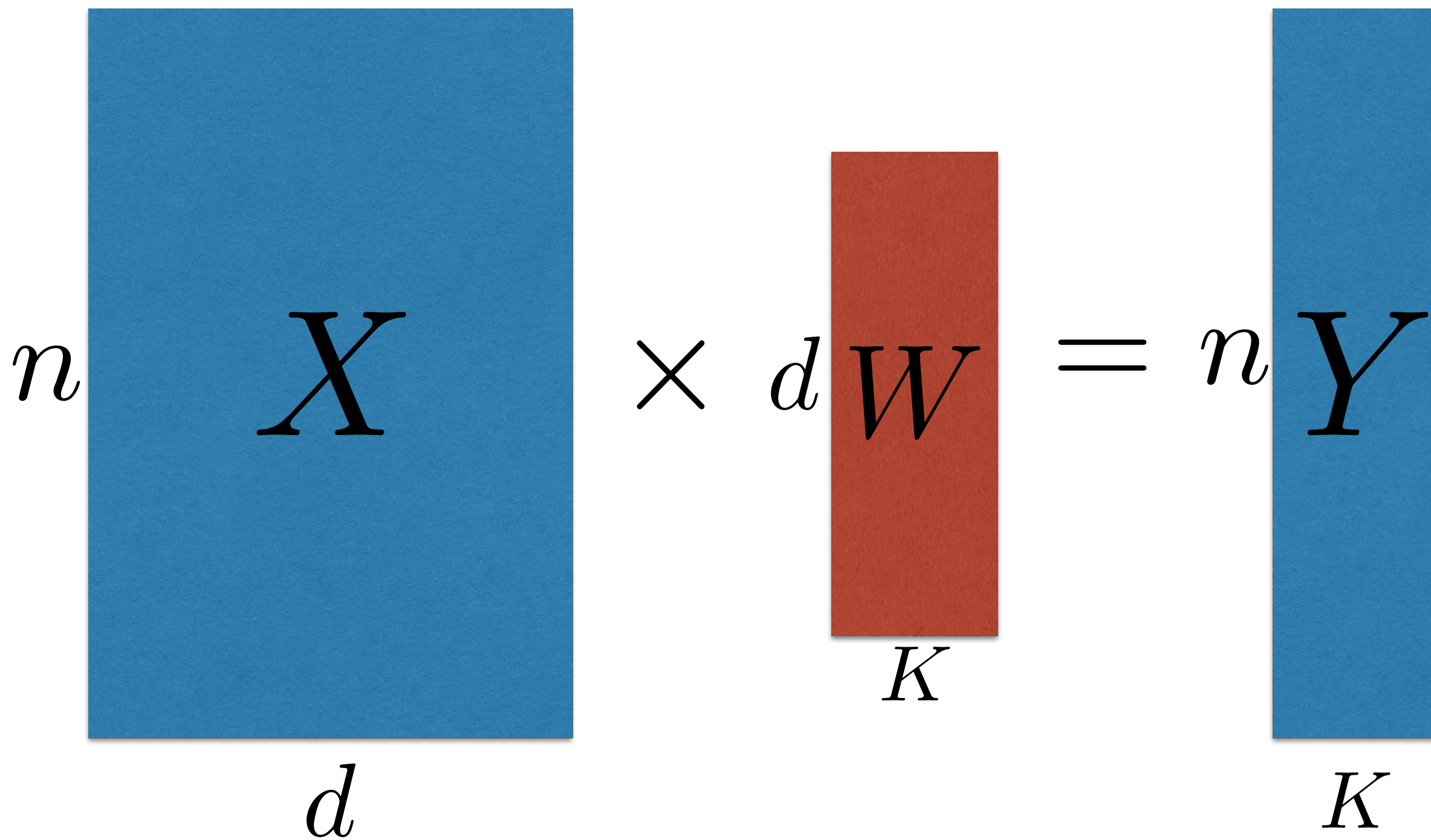


# Machine Learning for Data Science (CS4786)

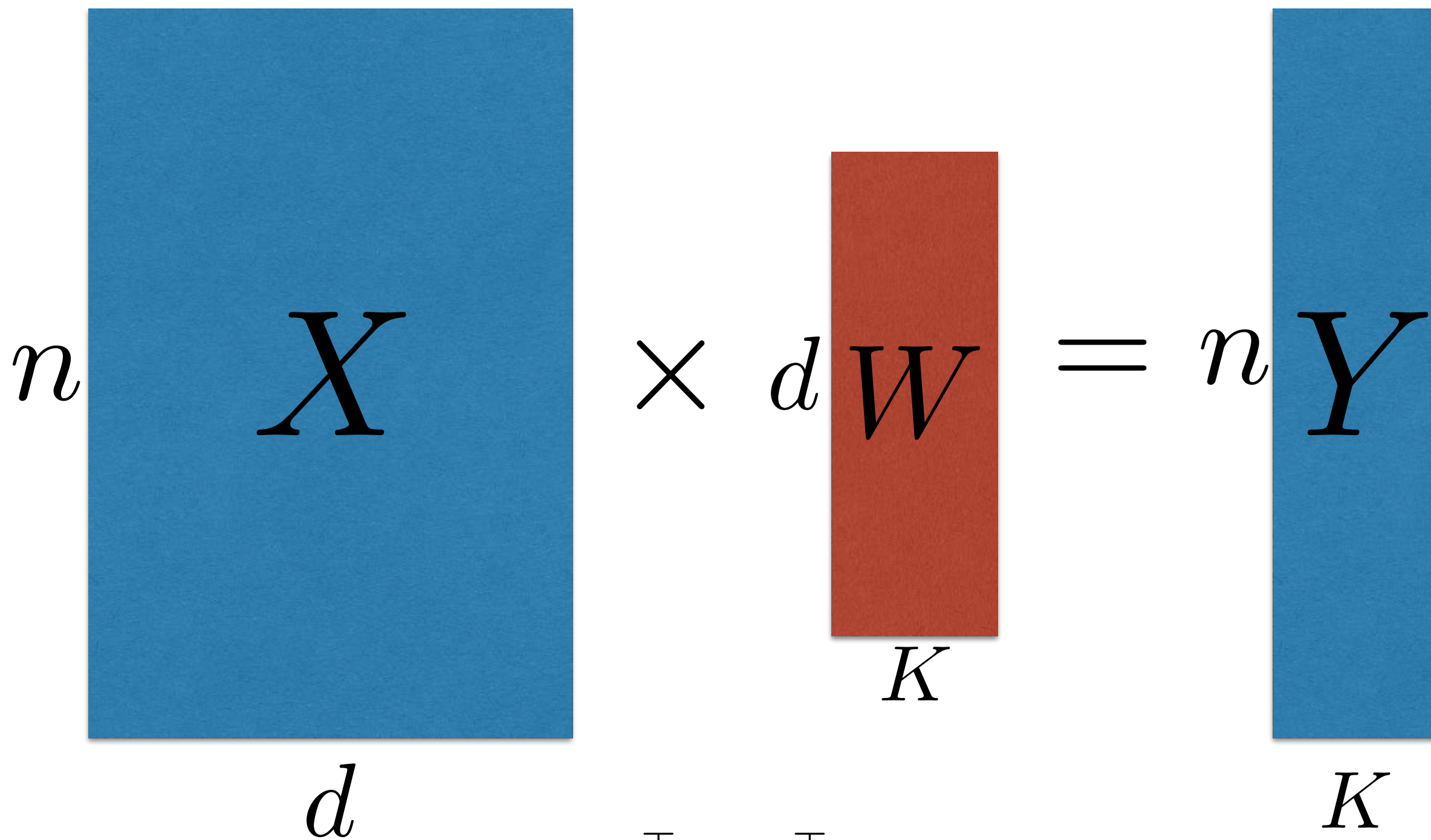
## Lecture 5

Random Projections & Canonical Correlation Analysis

# DIM REDUCTION: LINEAR TRANSFORMATION



# DIM REDUCTION: LINEAR TRANSFORMATION



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

# PRINCIPAL COMPONENT ANALYSIS

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

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

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

# 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?!

$$y_t - y_s = \mathbf{x}_t^\top \mathbf{u} - \mathbf{x}_s^\top \mathbf{u}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$\begin{aligned}y_t - y_s &= \mathbf{x}_t^\top \mathbf{u} - \mathbf{x}_s^\top \mathbf{u} \\ &= (\mathbf{x}_t - \mathbf{x}_s)^\top \mathbf{u}\end{aligned}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$\begin{aligned}y_t - y_s &= \mathbf{x}_t^\top \mathbf{u} - \mathbf{x}_s^\top \mathbf{u} \\ &= (\mathbf{x}_t - \mathbf{x}_s)^\top \mathbf{u} \\ &= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]\end{aligned}$$



# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$\begin{aligned}y_t - y_s &= \mathbf{x}_t^\top \mathbf{u} - \mathbf{x}_s^\top \mathbf{u} \\&= (\mathbf{x}_t - \mathbf{x}_s)^\top \mathbf{u} \\&= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]\end{aligned}$$

$$\mathbb{E}[y_t - y_s] = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbb{E}[\mathbf{u}[k]] = 0$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$y_t - y_s = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$y_t - y_s = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]$$

$$\begin{aligned} & (y_t - y_s)^2 \\ = & \left( \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k] \right)^2 \end{aligned}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$y_t - y_s = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]$$

$$\begin{aligned} & (y_t - y_s)^2 \\ &= \left( \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k] \right)^2 \\ &= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 \mathbf{u}[k]^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \mathbf{u}[k] \cdot \mathbf{u}[j] \end{aligned}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$y_t - y_s = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]$$

$$\begin{aligned} & (y_t - y_s)^2 \\ &= \left( \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k] \right)^2 \\ &= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 \mathbf{u}[k]^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \mathbf{u}[k] \cdot \mathbf{u}[j] \\ &= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \cdot [\mathbf{u}[k] \cdot \mathbf{u}[j]] \end{aligned}$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$y_t - y_s = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]$$

$$\mathbb{E}(y_t - y_s)^2$$

$$= \mathbb{E} \left( \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k] \right)^2$$

$$= \mathbb{E} \left[ \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 \mathbf{u}[k]^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \mathbf{u}[k] \cdot \mathbf{u}[j] \right]$$

$$= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \mathbb{E} [\mathbf{u}[k] \cdot \mathbf{u}[j]]$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

$$y_t - y_s = \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k]$$

$$\mathbb{E}(y_t - y_s)^2$$

$$= \mathbb{E} \left( \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \mathbf{u}[k] \right)^2$$

$$= \mathbb{E} \left[ \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 \mathbf{u}[k]^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \mathbf{u}[k] \cdot \mathbf{u}[j] \right]$$

$$= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 + 2 \sum_{j=1}^d \sum_{k=j+1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k]) \cdot (\mathbf{x}_t[j] - \mathbf{x}_s[j]) \mathbb{E} [\mathbf{u}[k] \cdot \mathbf{u}[j]]$$

$$= \sum_{k=1}^d (\mathbf{x}_t[k] - \mathbf{x}_s[k])^2 = \|\mathbf{x}_t - \mathbf{x}_s\|_2^2$$



# 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

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?!

- 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**

# 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  $i, j \in \{1, \dots, n\}$ ,

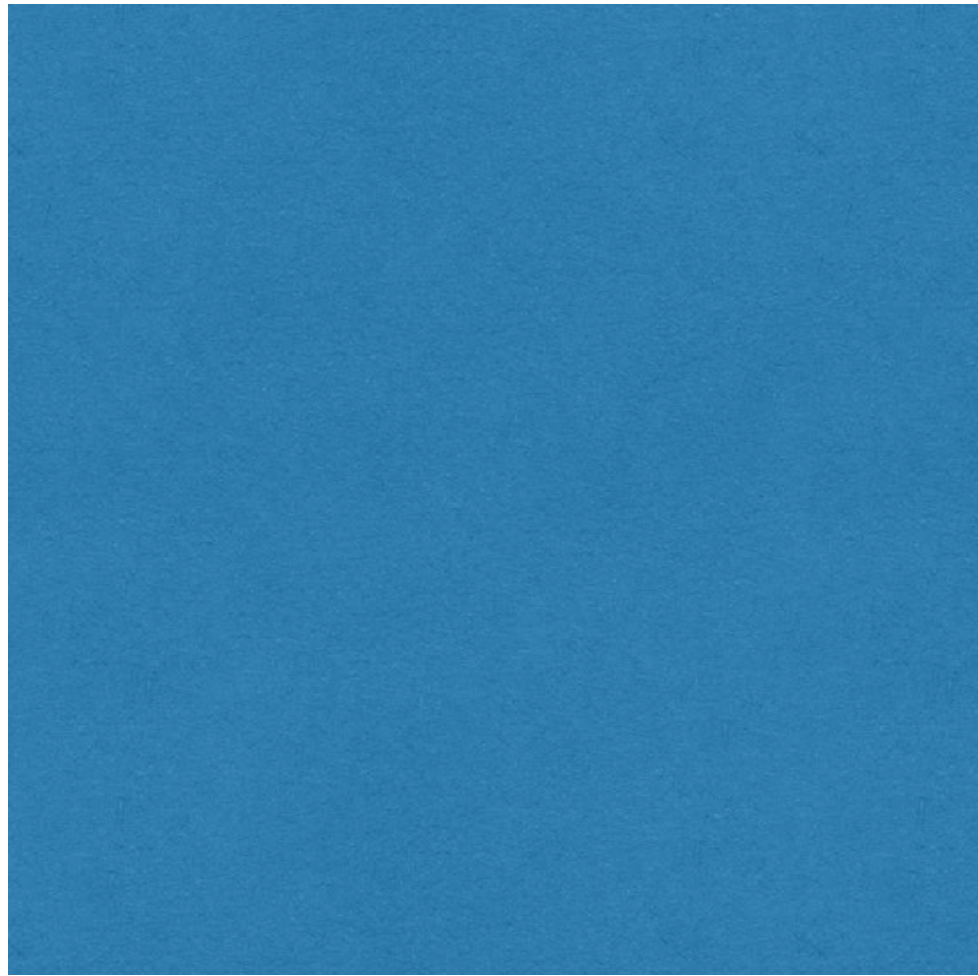
$$(1 - \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

Lets try ...

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

# WHY IS THIS SO RIDICULOUSLY MAGICAL?

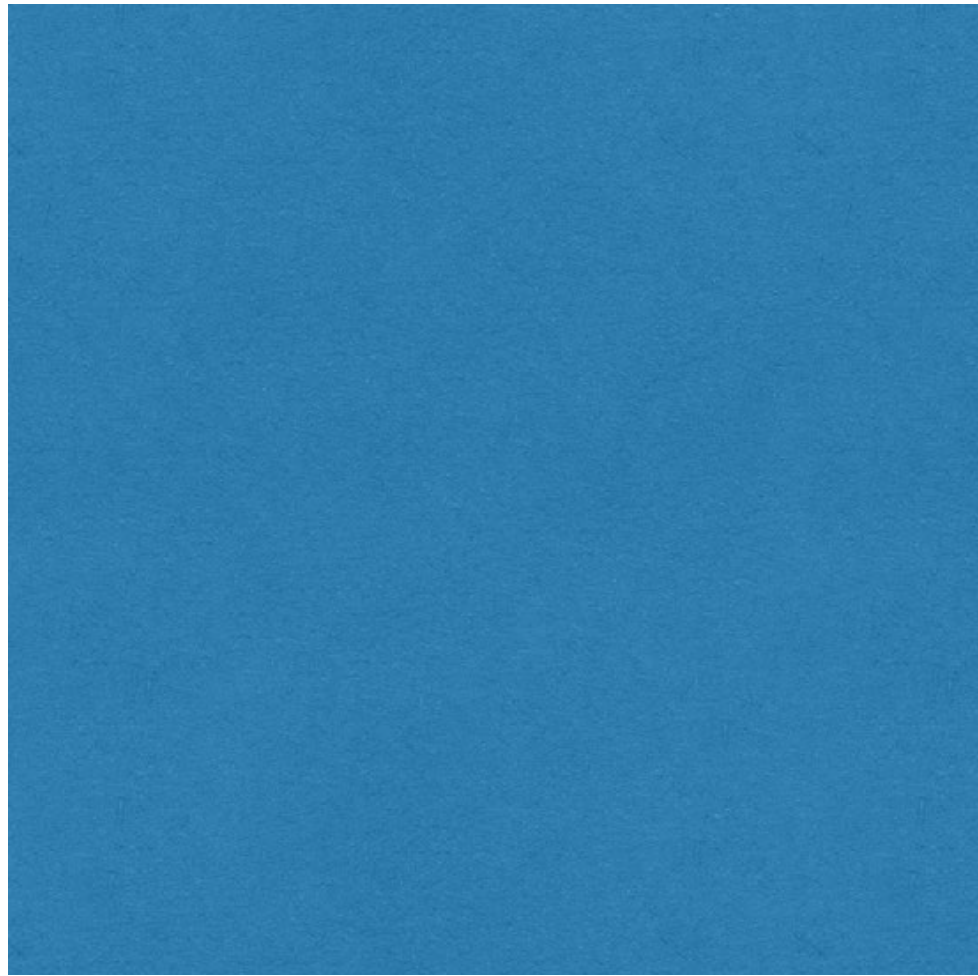
$n =$   
1000



$d = 1000$

# WHY IS THIS SO RIDICULOUSLY MAGICAL?

$n =$   
1000

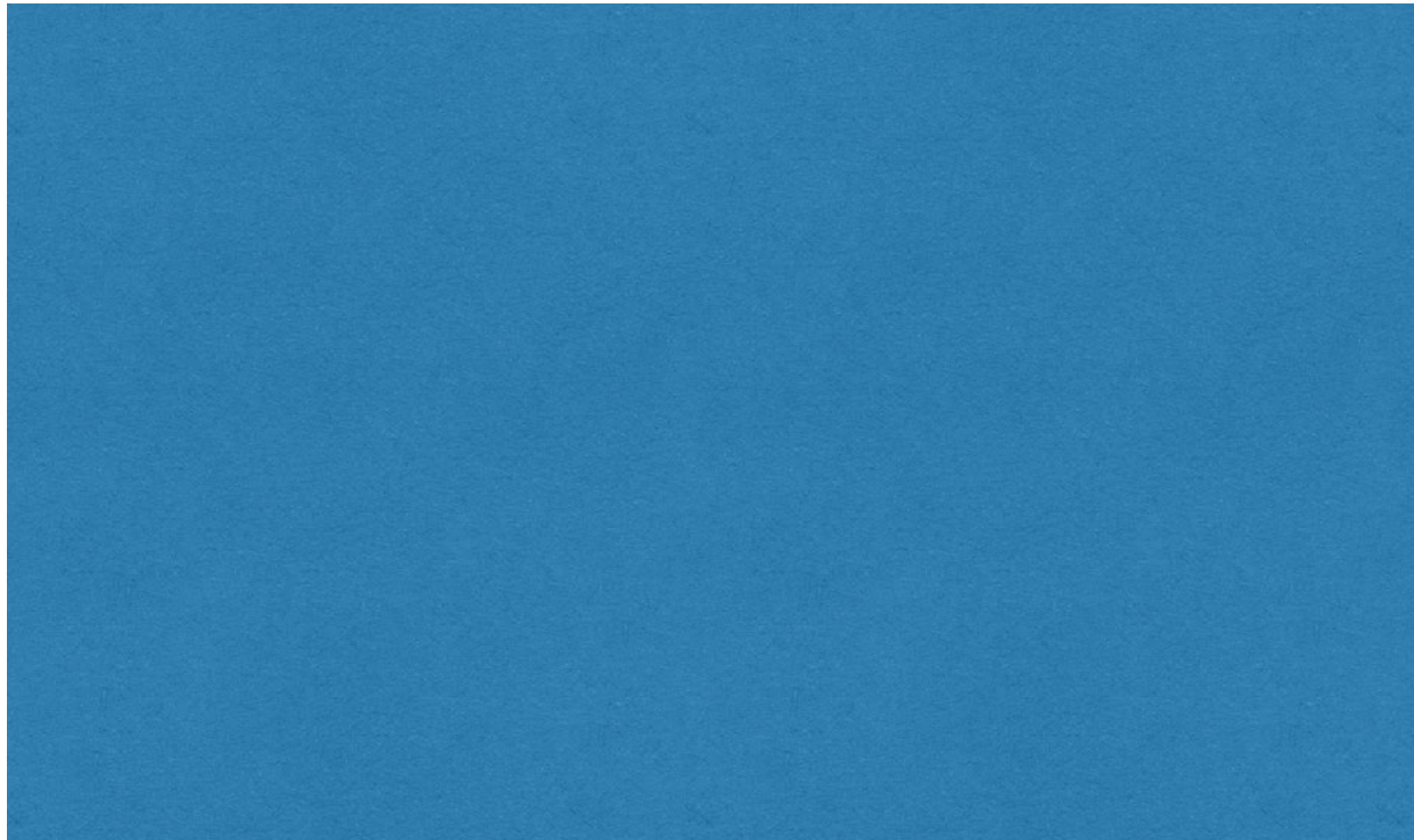


$d = 1000$

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  $\epsilon$

# TWO VIEW DIMENSIONALITY REDUCTION

- Data comes in pairs  $(\mathbf{x}_1, \mathbf{x}'_1), \dots, (\mathbf{x}_n, \mathbf{x}'_n)$  where  $\mathbf{x}_t$ 's are  $d$  dimensional and  $\mathbf{x}'_t$ 's are  $d'$  dimensional
- Goal: Compress say view one into  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , that are  $K$  dimensional vectors
  - Retain information redundant between the two views
  - Eliminate “noise” specific to only one of the views

# Canonical Correlation Analysis

## Analysis

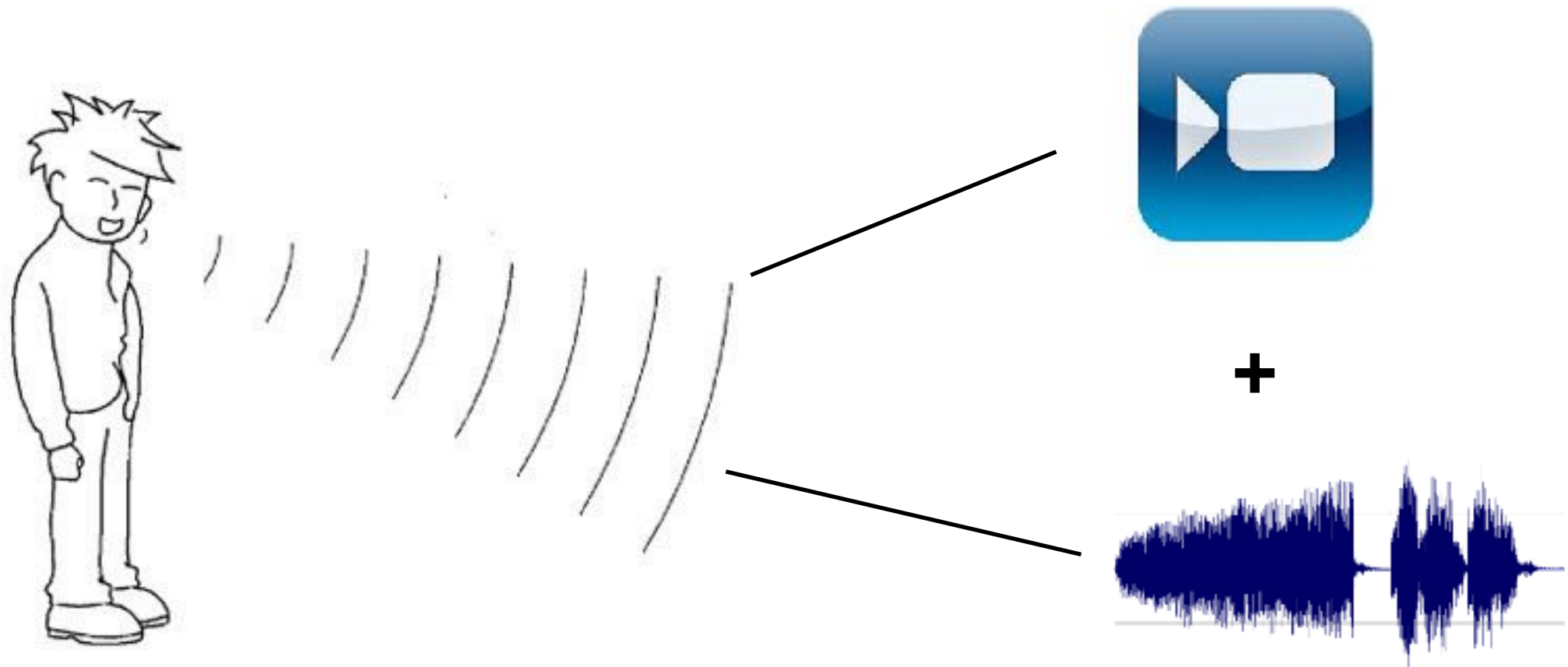


# Canonical Correlation Analysis



Age  
+ Gender  
Candies per week

# EXAMPLE I: SPEECH RECOGNITION



- Audio might have background sounds uncorrelated with video
- Video might have lighting changes uncorrelated with audio
- Redundant information between two views: the speech

# EXAMPLE II: COMBINING FEATURE EXTRACTIONS

- Method A and Method B are both equally good feature extraction techniques
- Concatenating the two features blindly yields large dimensional feature vector with redundancy
- Applying techniques like CCA extracts the key information between the two methods
- Removes extra unwanted information



How do we get the right direction? (say  $K = 1$ )



# How do we get the right direction? (say $K = 1$ )



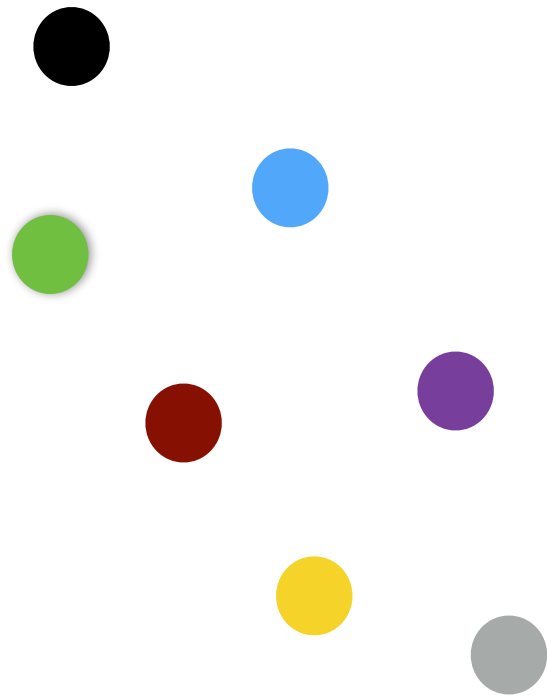
Age

+ Gender

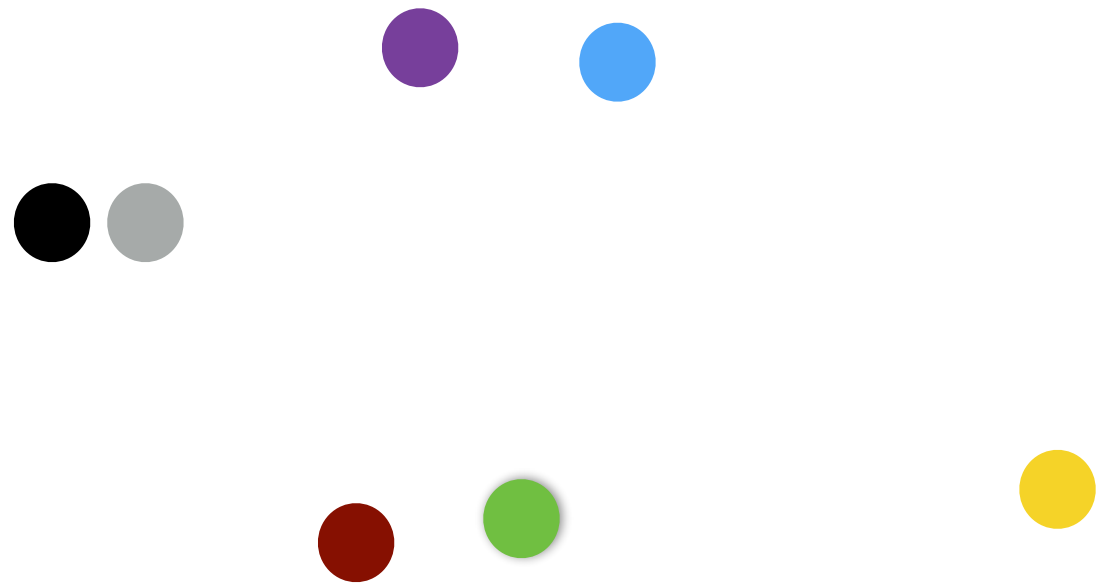
Candies per week



# WHICH DIRECTION TO PICK?

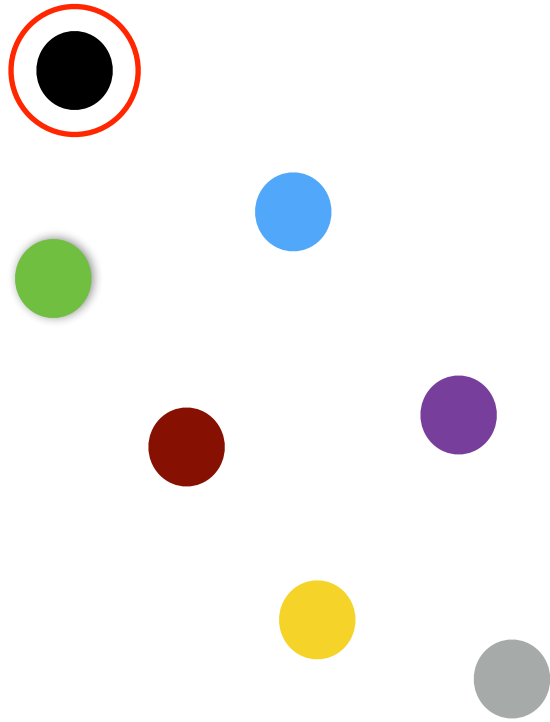


View I

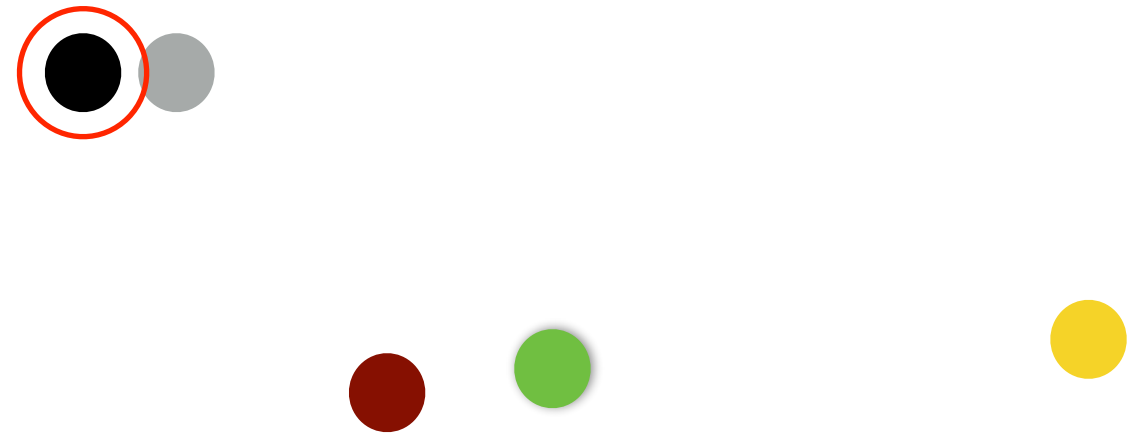


View II

# WHICH DIRECTION TO PICK?

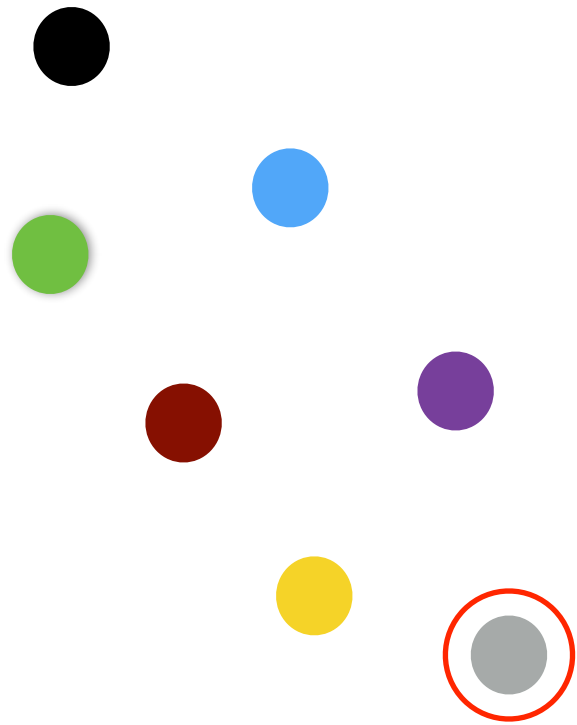


View I

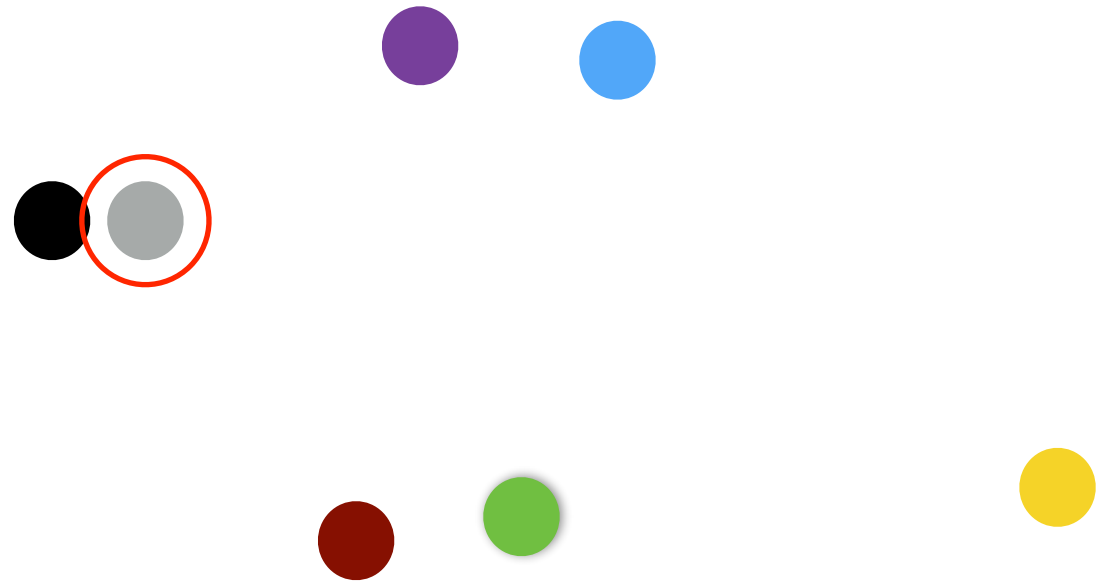


View II

# WHICH DIRECTION TO PICK?

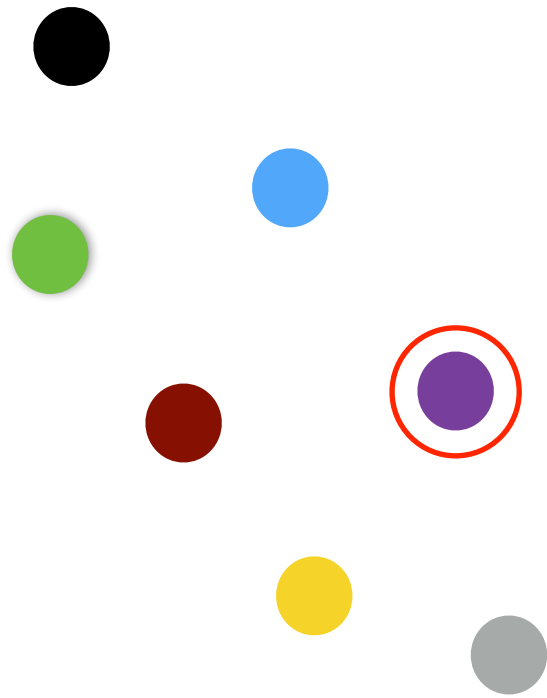


View I

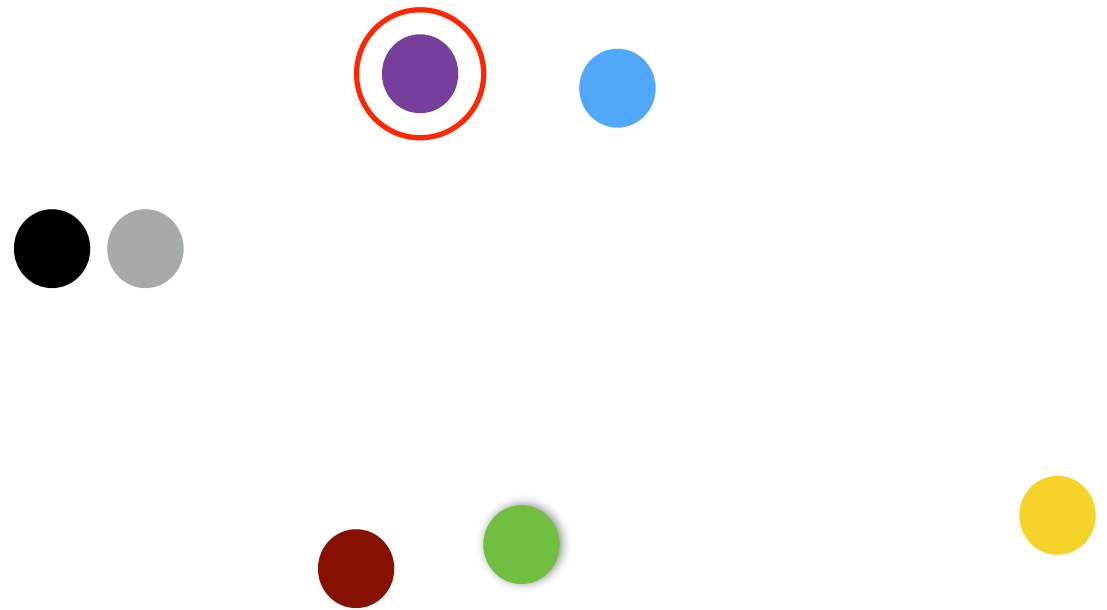


View II

# WHICH DIRECTION TO PICK?



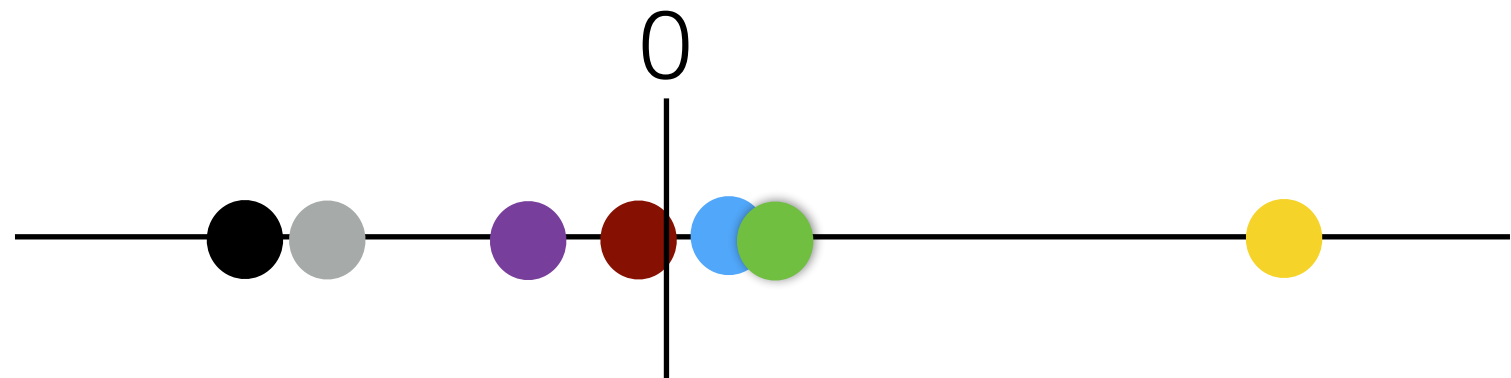
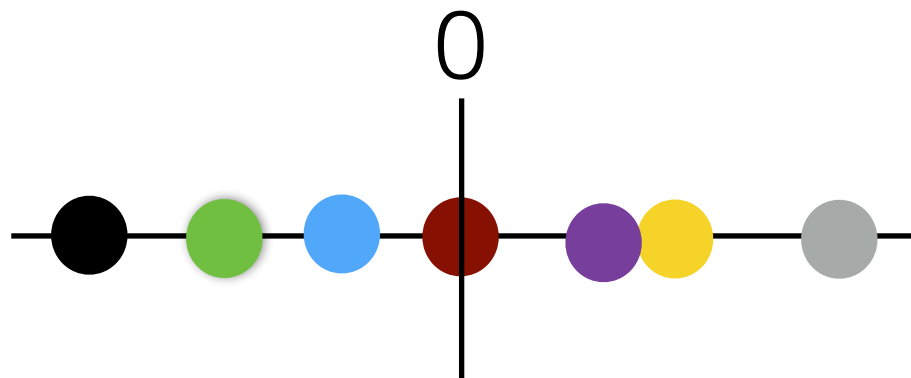
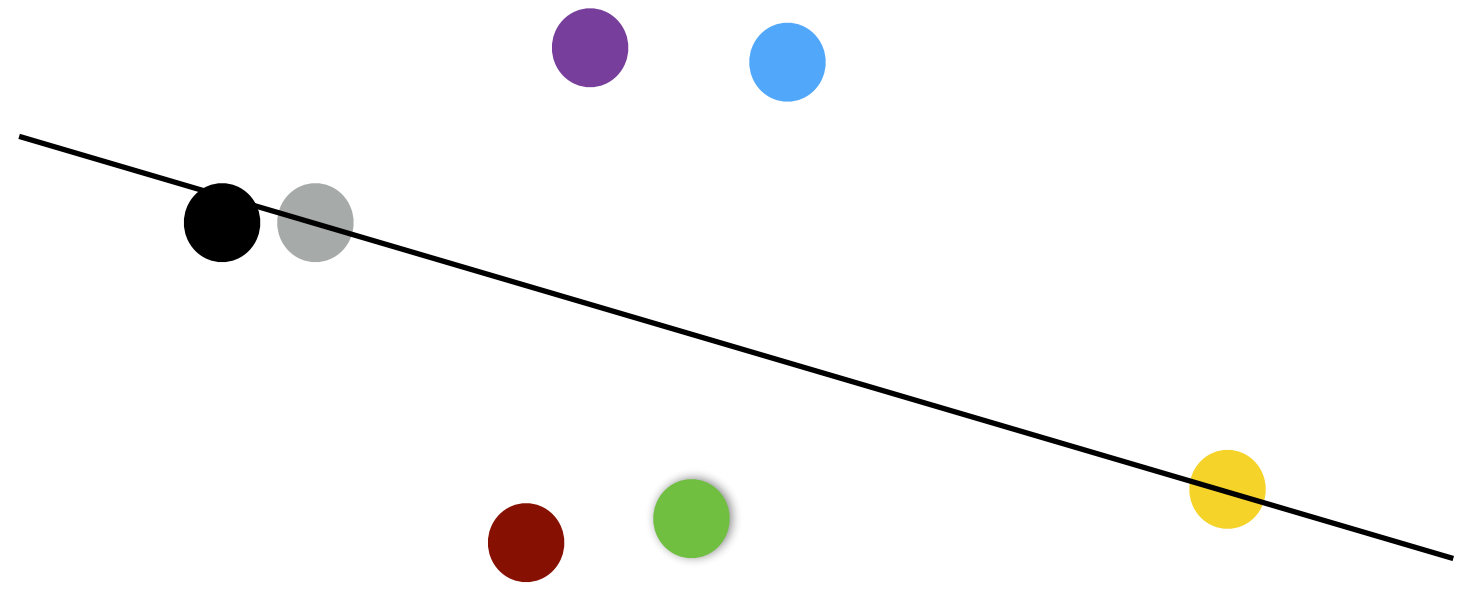
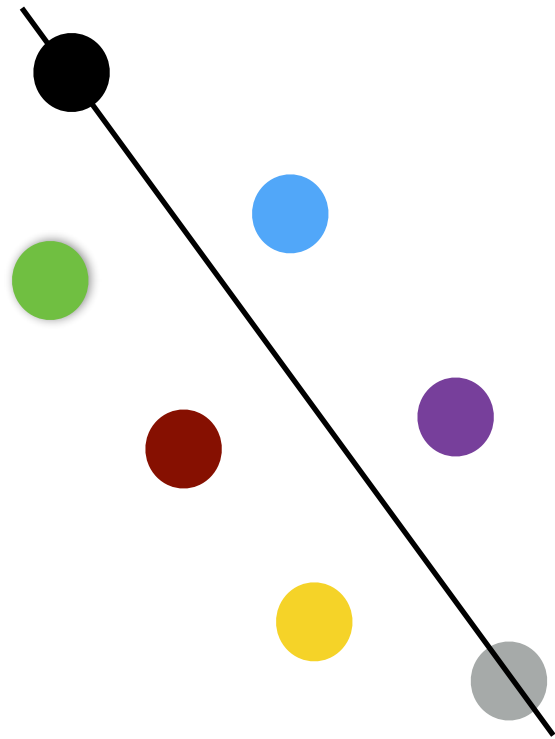
View I



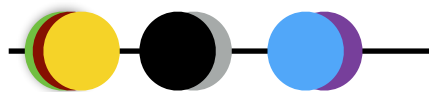
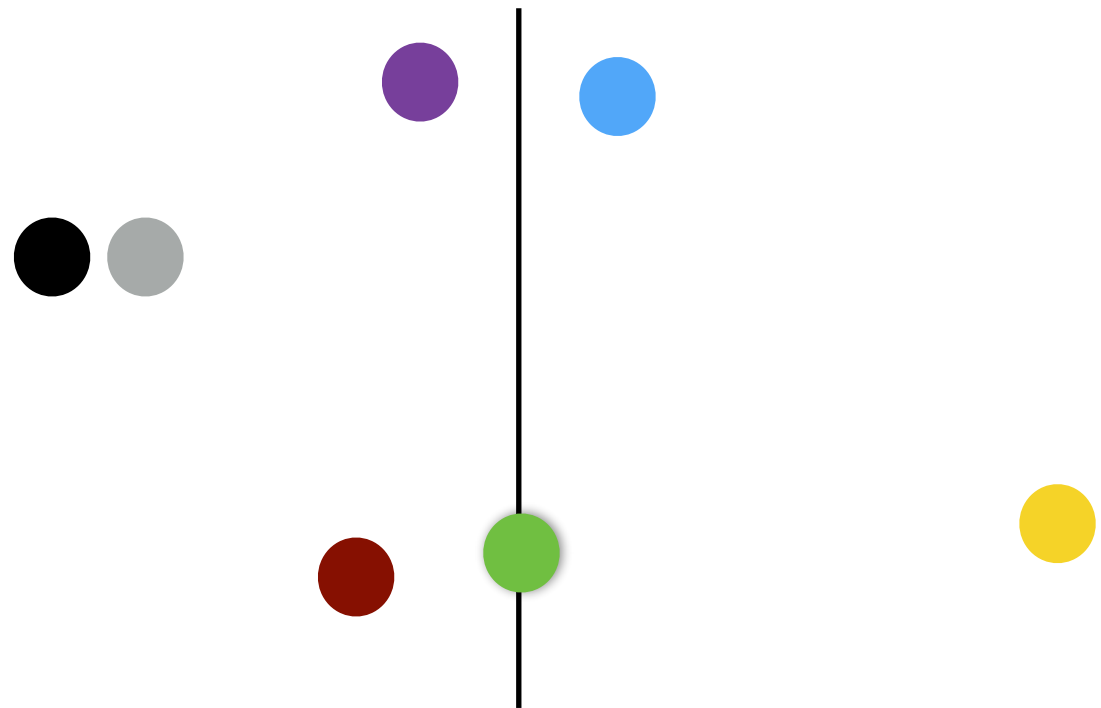
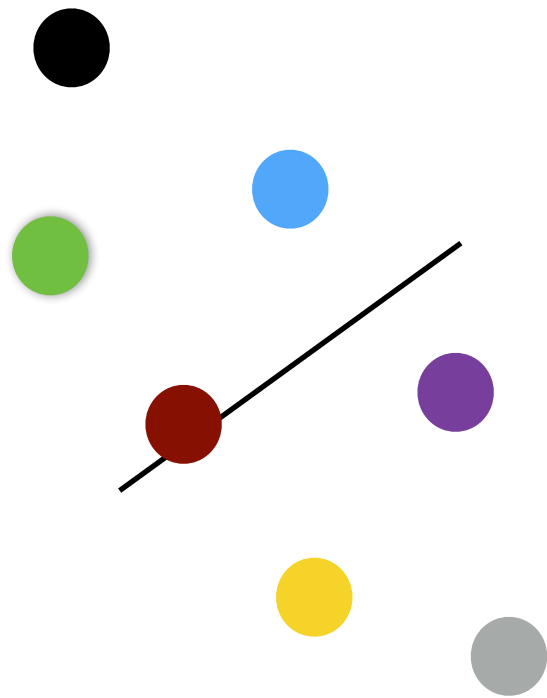
View II

# WHICH DIRECTION TO PICK?

PCA direction



# WHICH DIRECTION TO PICK?



Direction has large covariance

How do we pick the right direction to project to?

# MAXIMIZING CORRELATION COEFFICIENT

- Say  $\mathbf{w}_1$  and  $\mathbf{v}_1$  are the directions we choose to project in views 1 and 2 respectively we want these directions to maximize,

$$\frac{1}{n} \sum_{t=1}^n \left( \mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1] \right) \cdot \left( \mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1] \right)$$

where  $\mathbf{y}_t[1] = \mathbf{w}_1^\top \mathbf{x}_t$  and  $\mathbf{y}'_t[1] = \mathbf{v}_1^\top \mathbf{x}'_t$



What is the problem  
with the above?

How do we get the right direction? (single dimension  $K = 1$ )



# How do we get the right direction? (single dimension $K = 1$ )



Age

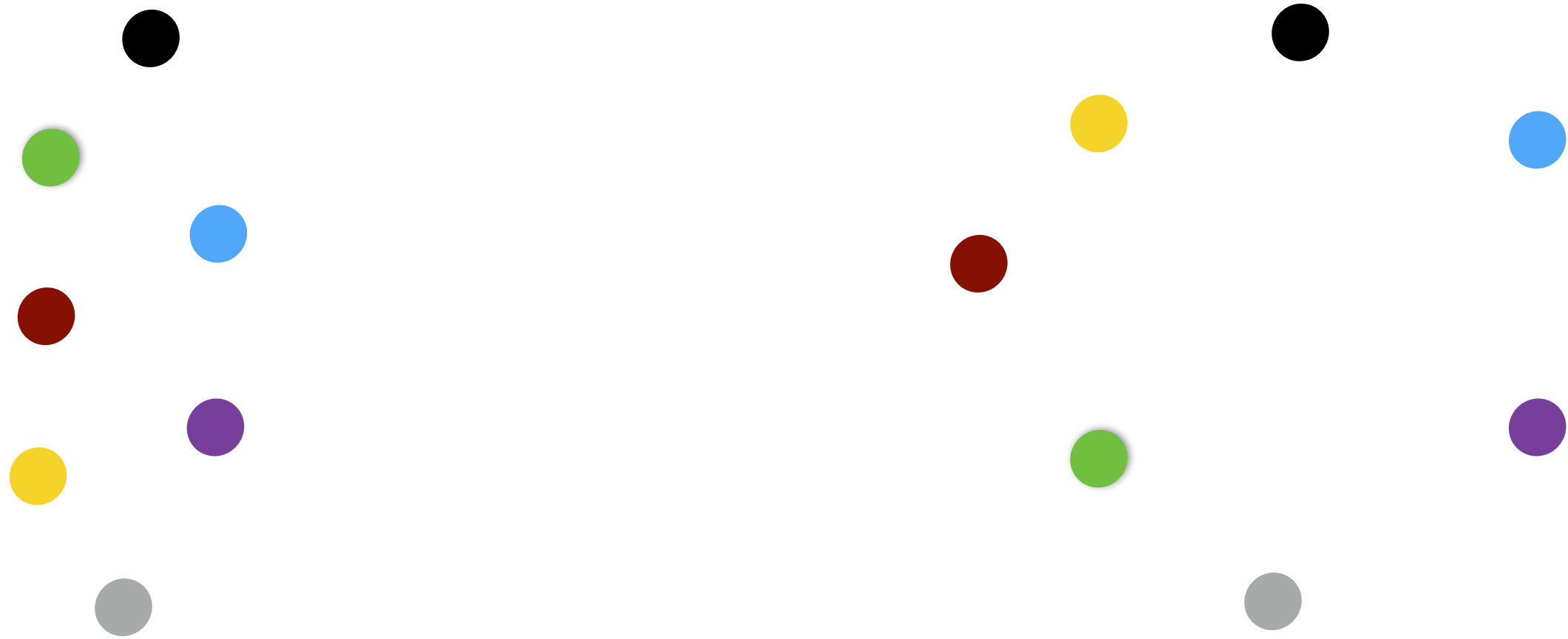
+ Gender

**Candies per week**

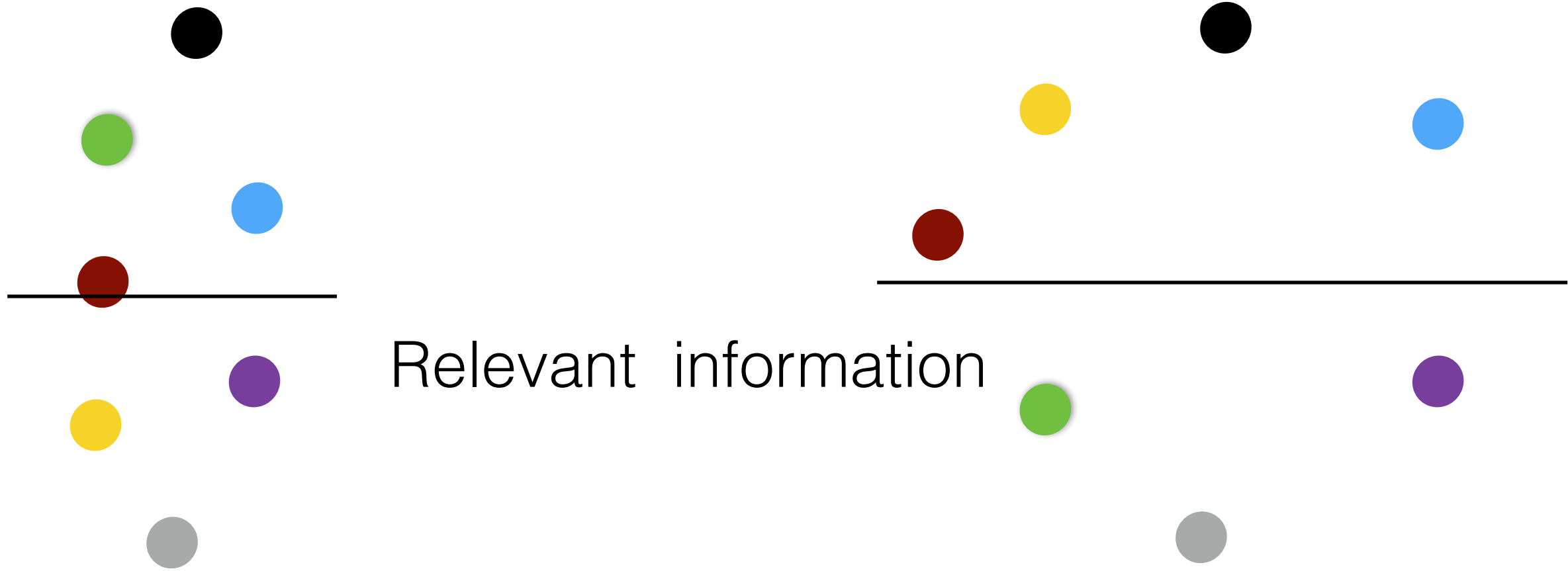
Favorite Cartoon

# WHY NOT MAXIMIZE COVARIANCE

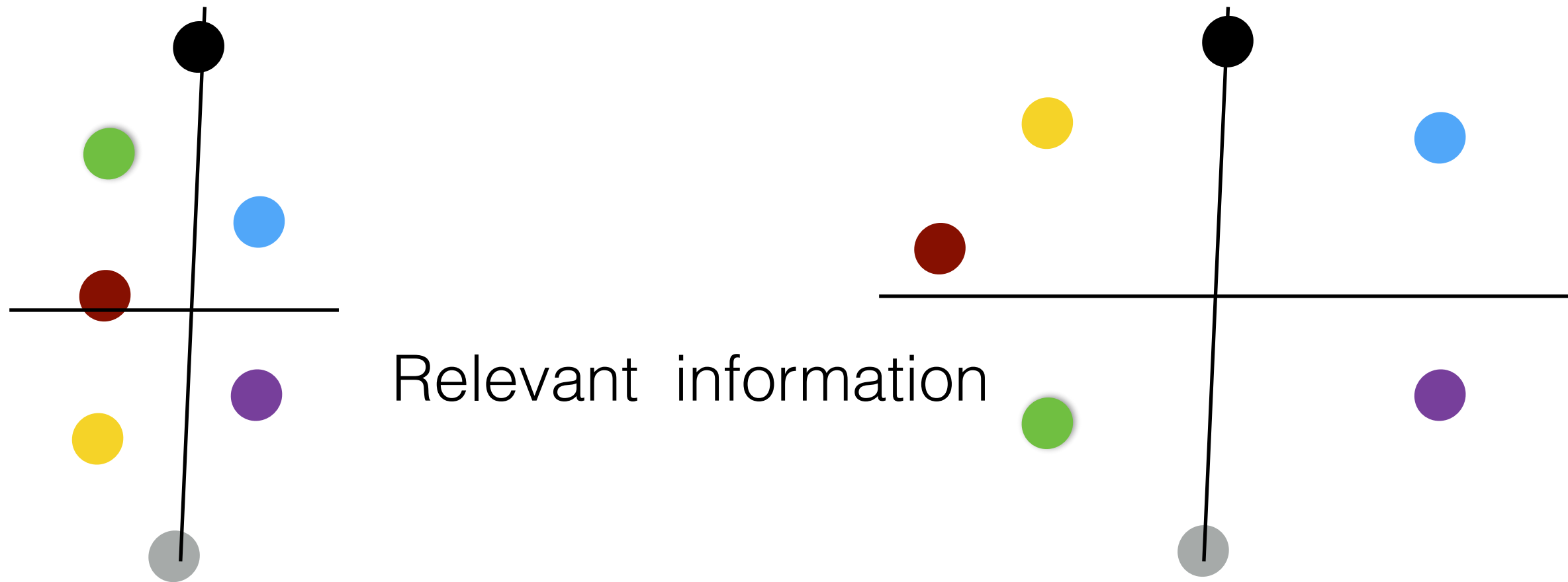
# WHY NOT MAXIMIZE COVARIANCE



# WHY NOT MAXIMIZE COVARIANCE



# WHY NOT MAXIMIZE COVARIANCE



**Say covariance in some coordinate just happens to be  $> 0$**

Scaling up this coordinate we can blow up covariance