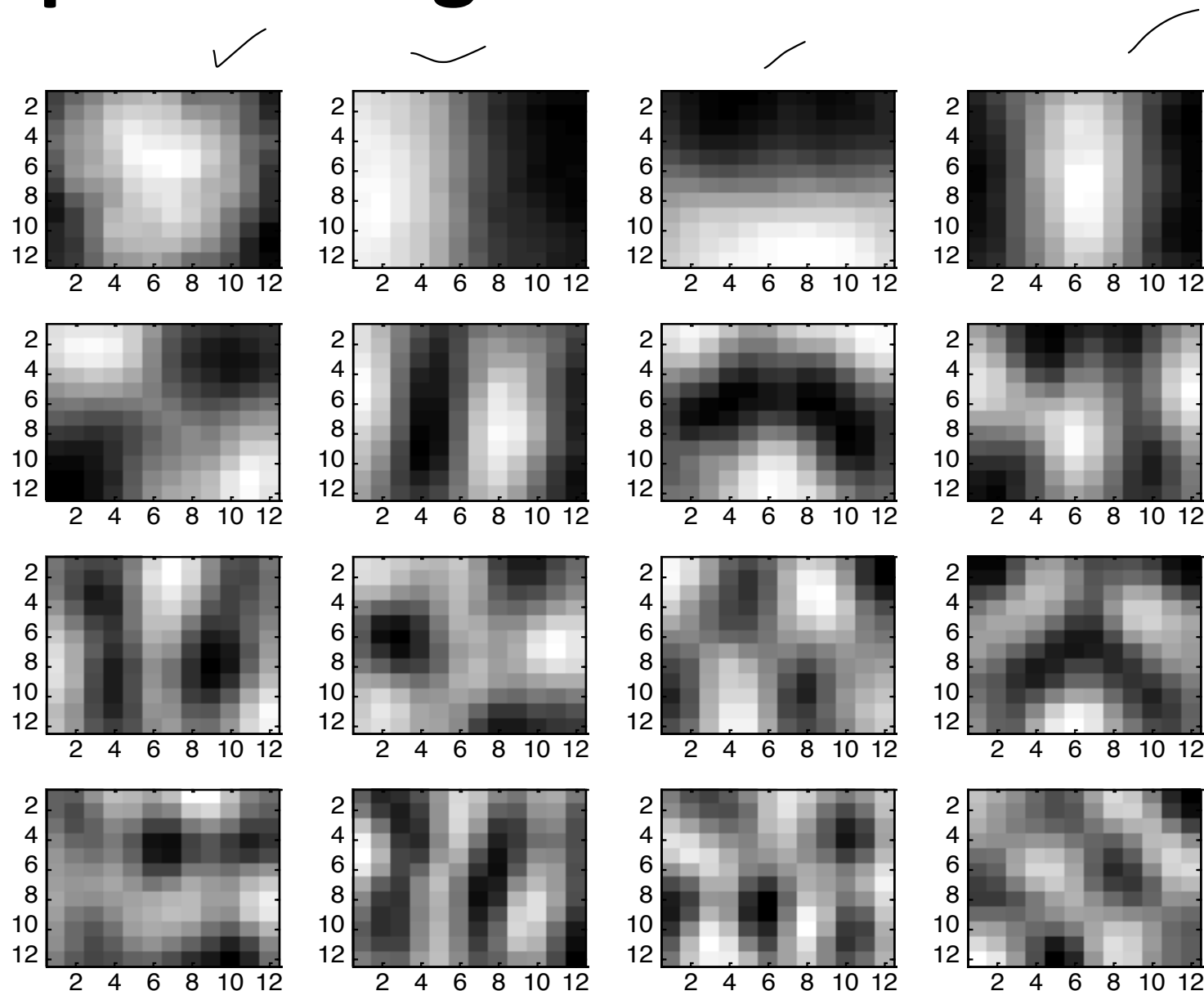
# PCA compression: 144D $\rightarrow$ 16D







# 16 most important eigenvectors



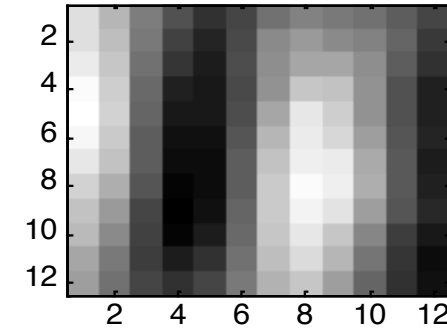
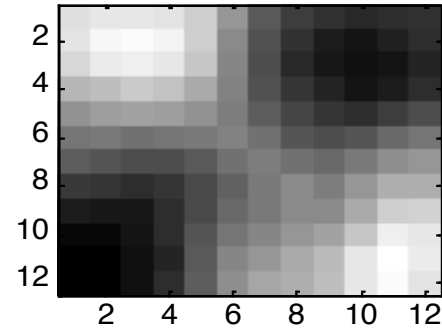
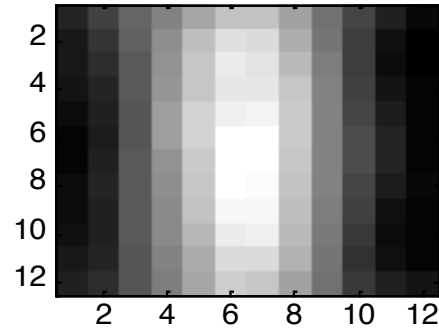
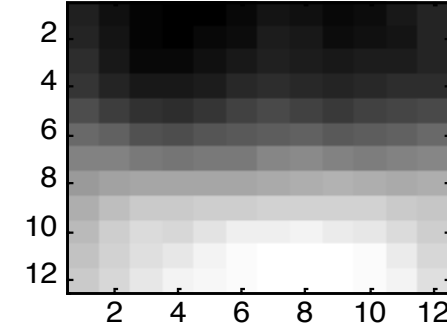
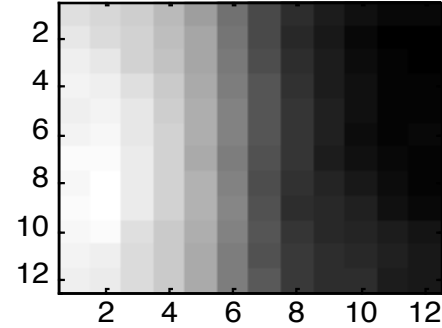
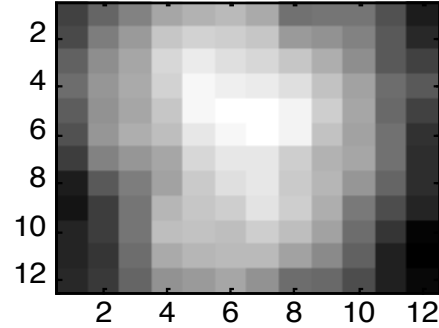


# PCA compression: 144D $\rightarrow$ 6D





# 6 most important eigenvectors



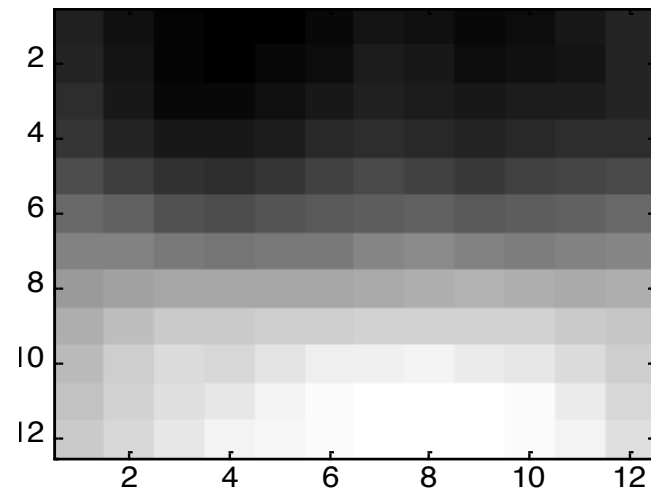
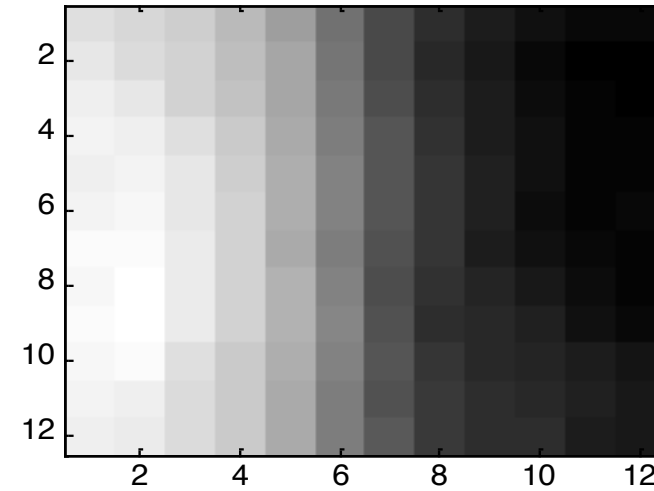
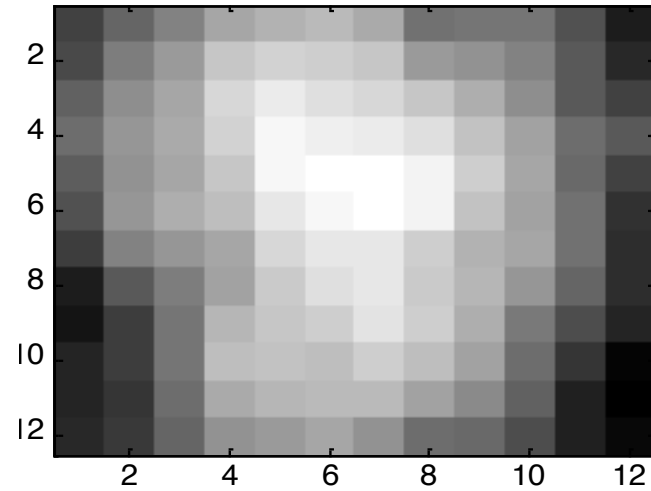


# PCA compression: 144D $\rightarrow$ 3D



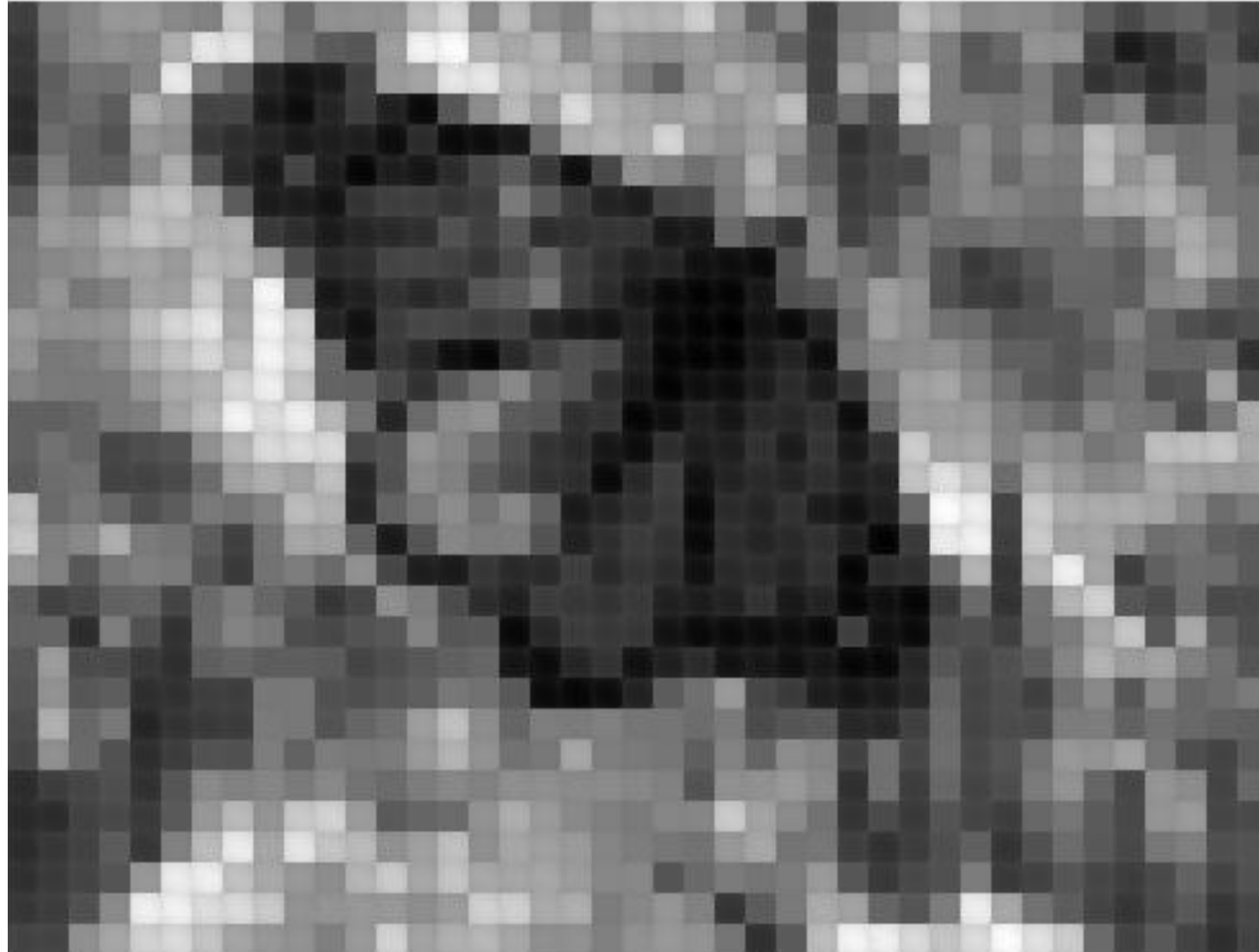


# 3 most important eigenvectors





# PCA compression: 144D $\rightarrow$ 1D

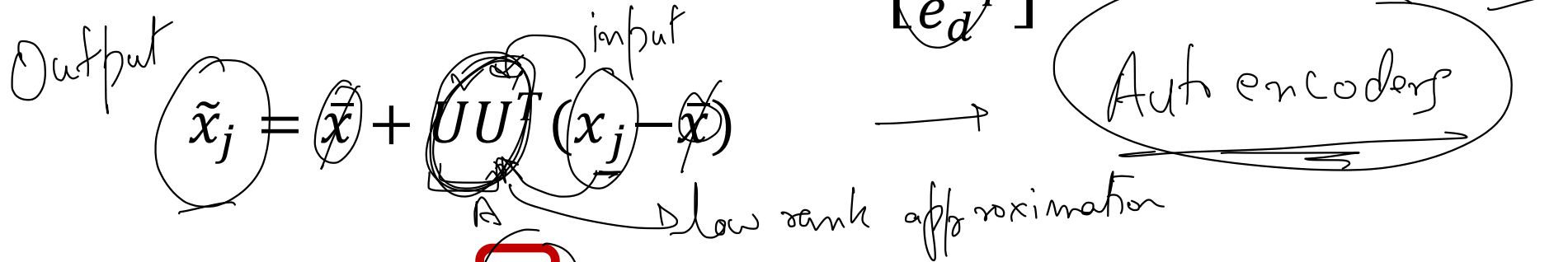




# PCA in Matrix Form

$$U U^T u$$

$$x_j \approx \bar{x} + [e_1 \quad e_2 \quad \dots \quad e_d] \begin{bmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_d^T \end{bmatrix} (x_j - \bar{x})$$



$$= \bar{x} + U z_j$$

Low dimensional representation of  $x_j$

$\| \tilde{x} - x \|_2$   
reconstruction

$$\tilde{x} \approx x$$

$$U \Leftrightarrow U^T$$



# PCA as Minimizing Reconstruction Error

- It can be shown that PCA minimizes the reconstruction error:

$$\min_U \sum_{j=1}^N \|x_j - \tilde{x}_j\|^2 = \sum_{j=1}^N \|x_j - Uz_j\|^2 \quad \swarrow \text{Assuming zero mean}$$

$$\min_U \|X - \overset{\check \check}{UZ}\|_F$$

$\|A\|_F$  shows Frobenius norm,  
sum of elements-wise squares












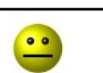



# Matrix Completion/Recommender Systems

- The Netflix problem
  - Movie recommendation: Users watch movies and rate them as good/bad.
  - Because users only rate a few items, one would like to infer their preference for unrated items

1000

1000 x 1000

User	Movie	Rating
	Thor	★ ☆ ☆ ☆ ☆
	Chained	★ ★ ☆ ☆ ☆
	Frozen	★ ★ ★ ☆ ☆
	Chained	★ ★ ★ ★ ☆
	Bambi	★ ★ ★ ★ ★
	Titanic	★ ★ ★ ☆ ☆
	Goodfellas	★ ★ ★ ★ ★
	Dumbo	★ ★ ★ ★ ★
	Twilight	★ ★ ☆ ☆ ☆
	Frozen	★ ★ ★ ★ ★
	Tangled	★ ☆ ☆ ☆ ☆



# Matrix Completion

- Matrix completion problem: Transform the table into a  $N$  users by  $M$  movies matrix called  $R$

Rating matrix

Ninja	2	3	?	?	?	?	?	1	?
Cat	4	?	5	?	?	?	?	?	?
Angel	?	?	?	3	5	5	?	?	?
Nursey	?	?	?	?	?	?	2	?	?
Tongey	?	5	?	?	?	?	?	?	?
Neutral	?	?	?	?	?	?	?	?	1
	Chained	Frozen	Bambi	Titanic	Goodfellas	Dumbo	Twilight	Thor	Tangled



# Matrix Completion

- Matrix completion problem: Transform the table into a  $N$  users by  $M$  movies matrix called  $R$
- **Data:** Users rate some movies.  $R$  matrix is very sparse.
- **Task:** Predict missing entries, i.e. how a user would rate a movie they haven't previously rated
- **Evaluation Metric:** Squared error (used by Netflix Competition)



# Matrix Completion

- Let the representation of user  $i$  in the  $K$ -dimensional space be  $u_i$  and the representation of movie  $j$  be  $z_j$ 
  - Intuition: maybe the first entry of  $u_i$  says how much the user likes horror films, and the first entry of  $z_j$  says how much movie  $j$  is a horror film.

*Frobenius*

- Assume the rating user  $i$  gives to movie  $j$  is given by a dot product:

$$\sum_i u_i^2$$

$$R_{ij} = u_i \cdot z_j$$

$u_i$  (User  $i$ )       $z_j$  (movie  $j$ )

$$\|R - UZ\|_F$$

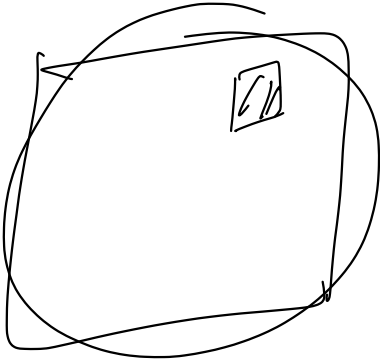
observed

$10^4 \times 10$        $10^4 \times 20$

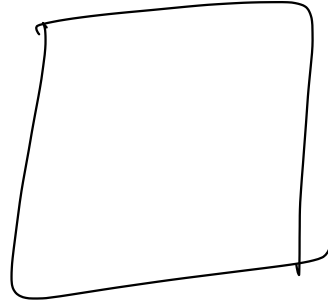
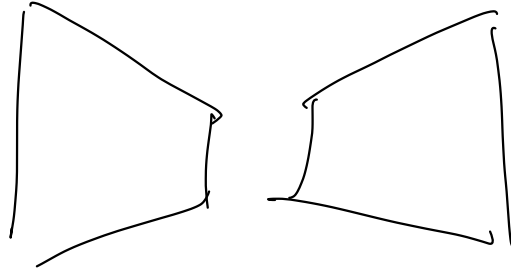
$$\|R - UZ\|_F$$

Sum of element wise squares

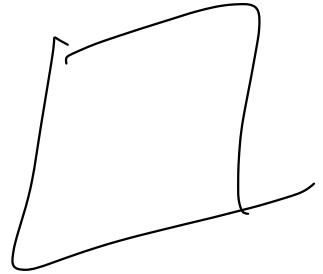
# Pseudo Labels



Masked  
AE



1m → Encoder  
↓  
1029



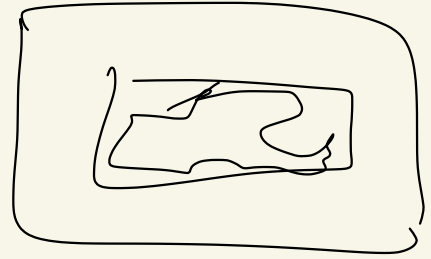
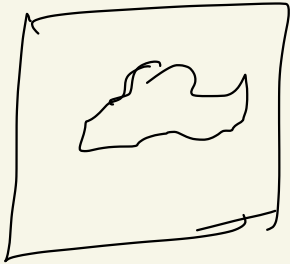
→ Supervised

→ Un-Supervised

Weakly Supervised →

Semi-supervised

└──────────────────┐  
└──────────────────┘ labeled →



Reinforcement Learning

→ unlabeled

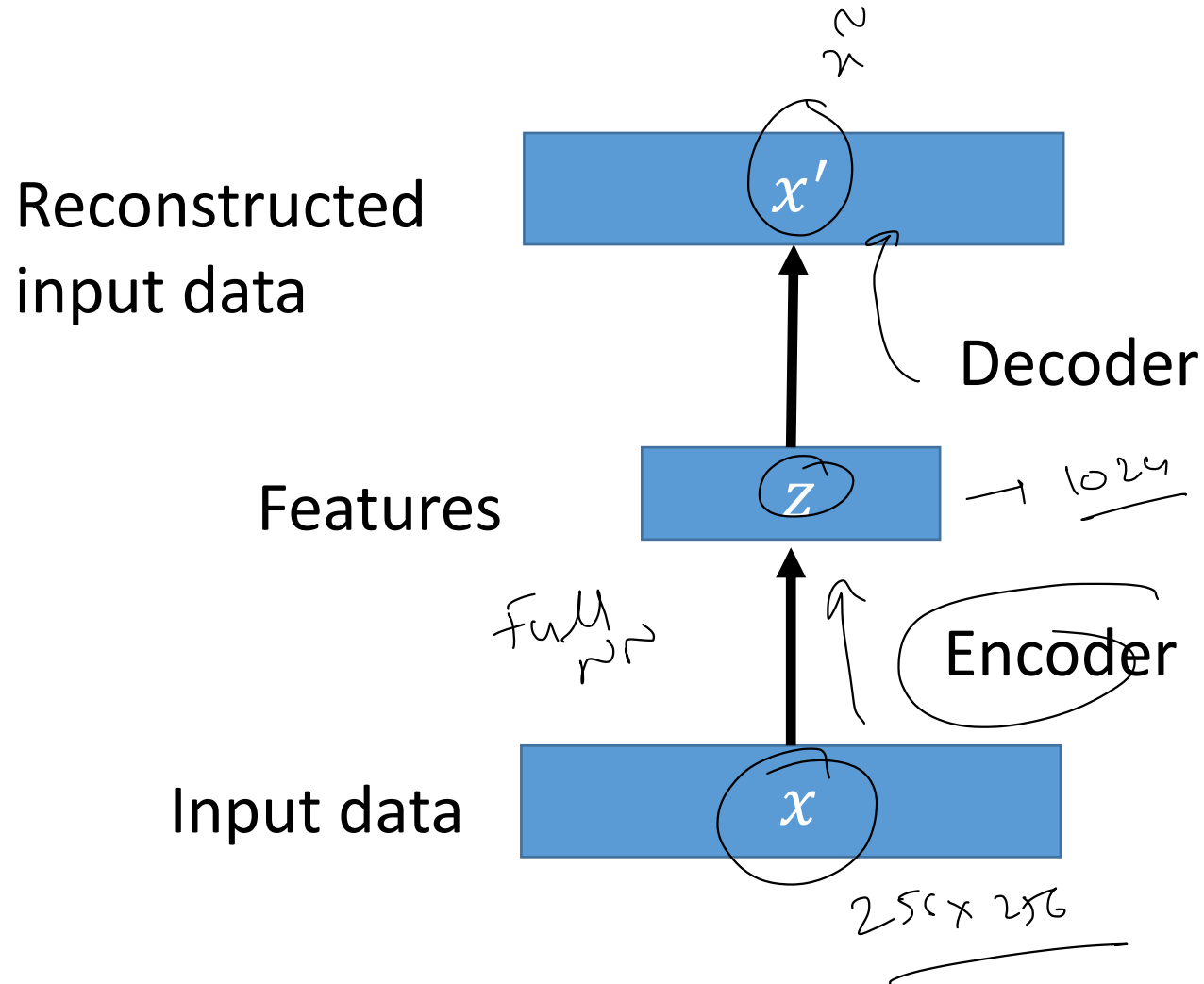


# Matrix Completion

- In matrix form:  $U^T = [u_1 \quad u_2 \quad \dots \quad u_N]$
- In matrix form:  $Z^T = [z_1 \quad z_2 \quad \dots \quad z_M]$
- $R = UZ^T$
- **Matrix completion problem:**  $\min_{U,Z} \|R - UZ^T\|_F$
- Frobenius norm computed over only observed values.



# Autoencoders as Non-Linear Factorization



Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data. Can we generate new images from an autoencoder?

$$\|z - z'\|_2 = U U^T z$$