# Machine Learning for Data Science (CS4786)
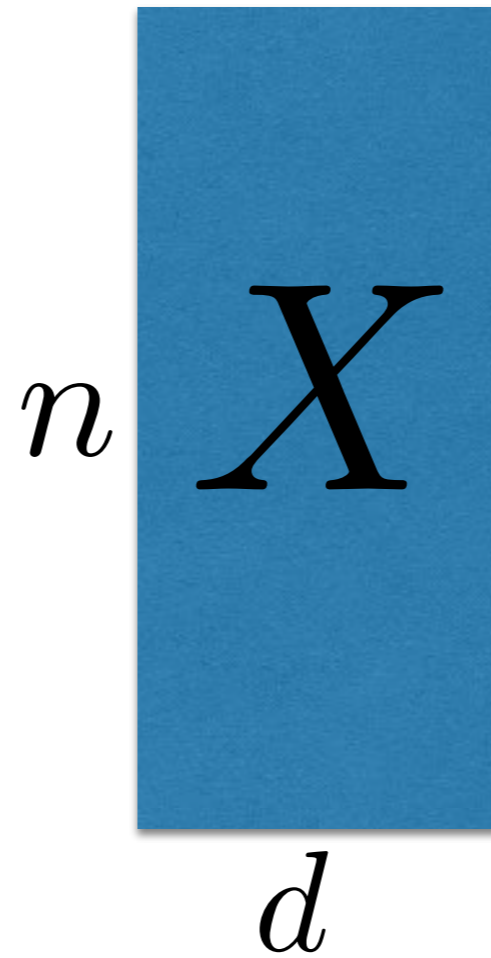# Lecture 11

Random Projections & Canonical Correlation Analysis

Course Webpage :
http://www.cs.cornell.edu/Courses/cs4786/2017fa/

$$d \boxed{X^\top}_n \times n \boxed{X}_d \Big/ n = d \boxed{\Sigma}^d$$

$$d \boxed{\frac{X^\top}{n}} \times n \boxed{\frac{X}{d}} \Big/ n = d \boxed{\Sigma}^d$$

$$d \boxed{W}_K = \mathrm{Eigs}\left( \boxed{\Sigma}, K \right)$$

$n$

$X$

$d$

$n$ $X$

$d$

$\mathrm{SVD}(X)$

$d$

$K$

$V^\top$

$n$ $U$ $\times$ $\times$ $d$

$n$

$$X$$

- $d$ 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?

$$Y = X \times \begin{bmatrix} +1 & \dots & -1 \\ -1 & \dots & +1 \\ +1 & \dots & -1 \\ & \cdot & \\ & \cdot & \\ & \cdot & \\ +1 & \dots & -1 \end{bmatrix} d \Big/ \sqrt{K}$$

$$K$$

- What does "it works" even mean?

- 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.

- 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 $i, j \in \{1, \ldots, n\}$,

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

Say $K = 1$. Consider any vector $\tilde{\mathbf{x}} \in \mathbb{R}^d$ and let $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} W$. Note that

Say $K = 1$. Consider any vector $\tilde{\mathbf{x}} \in \mathbb{R}^d$ and let $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} W$. Note that

$$\tilde{\mathbf{y}}^2 = \left( \sum_{i=1}^{d} W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2$$

Say $K = 1$. Consider any vector $\tilde{\mathbf{x}} \in \mathbb{R}^d$ and let $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} W$. Note that

$$\tilde{\mathbf{y}}^2 = \left( \sum_{i=1}^{d} W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2$$

$$= \sum_{i=1}^{d} \left( W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2 + 2 \sum_{i' > i} \left( W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right) \left( W[i', 1] \cdot \tilde{\mathbf{x}}[i'] \right)$$

Say $K = 1$. Consider any vector $\tilde{\mathbf{x}} \in \mathbb{R}^d$ and let $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} \, W$. Note that

$$
\begin{aligned}
\tilde{\mathbf{y}}^2 &= \left( \sum_{i=1}^{d} W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2 \\
&= \sum_{i=1}^{d} \left( W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2 + 2 \sum_{i'>i} \left( W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right) \left( W[i', 1] \cdot \tilde{\mathbf{x}}[i'] \right) \\
&= \sum_{i=1}^{d} W^2[i, 1] \tilde{\mathbf{x}}^2[i] + \sum_{i'>i} \left( W[i, 1] \cdot W[i', 1] \right) \cdot \left( \tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i'] \right)
\end{aligned}
$$

# WHY SHOULD RANDOM PROJECTIONS EVEN WORK?!

Say $K = 1$. Consider any vector $\tilde{\mathbf{x}} \in \mathbb{R}^d$ and let $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} W$. Note that

$$\tilde{\mathbf{y}}^2 = \left( \sum_{i=1}^{d} W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2$$

$$= \sum_{i=1}^{d} \left( W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right)^2 + 2 \sum_{i' > i} \left( W[i, 1] \cdot \tilde{\mathbf{x}}[i] \right) \left( W[i', 1] \cdot \tilde{\mathbf{x}}[i'] \right)$$

$$= \sum_{i=1}^{d} W^2[i, 1] \tilde{\mathbf{x}}^2[i] + \sum_{i' > i} \left( W[i, 1] \cdot W[i', 1] \right) \cdot \left( \tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i'] \right)$$

However $W^2[i, 1] = 1/K = 1$ when $K = 1$

$$= \sum_{i=1}^{d} \tilde{\mathbf{x}}^2[i] + \sum_{i' > i} \left( W[i, 1] \cdot W[i', 1] \right) \cdot \left( \tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i'] \right)$$

Hence,

$$\mathbb{E}\left[\tilde{\mathbf{y}}^2\right] = \sum_{i=1}^{d} \tilde{\mathbf{x}}^2[i] + \sum_{i' > i} \mathbb{E}\left[W[i, 1] \cdot W[i', 1]\right] \cdot \left(\tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i']\right)$$

Hence,

$$\mathbb{E}\left[\tilde{\mathbf{y}}^2\right] = \sum_{i=1}^{d} \tilde{\mathbf{x}}^2[i] + \sum_{i'>i} \mathbb{E}\left[W[i,1] \cdot W[i',1]\right] \cdot \left(\tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i']\right)$$

However $W[i,1]$ and $W[i',1]$ are independent and so

$$\mathbb{E}\left[W[i,1] \cdot W[i',1]\right] = \mathbb{E}\left[W[i,1]\right] \cdot \mathbb{E}\left[W[i',1]\right] = 0$$

Hence,

$$\mathbb{E}[\tilde{\mathbf{y}}^2] = \sum_{i=1}^{d} \tilde{\mathbf{x}}^2[i] + \sum_{i'>i} \mathbb{E}[W[i,1] \cdot W[i',1]] \cdot (\tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i'])$$

However $W[i,1]$ and $W[i',1]$ are independent and so

$$\mathbb{E}[W[i,1] \cdot W[i',1]] = \mathbb{E}[W[i,1]] \cdot \mathbb{E}[W[i',1]] = 0$$

Using this we conclude that

$$\mathbb{E}[\tilde{\mathbf{y}}^2] = \sum_{i=1}^{d} \tilde{\mathbf{x}}^2[i] = \|\tilde{\mathbf{x}}\|^2$$

Hence,

$$\mathbb{E}\left[\left|\tilde{\mathbf{y}}\right|^2\right] = \|\tilde{\mathbf{x}}\|_2^2$$

Hence,

$$\mathbb{E}\left[|\tilde{\mathbf{y}}|^2\right] = \|\tilde{\mathbf{x}}\|_2^2$$

If we let $\tilde{\mathbf{x}} = \mathbf{x}_s - \mathbf{x}_t$ then

$$\tilde{\mathbf{y}} = \tilde{\mathbf{x}}W = \mathbf{x}_s W - \mathbf{x}_t W = \mathbf{y}_s - \mathbf{y}_t$$

Hence,

$$\mathbb{E}\left[|\tilde{\mathbf{y}}|^2\right] = \|\tilde{\mathbf{x}}\|_2^2$$

If we let $\tilde{\mathbf{x}} = \mathbf{x}_s - \mathbf{x}_t$ then

$$\tilde{\mathbf{y}} = \tilde{\mathbf{x}}W = \mathbf{x}_s W - \mathbf{x}_t W = \mathbf{y}_s - \mathbf{y}_t$$

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

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

Hence,

$$\mathbb{E}\left[|\tilde{\mathbf{y}}|^2\right] = \|\tilde{\mathbf{x}}\|_2^2$$

If we let $\tilde{\mathbf{x}} = \mathbf{x}_s - \mathbf{x}_t$ then

$$\tilde{\mathbf{y}} = \tilde{\mathbf{x}}W = \mathbf{x}_s W - \mathbf{x}_t W = \mathbf{y}_s - \mathbf{y}_t$$

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

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

Lets try this in Matlab …

- Setting $K$ large is like getting $K$ samples.

- Setting $K$ large is like getting $K$ samples.
- Specifically since we take W to be random signs normalized by $\sqrt{K}$, for each $j \in [K]$, for any $\tilde{\mathbf{x}}$ if $\tilde{\mathbf{y}} = \tilde{\mathbf{x}}\, W$, then

$$\mathbb{E}\big[\tilde{\mathbf{y}}^2[j]\big] = \|\tilde{\mathbf{x}}\|_2^2 / K$$

- Setting $K$ large is like getting $K$ samples.
- Specifically since we take W to be random signs normalized by $\sqrt{K}$, for each $j \in [K]$, for any $\tilde{\mathbf{x}}$ if $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} \, W$, then

$$\mathbb{E}\left[\tilde{\mathbf{y}}^2[j]\right] = \|\tilde{\mathbf{x}}\|_2^2 / K$$

Hence we can conclude that

$$\mathbb{E}\left[\sum_{j=1}^{K} \tilde{\mathbf{y}}^2[j]\right] = \sum_{j=1}^{K} \mathbb{E}\left[\tilde{\mathbf{y}}^2[j]\right] = \sum_{j=1}^{K} \frac{\|\tilde{\mathbf{x}}\|_2^2}{K} = \|\tilde{\mathbf{x}}\|_2^2$$

- Setting $K$ large is like getting $K$ samples.
- Specifically since we take W to be random signs normalized by $\sqrt{K}$, for each $j \in [K]$, for any $\tilde{\mathbf{x}}$ if $\tilde{\mathbf{y}} = \tilde{\mathbf{x}} \, W$, then

$$\mathbb{E}\left[\tilde{\mathbf{y}}^2[j]\right] = \|\tilde{\mathbf{x}}\|_2^2 / K$$

Hence we can conclude that

$$\mathbb{E}\left[\sum_{j=1}^{K} \tilde{\mathbf{y}}^2[j]\right] = \sum_{j=1}^{K} \mathbb{E}\left[\tilde{\mathbf{y}}^2[j]\right] = \sum_{j=1}^{K} \frac{\|\tilde{\mathbf{x}}\|_2^2}{K} = \|\tilde{\mathbf{x}}\|_2^2$$

This is like taking an average of $K$ independent measurements whose expectations are $\|\tilde{\mathbf{x}}\|_2^2$

For large $K$, not only true in expectation but also with high probability

For large $K$, not only true in expectation but also with high probability

For any $\epsilon > 0$, if $K \approx \log\left(n/\delta\right)/\epsilon^2$, with probability $1 - \delta$ over draw of $W$, for all pairs of data points $i, j \in \{1, \ldots, n\}$,
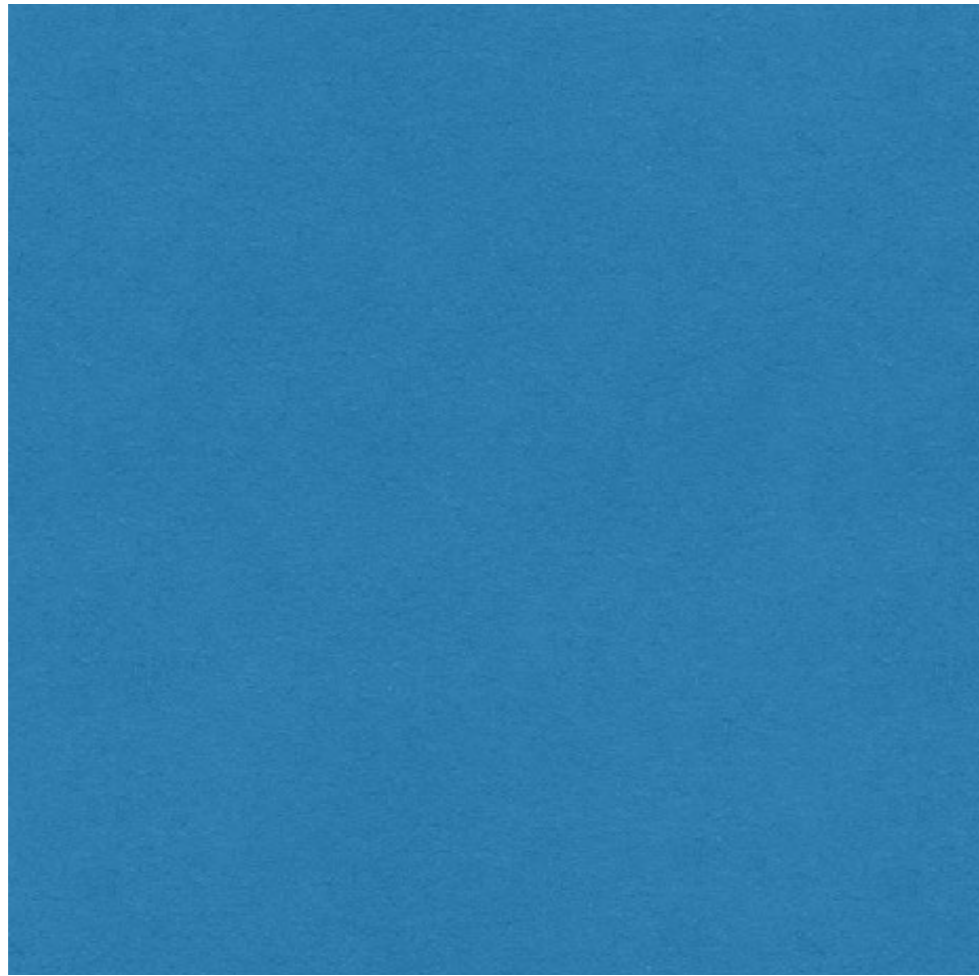
$$(1 - \epsilon)\left\|\mathbf{y}_i - \mathbf{y}_j\right\|_2^2 \leq \left\|\mathbf{x}_i - \mathbf{x}_j\right\|_2 \leq (1 + \epsilon)\left\|\mathbf{y}_i - \mathbf{y}_j\right\|_2^2$$

For large $K$, not only true in expectation but also with high probability

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 $i, j \in \{1, \ldots, n\}$,

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

Lets try on Matlab ...

For large $K$, not only true in expectation but also with high probability

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 $i, j \in \{1, \ldots, n\}$,

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

Lets try on Matlab ...

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

n=
1000

d = 1000

n=
1000

d = 1000

If we take $K = 69.1/\epsilon^2$, with probability

0.99 distances are preserved to accuracy $\epsilon$

n=
1000

d = 10000

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

n=
1000

d = 1000000

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

- Data comes in pairs $(\mathbf{x}_1, \mathbf{x}_1'), \ldots, (\mathbf{x}_n, \mathbf{x}_n')$ where $\mathbf{x}_t$'s are $d$ dimensional and $\mathbf{x}_t''$'s are $d'$ dimensional

- Goal: Compress say view one into $\mathbf{y}_1, \ldots, \mathbf{y}_n$, that are $K$ dimensional vectors

  - Retain information redundant between the two views

  - Eliminate "noise" specific to only one of the views

# Canonical Correlation Analysis

# Canonical Correlation Analysis



Age

+ Gender

Angle

# Canonical Correlation Analysis



Age

+ Gender

Angle

- Audio might have background sounds uncorrelated with video

- Video might have lighting changes uncorrelated with audio

- Redundant information between two views: the speech

- 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)



\+   Age

Gender

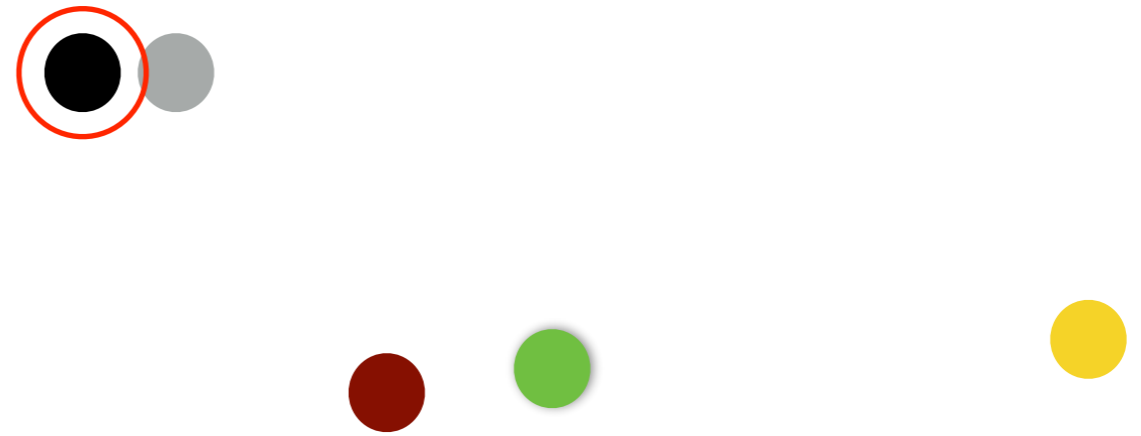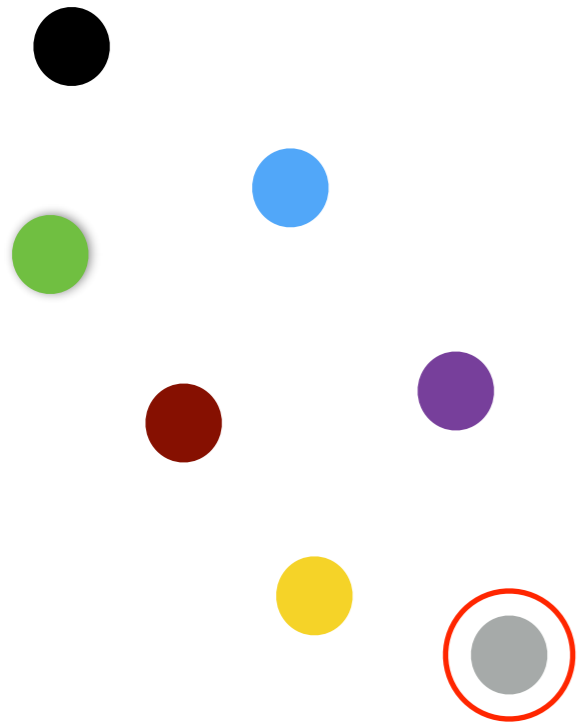<span style="color:red">Angle</span>

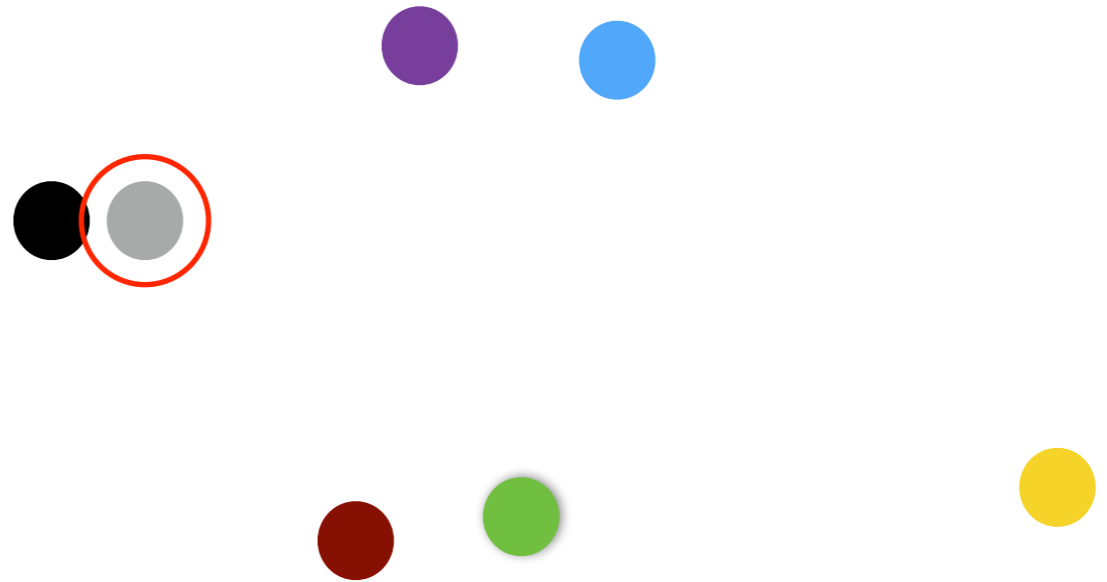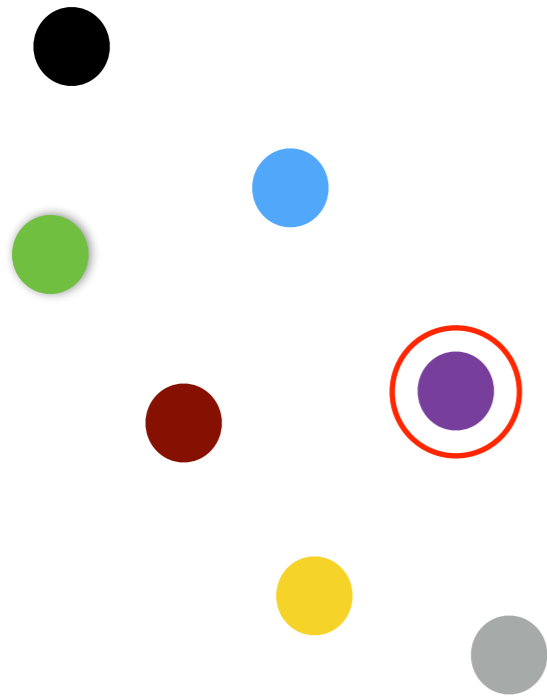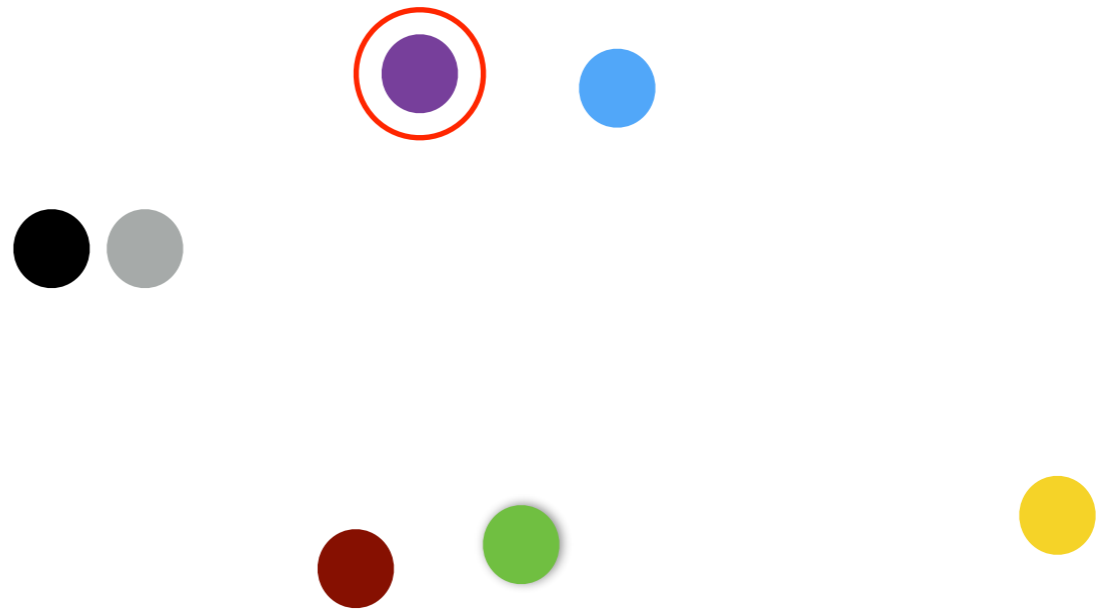# WHICH DIRECTION TO PICK?

View I

View II

# WHICH DIRECTION TO PICK?

View I

View II

# WHICH DIRECTION TO PICK?

View I

View II

View I

View II

# WHICH DIRECTION TO PICK?

PCA direction

Direction has large covariance
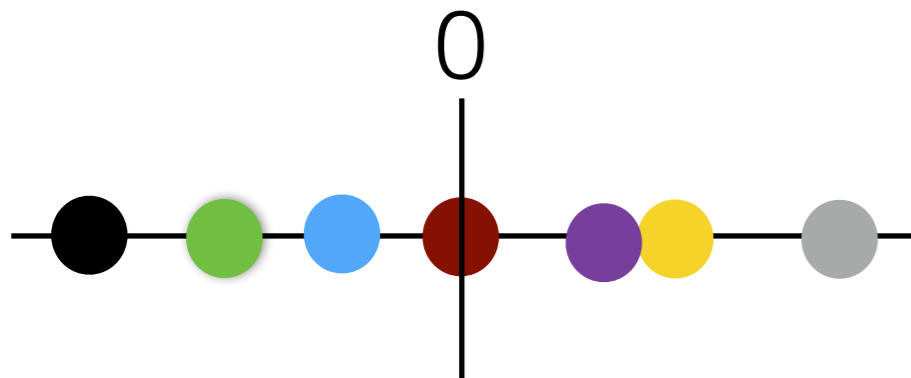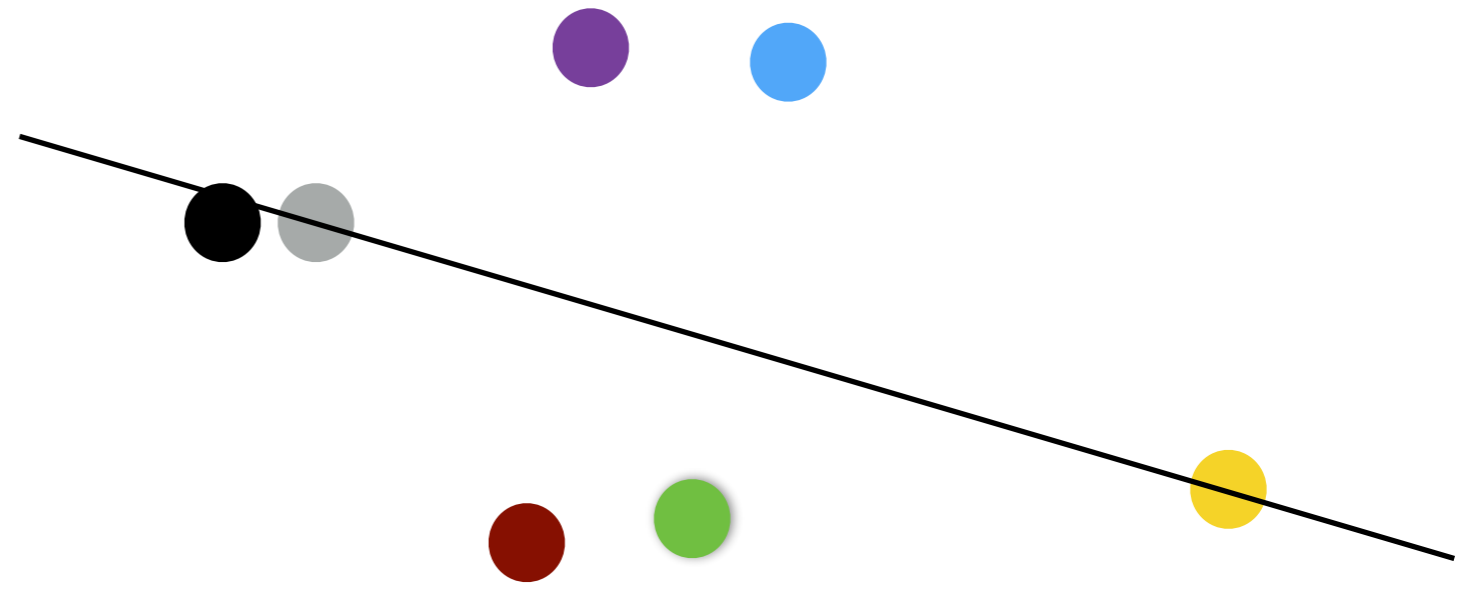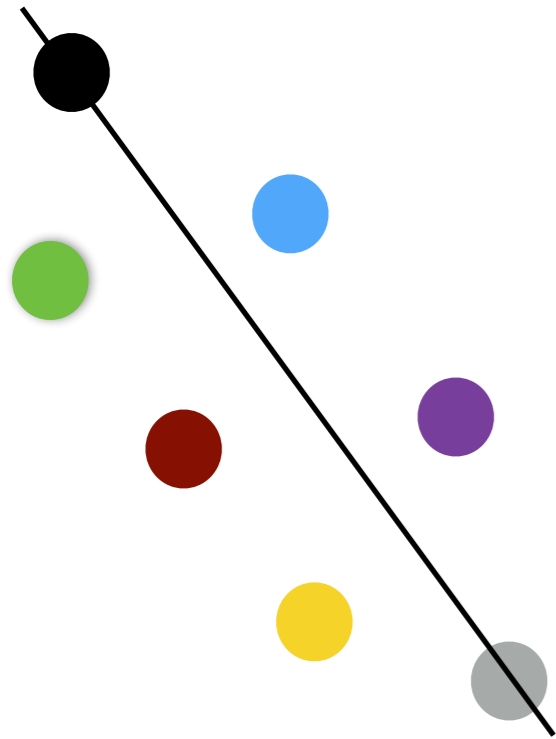
How do we pick the right direction to project to?

- 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?

$$\text{Say } \frac{1}{n} \sum_{t=1}^{n} \mathbf{x}_t[2] \cdot \mathbf{x}'_t[2] > 0$$

Scaling up this coordinate we can blow up covariance

$$\text{Say } \frac{1}{n}\sum_{t=1}^{n}\mathbf{x}_t[2] \cdot \mathbf{x'}_t[2] > 0$$

Scaling up this coordinate we can blow up covariance

Relevant information

$$\text{Say } \frac{1}{n} \sum_{t=1}^{n} \mathbf{x}_t[2] \cdot \mathbf{x'}_t[2] > 0$$

Scaling up this coordinate we can blow up covariance

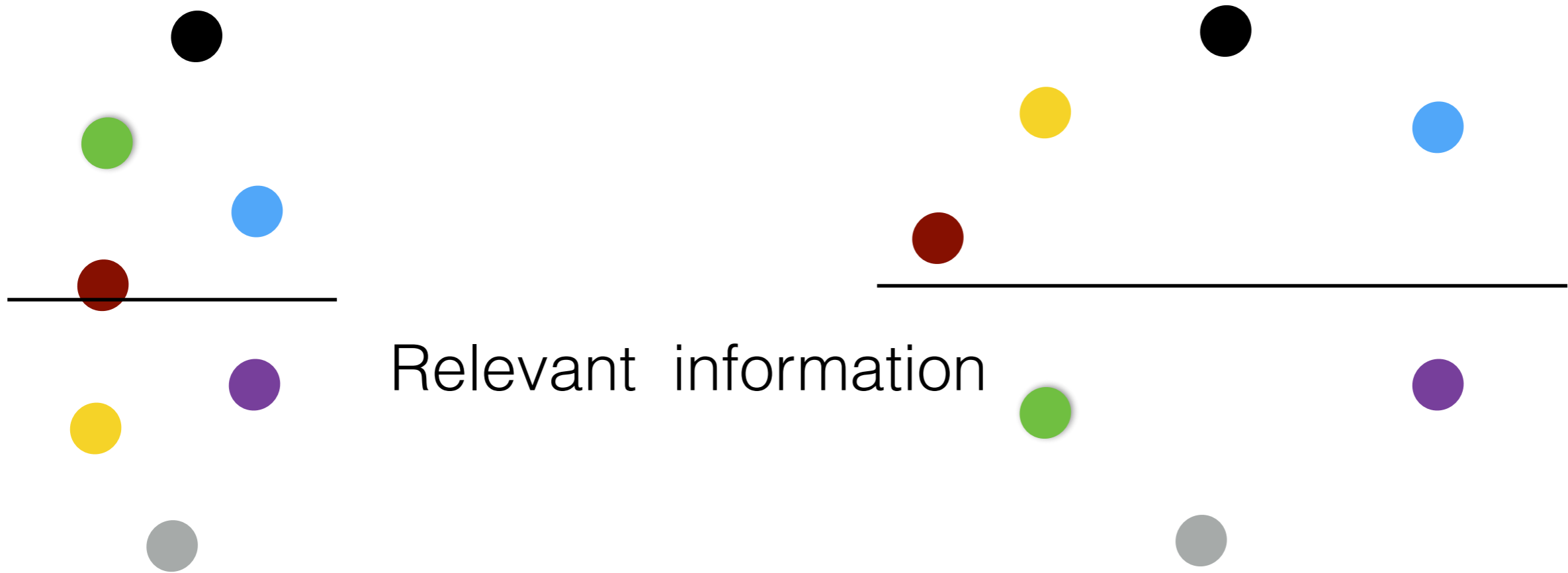- Say $\mathbf{w}_1$ and $\mathbf{v}_1$ are the directions we choose to project in views 1 and 2 respectively we want these directions to maximize,

$$\frac{\frac{1}{n} \sum_{t=1}^{n} \left( \mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^{n} \mathbf{y}_t[1] \right) \cdot \left( \mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^{n} \mathbf{y}'_t[1] \right)}{\sqrt{\frac{1}{n} \sum_{t=1}^{n} \left( \mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^{n} \mathbf{y}_t[1] \right)^2} \sqrt{\frac{1}{n} \sum_{t=1}^{n} \left( \mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^{n} \mathbf{y}'_t[1] \right)}}$$

- Normalize variance in chosen direction to be constant (say 1)

- Then maximize covariance

- This is same as maximizing "correlation coefficient"

- $\text{Covariance}(A, B) = \mathbb{E}[(A - \mathbb{E}[A]) \cdot (B - \mathbb{E}[B])]$

  Depends on the scale of $A$ and $B$. If $B$ is rescaled, covariance shifts.

- $\text{Corelation}(A, B) = \dfrac{\mathbb{E}[(A - \mathbb{E}[A]) \cdot (B - \mathbb{E}[B])]}{\sqrt{\text{Var}(A)} \sqrt{\text{Var}(B)}}$

  Scale free.

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

- 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,

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

s.t. $\frac{1}{n}\sum_{t=1}^{n}\left(\mathbf{y}_t[1] - \frac{1}{n}\sum_{t=1}^{n}\mathbf{y}_t[1]\right)^2 = \frac{1}{n}\sum_{t=1}^{n}\left(\mathbf{y}_t'[1] - \frac{1}{n}\sum_{t=1}^{n}\mathbf{y}_t'[1]\right) = 1$

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

- Hence we want to solve for projection vectors $\mathbf{w}_1$ and $\mathbf{v}_1$ that

$$\text{maximize } \frac{1}{n} \sum_{t=1}^{n} \mathbf{w}_1^\top (\mathbf{x}_t - \mu) \cdot \mathbf{v}_1^\top (\mathbf{x}_t' - \mu')$$

$$\text{subject to } \frac{1}{n} \sum_{t=1}^{n} (\mathbf{w}_1^\top (\mathbf{x}_t - \mu))^2 = \frac{1}{n} \sum_{t=1}^{n} (\mathbf{v}_1^\top (\mathbf{x}_t' - \mu'))^2 = 1$$

where $\mu = \frac{1}{n} \sum_{t=1}^{n} \mathbf{x}_t$ and $\mu' = \frac{1}{n} \sum_{t=1}^{n} \mathbf{x}_t'$

- Hence we want to solve for projection vectors $\mathbf{w}_1$ and $\mathbf{v}_1$ that

$$\text{maximize } \mathbf{w}_1^\top \Sigma_{1,2} \mathbf{v}_1$$
$$\text{subject to } \mathbf{w}_1^\top \Sigma_{1,1} \mathbf{w}_1 = \mathbf{v}_1^\top \Sigma_{2,2} \mathbf{v}_1 = 1$$

- Hence we want to solve for projection vectors $\mathbf{w}_1$ and $\mathbf{v}_1$ that

$$\text{maximize } \mathbf{w}_1^\top \Sigma_{1,2} \mathbf{v}_1$$

$$\text{subject to } \mathbf{w}_1^\top \Sigma_{1,1} \mathbf{w}_1 = \mathbf{v}_1^\top \Sigma_{2,2} \mathbf{v}_1 = 1$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \text{cov}\left( X \, X' \right)$$

$$W_1 = \text{eigs}\left(\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}, K\right)$$

$$W_2 = \text{eigs}\left(\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}, K\right)$$

1. $X = \left( \begin{array}{c} \text{n} \end{array} \boxed{X_1} \boxed{X_2} \right)$

$$\underbrace{\phantom{X_1}}_{d_1} \quad \underbrace{\phantom{X_2}}_{d_2}$$

1. $X = \left( n \quad X_1 \quad X_2 \right)$

   $\quad\quad\quad\quad\; d_1 \quad\;\; d_2$

2. $\Sigma = \begin{matrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{matrix} = \mathrm{cov}\left( X \right)$

1. $X = \left( n \begin{array}{cc} X_1 & X_2 \end{array} \right)$
   $\phantom{X = (} \underbrace{\phantom{X_1}}_{d_1} \quad \underbrace{\phantom{X_2}}_{d_2}$

2. $\Sigma = \begin{array}{cc} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{array} = \text{cov}\left( X \right)$

3. $W_1 = \text{eigs}\left( \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}, K \right)$

1. $X = \left( \begin{array}{cc} \text{n} \boxed{X_1} & \boxed{X_2} \end{array} \right)$
   $\quad\quad\quad\quad d_1 \quad\quad\quad d_2$

2. $\boxed{\sum} = \boxed{\begin{array}{cc} \sum_{11} & \sum_{12} \\ \sum_{21} & \sum_{22} \end{array}} = \text{cov}\left( \boxed{X} \right)$

3. $\boxed{W_1} = \text{eigs}\left( \boxed{\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}}, K \right)$

4. $\boxed{Y_1} = \boxed{X_1 - \mu_1} \times \boxed{W_1}$