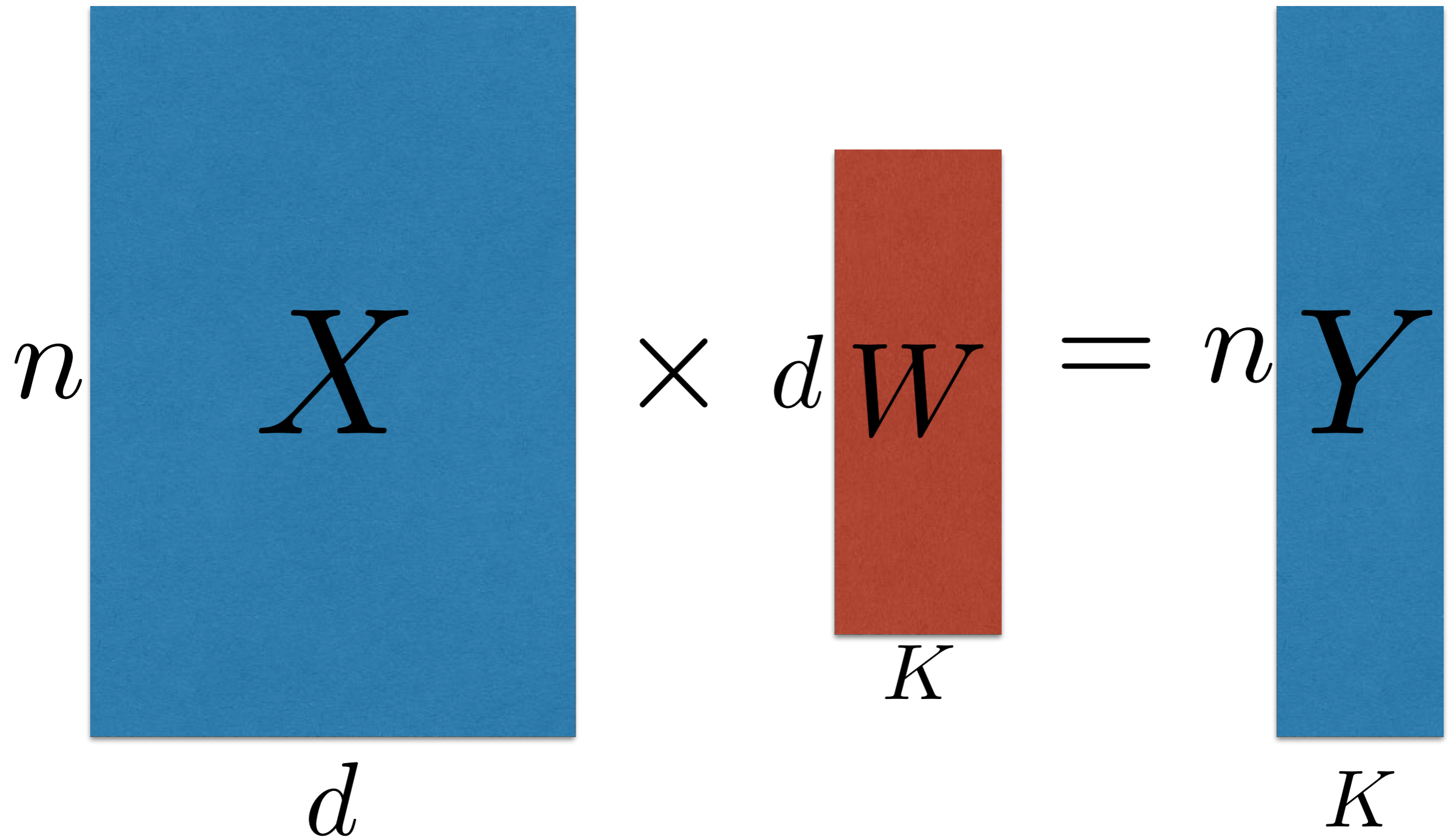


# Machine Learning for Data Science (CS4786)

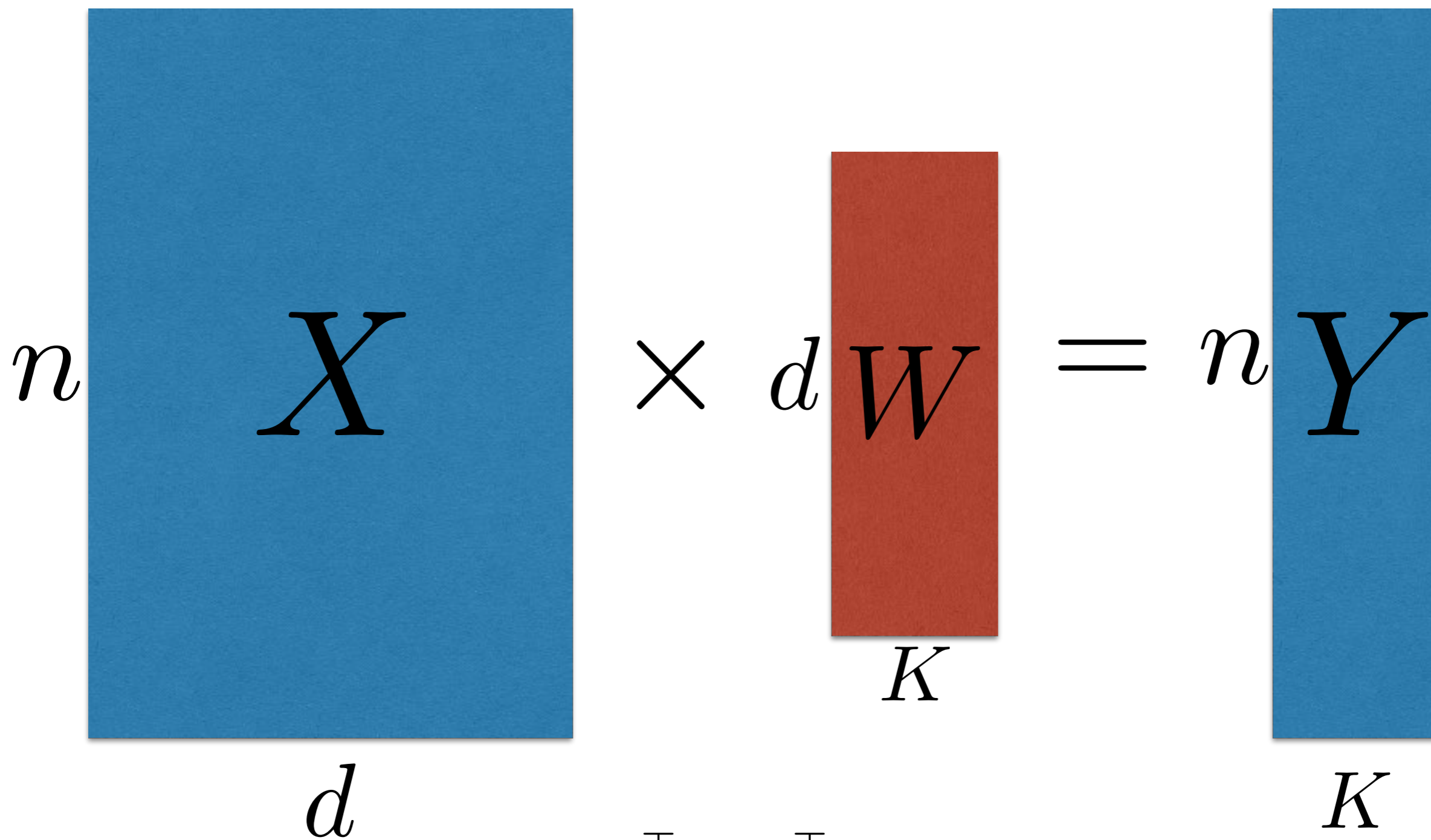
## Lecture 4

PCA and Random Projections

# DIM REDUCTION: LINEAR TRANSFORMATION



# DIM REDUCTION: LINEAR TRANSFORMATION



$$\mathbf{y}_i^\top = \mathbf{x}_i^\top W$$

# PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most

$$\text{Spread} = \frac{1}{n} \sum_{t=1}^n \text{dist}^2 \left( y_t, \frac{1}{n} \sum_{t=1}^n y_t \right) = \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j$$

$$\text{Maximize } \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j$$

$$\text{s.t. } \forall j, \|\mathbf{w}_j\|_2^2 = 1 \ \& \ \mathbf{w}_j \perp \mathbf{w}_i$$

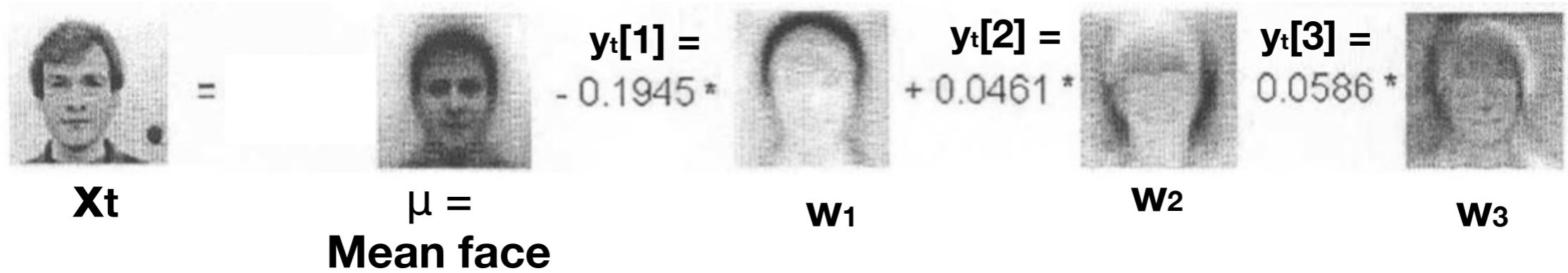
$\Sigma$  is the covariance matrix

This solutions is given by  $W =$  Top  $K$  eigenvectors of  $\Sigma$

# PRINCIPAL COMPONENT ANALYSIS (PCA)

Turk & Pentland'91

Eigen Face:



The diagram illustrates the Eigenface reconstruction process. It shows a target face image  $X_t$  (a man's face) on the left, followed by an equals sign. To the right of the equals sign is the mean face  $\mu$  (a woman's face), labeled "Mean face". This is followed by a plus sign and the first eigenface  $w_1$  (a face with a dark shadow on the left side), with the coefficient  $y_t[1] = -0.1945 *$  above it. This is followed by another plus sign and the second eigenface  $w_2$  (a face with a dark shadow on the right side), with the coefficient  $y_t[2] = +0.0461 *$  above it. Finally, there is a plus sign and the third eigenface  $w_3$  (a face with a dark shadow on the top), with the coefficient  $y_t[3] = 0.0586 *$  above it.

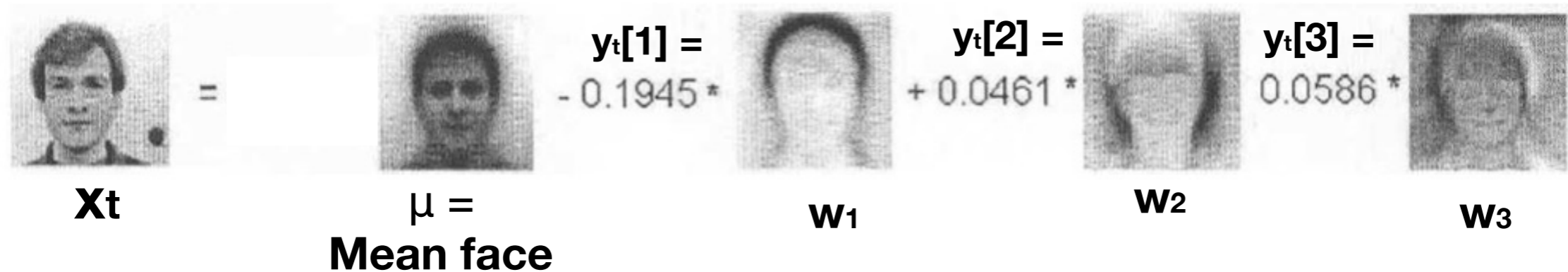
$$X_t = \mu + y_t[1] w_1 + y_t[2] w_2 + y_t[3] w_3$$

$X_t$        $\mu =$  Mean face       $w_1$        $w_2$        $w_3$

# PRINCIPAL COMPONENT ANALYSIS (PCA)

Turk & Pentland'91

Eigen Face:

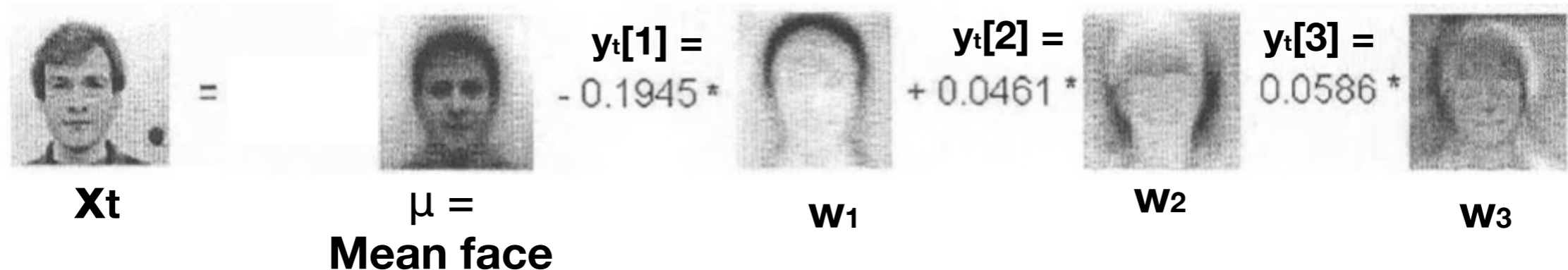


- Each  $x_t$  (each row of  $X$ ) is a face image (vectorized version)

# PRINCIPAL COMPONENT ANALYSIS (PCA)

Turk & Pentland'91

Eigen Face:

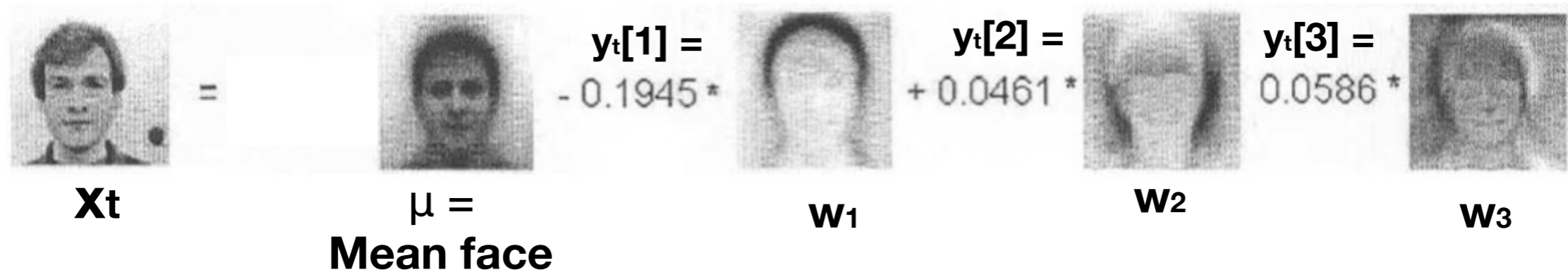


- Each  $x_t$  (each row of  $X$ ) is a face image (vectorized version)
- Each  $y_t$  is the set of coefficients we multiply to the eigen face

# PRINCIPAL COMPONENT ANALYSIS (PCA)

Turk & Pentland'91

Eigen Face:



- Each  $x_t$  (each row of  $X$ ) is a face image (vectorized version)
- Each  $y_t$  is the set of coefficients we multiply to the eigen face
- $w_i$ 's are orthogonal to each other and of unit length



# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,
- (Centered) Data-points as linear combination of some orthonormal basis, i.e.



$$\mathbf{x}_t = \boldsymbol{\mu} + \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j$$

where  $\mathbf{w}_1, \dots, \mathbf{w}_d \in \mathbb{R}^d$  are the orthonormal basis and  $\boldsymbol{\mu} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t$ .

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,
- (Centered) Data-points as linear combination of some orthonormal basis, i.e.



$$\mathbf{x}_t = \boldsymbol{\mu} + \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j$$

where  $\mathbf{w}_1, \dots, \mathbf{w}_d \in \mathbb{R}^d$  are the orthonormal basis and  $\boldsymbol{\mu} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t$ .

- Represent data as linear combination of just  $K$  orthonormal basis,



$$\hat{\mathbf{x}}_t = \boldsymbol{\mu} + \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2 \\ &= \left\| \sum_{j=K+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2 \\ &= \left\| \sum_{j=K+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{but } \|a + b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2a^\top b)\end{aligned}$$



# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2 \\ &= \left\| \sum_{j=K+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{but } \|a + b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2a^\top b) \\ &= \left( \sum_{j=K+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 + 2 \sum_{j=K+1}^d \sum_{i=j+1}^d \mathbf{y}_t[j] \mathbf{y}_t[i] \mathbf{w}_j^\top \mathbf{w}_i \right)\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2 \\ &= \left\| \sum_{j=K+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{but } \|a + b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2a^\top b) \\ &= \left( \sum_{j=K+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 + 2 \sum_{j=K+1}^d \sum_{i=j+1}^d \mathbf{y}_t[j] \mathbf{y}_t[i] \mathbf{w}_j^\top \mathbf{w}_i \right) \\ &= \sum_{j=K+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \left\| \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2 \\ &= \left\| \sum_{j=K+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{but } \|a + b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2a^\top b) \\ &= \left( \sum_{j=K+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 + 2 \sum_{j=K+1}^d \sum_{i=j+1}^d \mathbf{y}_t[j] \mathbf{y}_t[i] \mathbf{w}_j^\top \mathbf{w}_i \right) \\ &= \sum_{j=K+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 \quad (\text{last step because } \mathbf{w}_j \perp \mathbf{w}_i)\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 \quad (\text{but } \|\mathbf{w}_j\| = 1)$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\begin{aligned} \frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 \quad (\text{but } \|\mathbf{w}_j\| = 1) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\begin{aligned}\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 \quad (\text{but } \|\mathbf{w}_j\| = 1) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \quad (\text{now } \mathbf{y}_j = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d (\mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}))^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\begin{aligned}\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 \quad (\text{but } \|\mathbf{w}_j\| = 1) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \quad (\text{now } \mathbf{y}_j = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d (\mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}))^2 \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\begin{aligned}\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \|\mathbf{w}_j\|_2^2 \quad (\text{but } \|\mathbf{w}_j\| = 1) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{y}_t[j]^2 \quad (\text{now } \mathbf{y}_j = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d (\mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}))^2 \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j \\ &= \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j\end{aligned}$$



# PCA: MINIMIZING RECONSTRUCTION ERROR

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\begin{aligned} \text{Minimize} \quad & \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j \\ \text{s.t.} \quad & \forall j, \|\mathbf{w}_j\|_2^2 = 1 \ \& \ \mathbf{w}_j \perp \mathbf{w}_i \end{aligned}$$

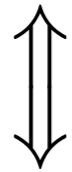


Maximize Total Spread

$$\frac{1}{n} \sum_{t=1}^n \text{dist}^2 \left( y_t, \frac{1}{n} \sum_{t=1}^n y_t \right)$$

Maximize Total Spread

$$\frac{1}{n} \sum_{t=1}^n \text{dist}^2 \left( y_t, \frac{1}{n} \sum_{t=1}^n y_t \right)$$

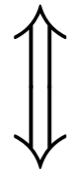


$$\text{Maximize } \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j$$

$$\text{s.t. } \forall j, \|\mathbf{w}_j\|_2^2 = 1 \ \& \ \mathbf{w}_j \perp \mathbf{w}_i$$

Maximize Total Spread

$$\frac{1}{n} \sum_{t=1}^n \text{dist}^2 \left( y_t, \frac{1}{n} \sum_{t=1}^n y_t \right)$$



Maximize  $\sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j$

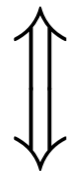
s.t.  $\forall j, \|\mathbf{w}_j\|_2^2 = 1 \ \& \ \mathbf{w}_j \perp \mathbf{w}_i$

Minimize Reconstruction Error

$$\frac{1}{n} \sum_{t=1}^n \|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2$$

Maximize Total Spread

$$\frac{1}{n} \sum_{t=1}^n \text{dist}^2 \left( y_t, \frac{1}{n} \sum_{t=1}^n y_t \right)$$

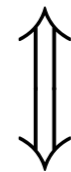


Maximize  $\sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j$

s.t.  $\forall j, \|\mathbf{w}_j\|_2^2 = 1 \ \& \ \mathbf{w}_j \perp \mathbf{w}_i$

Minimize Reconstruction Error

$$\frac{1}{n} \sum_{t=1}^n \|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2$$



Minimize  $\sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j$

s.t.  $\forall j, \|\mathbf{w}_j\|_2^2 = 1 \ \& \ \mathbf{w}_j \perp \mathbf{w}_i$

Claim



# Claim

Maximize Total Spread = Minimize Reconstruction  
Error

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\text{Claim: } \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j = \text{Constant} = \frac{1}{n} \sum_{t=1}^n \|x_t - \mu\|_2^2$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\text{Claim: } \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j = \text{Constant} = \frac{1}{n} \sum_{t=1}^n \|x_t - \mu\|_2^2$$

$$\text{Recall that: } \frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

$$\text{Claim: } \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j = \text{Constant} = \frac{1}{n} \sum_{t=1}^n \|x_t - \mu\|_2^2$$

$$\text{Recall that: } \frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j$$

Take  $K = 0$  so that  $\hat{\mathbf{x}}_t = \mu$

Maximize Total Spread = Minimize Reconstruction Error

$$\text{Minimize } \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

Maximize Total Spread = Minimize Reconstruction Error

$$\text{Minimize } \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

Maximize Total Spread = Minimize Reconstruction Error

$$\text{Minimize } \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \frac{1}{n} \sum_{t=1}^n \|\mathbf{x}_t - \mu\|_2^2 - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

Maximize Total Spread = Minimize Reconstruction Error

$$\text{Minimize } \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \frac{1}{n} \sum_{t=1}^n \|\mathbf{x}_t - \mu\|_2^2 - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$



Maximize Total Spread = Minimize Reconstruction Error

$$\text{Minimize } \sum_{j=K+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \sum_{j=1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( \frac{1}{n} \sum_{t=1}^n \|\mathbf{x}_t - \mu\|_2^2 - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Minimize } \left( - \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \right) \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

$$\iff \text{Maximize } \sum_{j=1}^K \mathbf{w}_j^\top \Sigma \mathbf{w}_j \quad \text{s.t. } \|\mathbf{w}_j\|_2 = 1, \mathbf{w}_j \perp \mathbf{w}_k$$

# PRINCIPAL COMPONENT ANALYSIS

1.  $\Sigma = \text{COV}(X)$

2.  $W = \text{eigs}(\Sigma, K)$

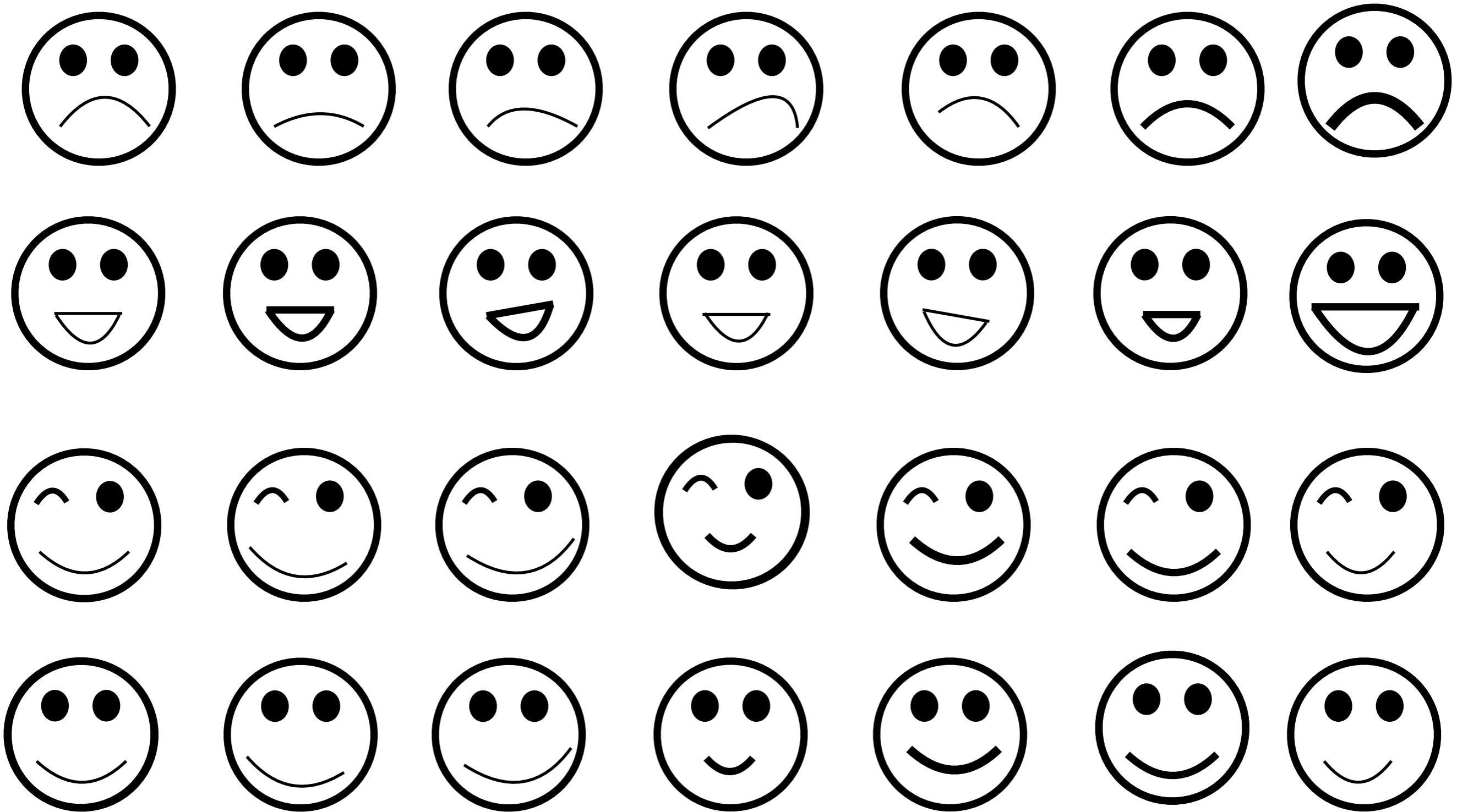
3.  $Y = (X - \mu) \times W$

# RECONSTRUCTION

4.

$$\hat{X} = Y \times W^T + \mu$$

# PRINCIPAL COMPONENT ANALYSIS: DEMO



# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large

# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large
- But we only need top  $K$  eigen vectors.

# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large
- But we only need top  $K$  eigen vectors.
- Idea: use SVD

$$X - \mu = UDV^T$$

# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large
- But we only need top  $K$  eigen vectors.
- Idea: use SVD

$$X - \mu = UDV^T$$

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



# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large
- But we only need top  $K$  eigen vectors.
- Idea: use SVD

$$X - \mu = UDV^T$$

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

Then note that,  $\Sigma = (X - \mu)^T (X - \mu) = VD^2V$

# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large
- But we only need top  $K$  eigen vectors.
- Idea: use SVD

$$X - \mu = UDV^T$$

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

Then note that,  $\Sigma = (X - \mu)^T (X - \mu) = VD^2V$

- Hence, matrix  $V$  is the same as matrix  $W$  got from eigen decomposition of  $\Sigma$ , eigenvalues are diagonal elements of  $D^2$

# WHEN $d \gg n$

- If  $d \gg n$  then  $\Sigma$  is large
- But we only need top  $K$  eigen vectors.
- Idea: use SVD

$$X - \mu = UDV^T$$

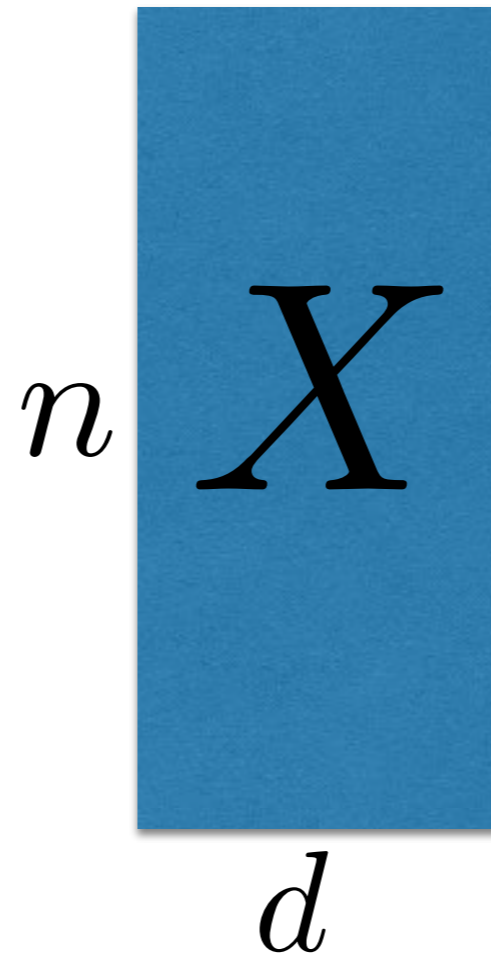
$$\begin{aligned} V^T V &= I \\ U^T U &= I \end{aligned}$$

Then note that,  $\Sigma = (X - \mu)^T (X - \mu) = VD^2V$

- Hence, matrix  $V$  is the same as matrix  $W$  got from eigen decomposition of  $\Sigma$ , eigenvalues are diagonal elements of  $D^2$
- Alternative algorithm:

$$[U, V] = \text{SVD}(X - \mu, K) \quad W = V$$

# The Tall, THE FAT AND THE UGLY



# The Tall, THE FAT AND THE UGLY

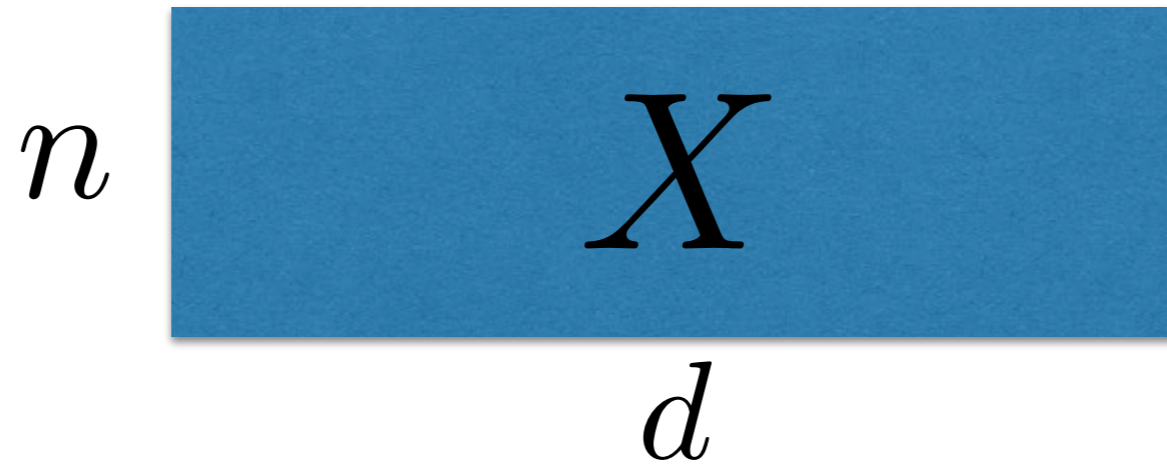
$$\begin{matrix} d \\ \times \\ X^T \\ n \end{matrix} \times \begin{matrix} n \\ \times \\ X \\ d \end{matrix} \Big/ n = d \begin{matrix} d \\ \times \\ \Sigma \end{matrix}$$

# The Tall, THE FAT AND THE UGLY

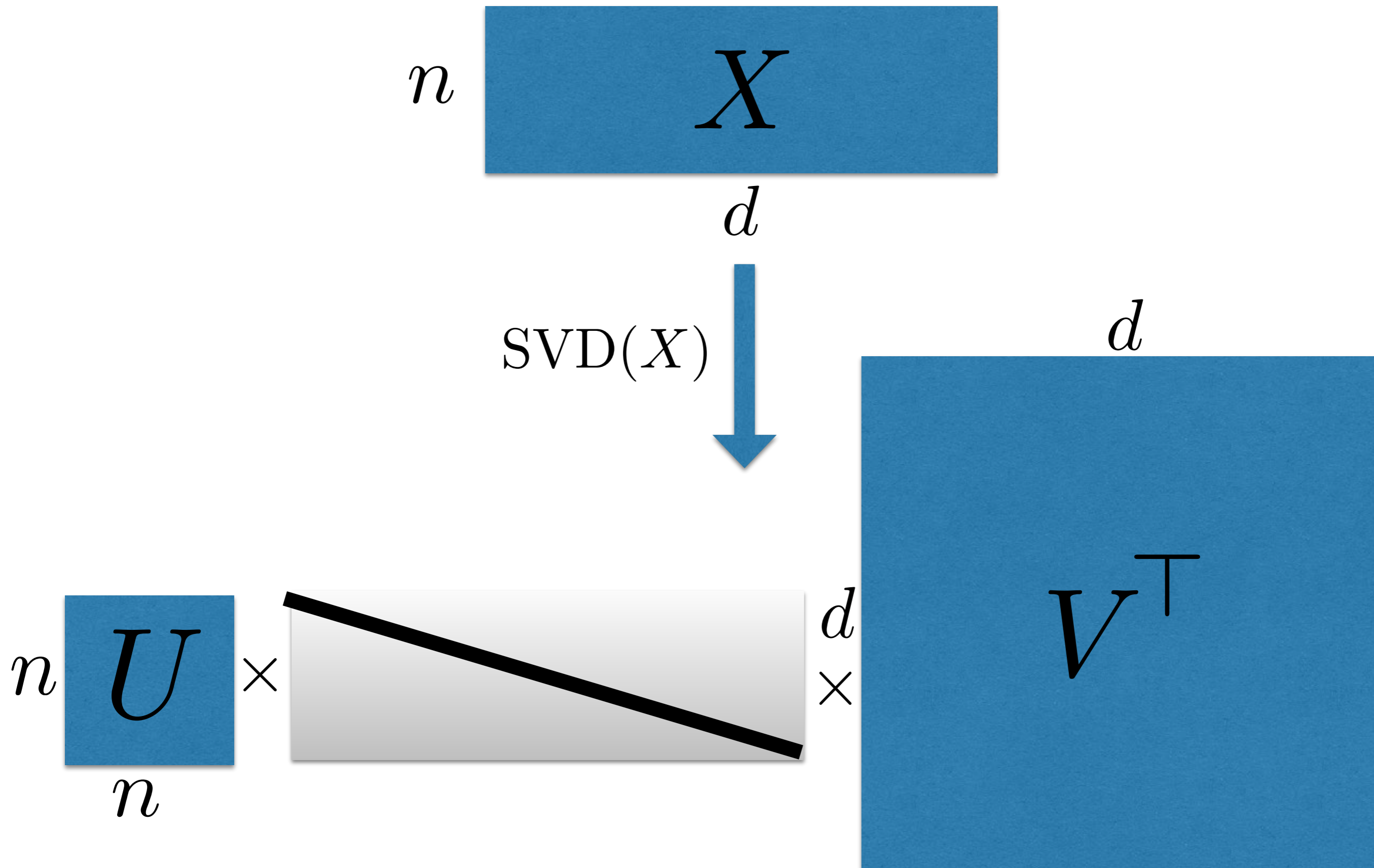
$$\begin{array}{c} d \\ \times \\ n \end{array} X^T \times \begin{array}{c} n \\ \times \\ d \end{array} X \Big/ n = \begin{array}{c} d \\ \times \\ d \end{array} \Sigma$$

$$\begin{array}{c} d \\ \times \\ K \end{array} W = \text{Eigs} \left( \begin{array}{c} d \\ \times \\ d \end{array} \Sigma, K \right)$$

# THE TALL, the Fat AND THE UGLY

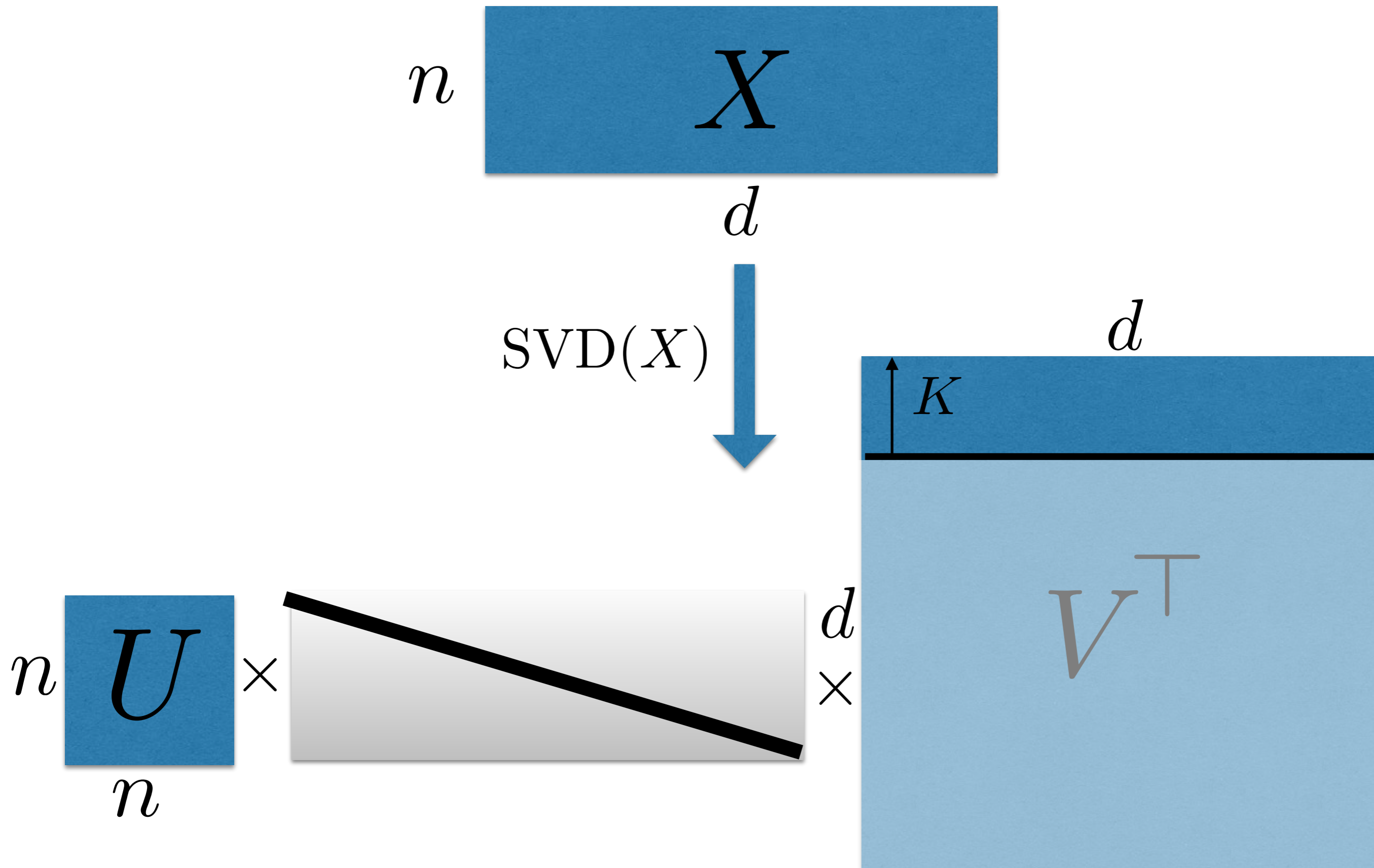


# THE TALL, the Fat AND THE UGLY





# THE TALL, the Fat AND THE UGLY



# THE TALL, THE FAT AND the Ugly

$X$



- $d$  and  $n$  so large we can't even store in memory
- Only have time to be linear in  $\text{size}(X) = n \times d$

I there any hope?

# PICK A RANDOM $W$

# PICK A RANDOM $W$

$$Y = X \times \left[ \begin{array}{ccc} +1 & \dots & -1 \\ -1 & \dots & +1 \\ +1 & \dots & -1 \\ & \cdot & \\ & \cdot & \\ & \cdot & \\ +1 & \dots & -1 \end{array} \right] \Bigg/ \sqrt{K}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

# RANDOM PROJECTION

- What does “it works” even mean?

# RANDOM PROJECTION

- What does “it works” even mean?

Distances between all pairs of data-points in low dim. projection is roughly the same as their distances in the high dim. space.

# RANDOM PROJECTION

- What does “it works” even mean?

Distances between all pairs of data-points in low dim. projection is roughly the same as their distances in the high dim. space.

That is, when  $K$  is “large enough”, with “high probability”, for all pairs of data points  $i, j \in \{1, \dots, n\}$ ,

$$(1 - \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2$$



# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

- Lets start with a one dimensional projection ( $K = 1$ )

$$y_t = \mathbf{x}_t^\top \mathbf{u} \quad \text{where each } \mathbf{u}[i] = \text{random } \pm 1$$

- What is the expected value of:

1.  $y_t - y_s$ ?

2.  $(y_t - y_s)^2$ ?

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

Hence for any  $s, t \in \{1, \dots, n\}$ ,

$$\mathbb{E}[|\mathbf{y}_s - \mathbf{y}_t|^2] = \|\mathbf{x}_s - \mathbf{x}_t\|_2^2$$

Lets try ...

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

Hence for any  $s, t \in \{1, \dots, n\}$ ,

$$\mathbb{E}[|\mathbf{y}_s - \mathbf{y}_t|^2] = \|\mathbf{x}_s - \mathbf{x}_t\|_2^2$$

Lets try ...

Law of large numbers says that average over multiple draws is close to expectation

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

- Like repeating the experiment  $K$  times and averaging

$$\mathbf{y}_t[k] = \mathbf{x}_t^\top \mathbf{u}_k / \sqrt{K} \quad \text{where each } \mathbf{u}_k[i] = \text{random } \pm 1$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

- Like repeating the experiment  $K$  times and averaging

$$\mathbf{y}_t[k] = \mathbf{x}_t^\top \mathbf{u}_k / \sqrt{K} \quad \text{where each } \mathbf{u}_k[i] = \text{random } \pm 1$$

$$(\mathbf{y}_s[k] - \mathbf{y}_t[k])^2 = (\mathbf{x}_t^\top \mathbf{u}_k - \mathbf{x}_s^\top \mathbf{u}_k)^2 / K$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

- Like repeating the experiment  $K$  times and averaging

$$\mathbf{y}_t[k] = \mathbf{x}_t^\top \mathbf{u}_k / \sqrt{K} \quad \text{where each } \mathbf{u}_k[i] = \text{random } \pm 1$$

$$(\mathbf{y}_s[k] - \mathbf{y}_t[k])^2 = (\mathbf{x}_t^\top \mathbf{u}_k - \mathbf{x}_s^\top \mathbf{u}_k)^2 / K$$

$$\|\mathbf{y}_t - \mathbf{y}_s\|_2^2 = \sum_{k=1}^K (\mathbf{y}_s[k] - \mathbf{y}_t[k])^2 = \frac{1}{K} \sum_{k=1}^K (\mathbf{x}_t^\top \mathbf{u}_k - \mathbf{x}_s^\top \mathbf{u}_k)^2$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

- Like repeating the experiment  $K$  times and averaging

$$\mathbf{y}_t[k] = \mathbf{x}_t^\top \mathbf{u}_k / \sqrt{K} \quad \text{where each } \mathbf{u}_k[i] = \text{random } \pm 1$$

$$(\mathbf{y}_s[k] - \mathbf{y}_t[k])^2 = (\mathbf{x}_t^\top \mathbf{u}_k - \mathbf{x}_s^\top \mathbf{u}_k)^2 / K$$

$$\|\mathbf{y}_t - \mathbf{y}_s\|_2^2 = \sum_{k=1}^K (\mathbf{y}_s[k] - \mathbf{y}_t[k])^2 = \frac{1}{K} \sum_{k=1}^K (\mathbf{x}_t^\top \mathbf{u}_k - \mathbf{x}_s^\top \mathbf{u}_k)^2$$

**This is an average over  $K$  trials**

PICK A RANDOM  $W$



# PICK A RANDOM $W$

$$Y = X \times \left[ \begin{array}{ccc} +1 & \dots & -1 \\ -1 & \dots & +1 \\ +1 & \dots & -1 \\ \cdot & & \\ \cdot & & \\ \cdot & & \\ +1 & \dots & -1 \\ K & & \end{array} \right] \Bigg/ \sqrt{K}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

For any  $\epsilon > 0$ , if  $K \approx \log(n/\delta) / \epsilon^2$ , with probability  $1 - \delta$  over draw of  $W$ , for all pairs of data points  $i, j \in \{1, \dots, n\}$ ,

$$(1 - \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

For any  $\epsilon > 0$ , if  $K \approx \log(n/\delta) / \epsilon^2$ , with probability  $1 - \delta$  over draw of  $W$ , for all pairs of data points  $i, j \in \{1, \dots, n\}$ ,

$$(1 - \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

Lets try ...

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

For any  $\epsilon > 0$ , if  $K \approx \log(n/\delta) / \epsilon^2$ , with probability  $1 - \delta$  over draw of  $W$ , for all pairs of data points  $i, j \in \{1, \dots, n\}$ ,

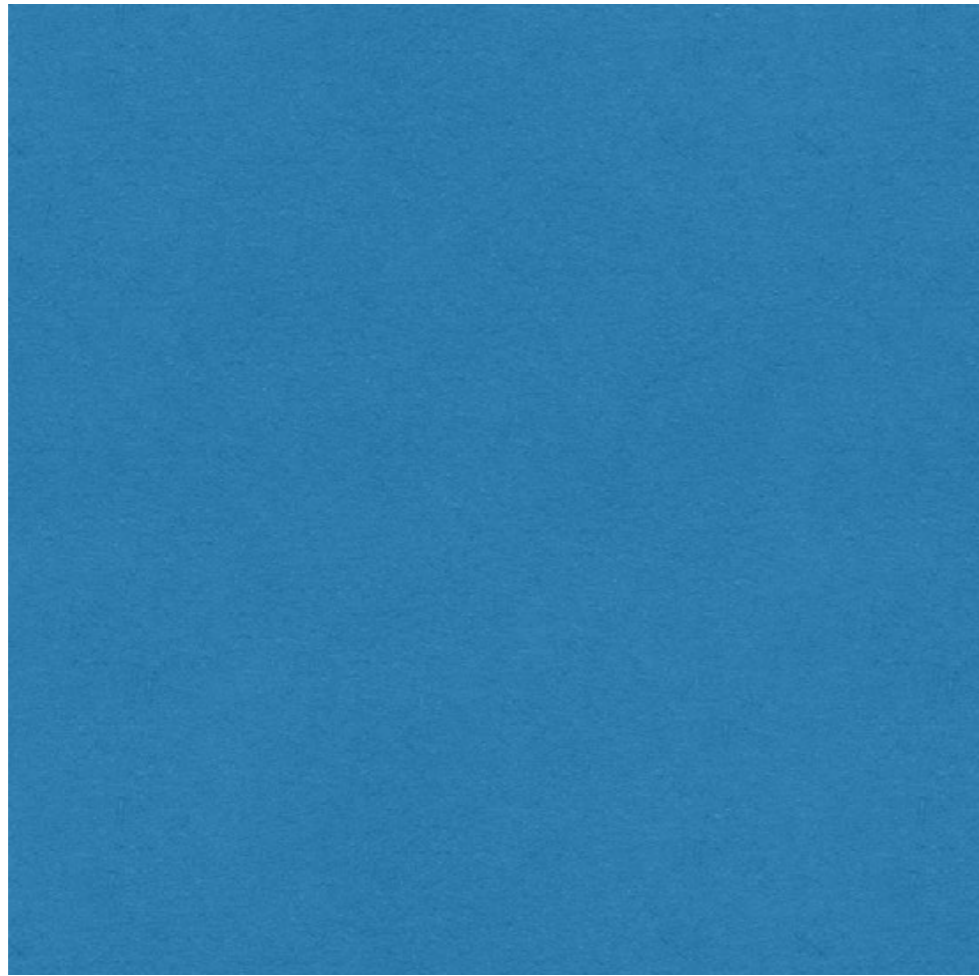
$$(1 - \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

Lets try ...

This is called the Johnson-Lindenstrauss lemma or JL lemma for short.

# WHY IS THIS SO RIDICULOUSLY MAGICAL?

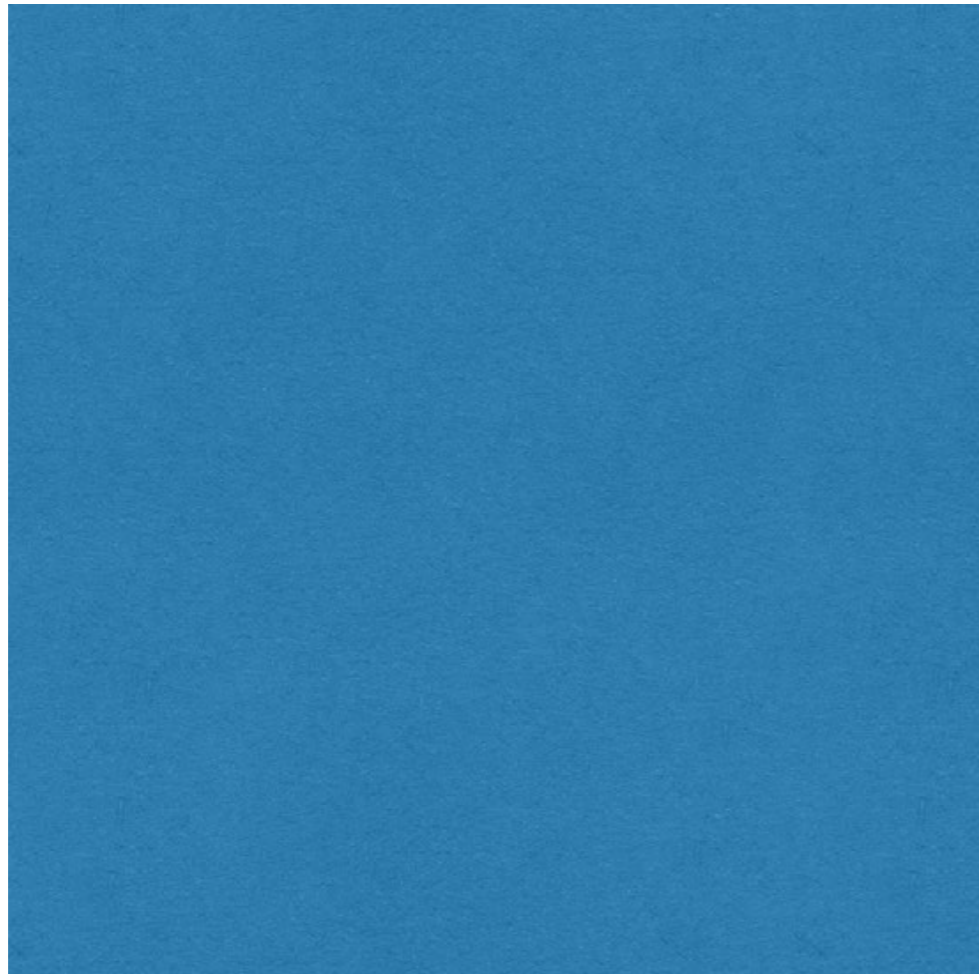
$n =$   
1000



$d = 1000$

# WHY IS THIS SO RIDICULOUSLY MAGICAL?

$n =$   
1000

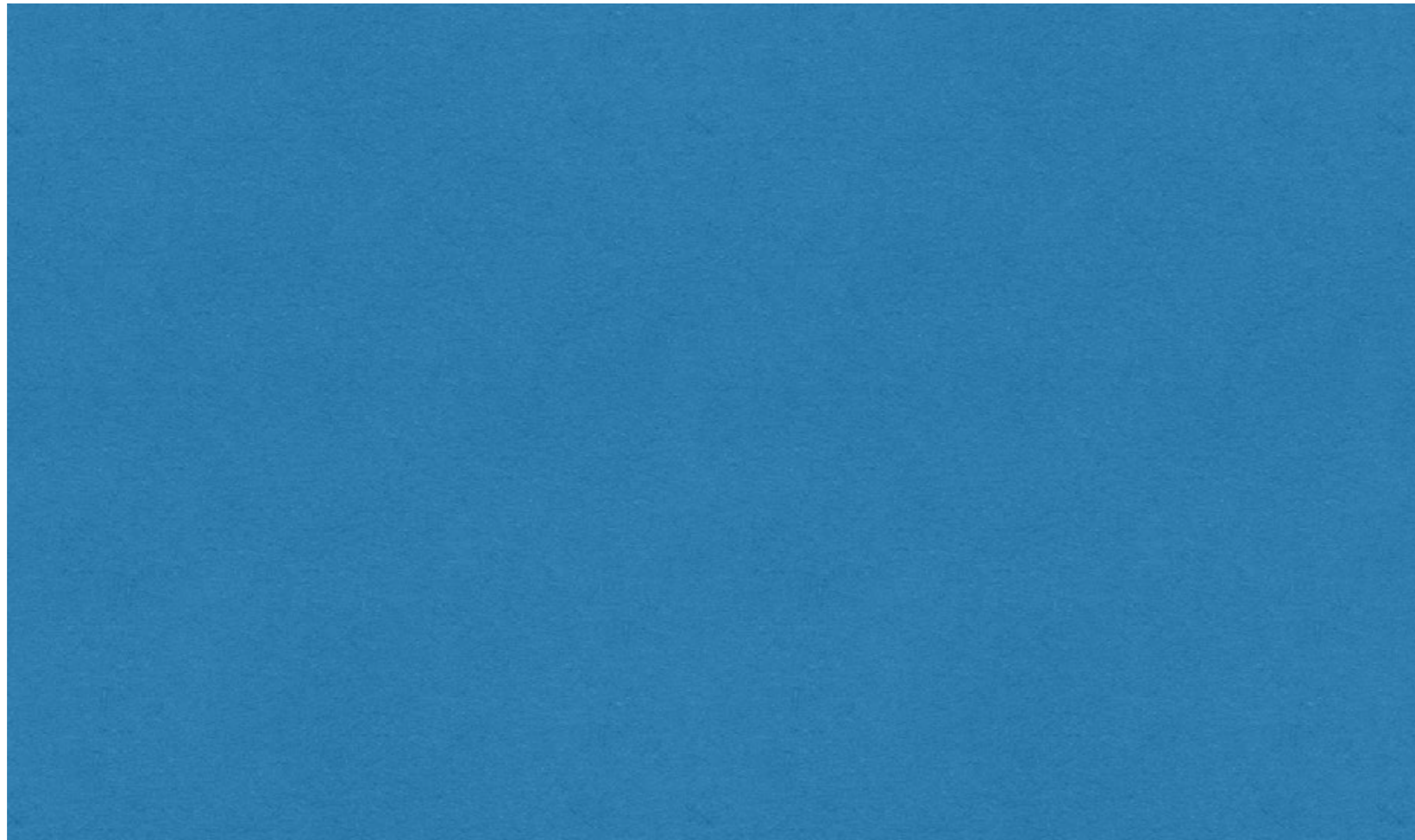


$d = 1000$

If we take  $K = 69.1/\epsilon^2$ , with probability 0.99 distances are preserved to accuracy  $\epsilon$

# WHY IS THIS SO RIDICULOUSLY MAGICAL?

$n =$   
1000



$d = 10000$

If we take  $K = 69.1/\epsilon^2$ , with probability 0.99 distances are preserved to accuracy  $\epsilon$

# WHY IS THIS SO RIDICULOUSLY MAGICAL?

$n =$   
1000

$d = 1000000$

If we take  $K = 69.1/\epsilon^2$ , with probability 0.99 distances are preserved to accuracy  $\epsilon$