# Machine Learning for Data Science (CS4786)
# Lecture 5

Random Projections

Course Webpage :
http://www.cs.cornell.edu/Courses/cs4786/2016sp/

# Recap

Given feature vectors $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbb{R}^d$, compress the data points into low dimensional representation $\mathbf{y}_1, \ldots, \mathbf{y}_n \in \mathbb{R}^K$ where $K << d$

**Principal Component Analysis:**
- **Find directions that maximize variance (spread)**
- **Find directions that minimize reconstruction error**

1. $$\Sigma = \mathrm{cov}(X)$$

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

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

4. $$\widehat{X} = Y \times W^\top + \mu$$

- Data can be split into pairs $(\mathbf{x}_1, \mathbf{x}_1'), \ldots, (\mathbf{x}_n, \mathbf{x}_n')$ where $\mathbf{x}_t$'s are $d_1$ dimensional and $\mathbf{x}_t''$'s are $d_2$ dimensional

- Goal: Compress $\mathbf{x}_1, \ldots, \mathbf{x}_n$ into $K$ dimensional vectors $\mathbf{y}_1, \ldots, \mathbf{y}_n$ (or $\mathbf{x}_1', \ldots, \mathbf{x}_n'$ into $\mathbf{y}_1', \ldots, \mathbf{y}_n'$ or both)

  - Retain information redundant between the two views

**Canonical Correlation Analysis:**
- **Find directions that maximize correlations between the projections in the two views**

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

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

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

4. $Y_1 = X_1 - \mu_1 \times W_1$

$$n \quad X \quad \times d\ W = n\ Y$$
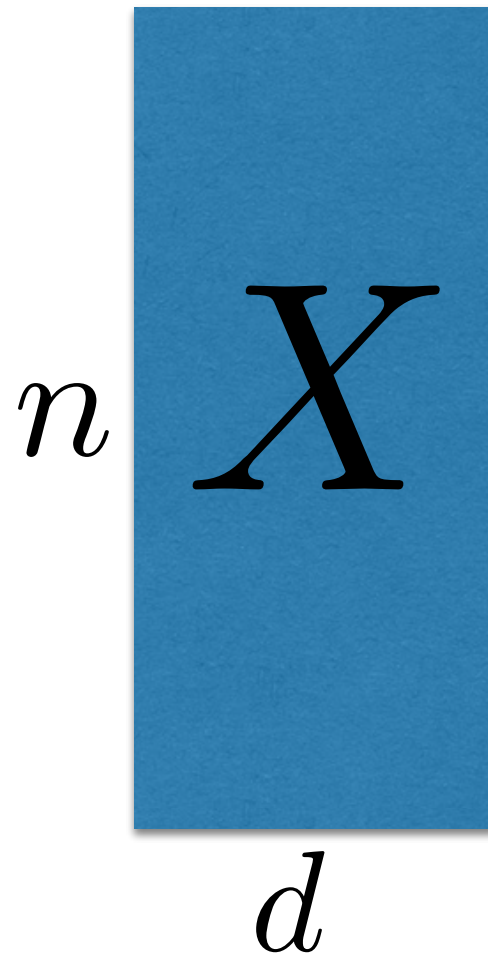
with row vectors $\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top$ of $X$ (width $d$), $W$ of height $d$ and width $K$, and $Y$ with rows $\mathbf{y}_1^\top, \dots, \mathbf{y}_n^\top$ (height $n$, width $K$).

$$n \quad X$$

$$d$$

$$\underset{n}{\underset{d}{\boxed{X^\top}}} \times \underset{d}{\underset{}{\boxed{n \, X}}} \Big/ n = \overset{d}{\underset{}{d \, \boxed{\Sigma}}}$$

$$\underset{K}{d\,\boxed{W}} = \mathrm{Eigs}\Big(\boxed{\Sigma}, K\Big)$$

$$n \quad \boxed{X} \quad d$$

$$n \quad X$$

$$d$$

$$\text{SVD}(X)$$

$$n \quad V \quad \times \quad \times \quad W^\top$$

$$n \qquad\qquad d$$

$$X$$

- $d$ 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?

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

$$K$$

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

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

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

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

$$
= \sum_{i,i'} \left( W[i,j] \cdot W[i',j] \right) \cdot \left( \tilde{\mathbf{x}}[i] \cdot \tilde{\mathbf{x}}[i'] \right)
$$

Hence,

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

if $i \neq i'$, $W[i,j]$ and $W[i',j]$ are independent and so

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

$$= \sum_{i=1}^{d} \tilde{\mathbf{x}}[i]^2 / \sqrt{K}^2 = \|\tilde{\mathbf{x}}\|_2^2 / K$$

Hence,

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

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

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

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

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

Lets try this in Matlab …

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 $i, j \in \{1, \ldots, n\}$,
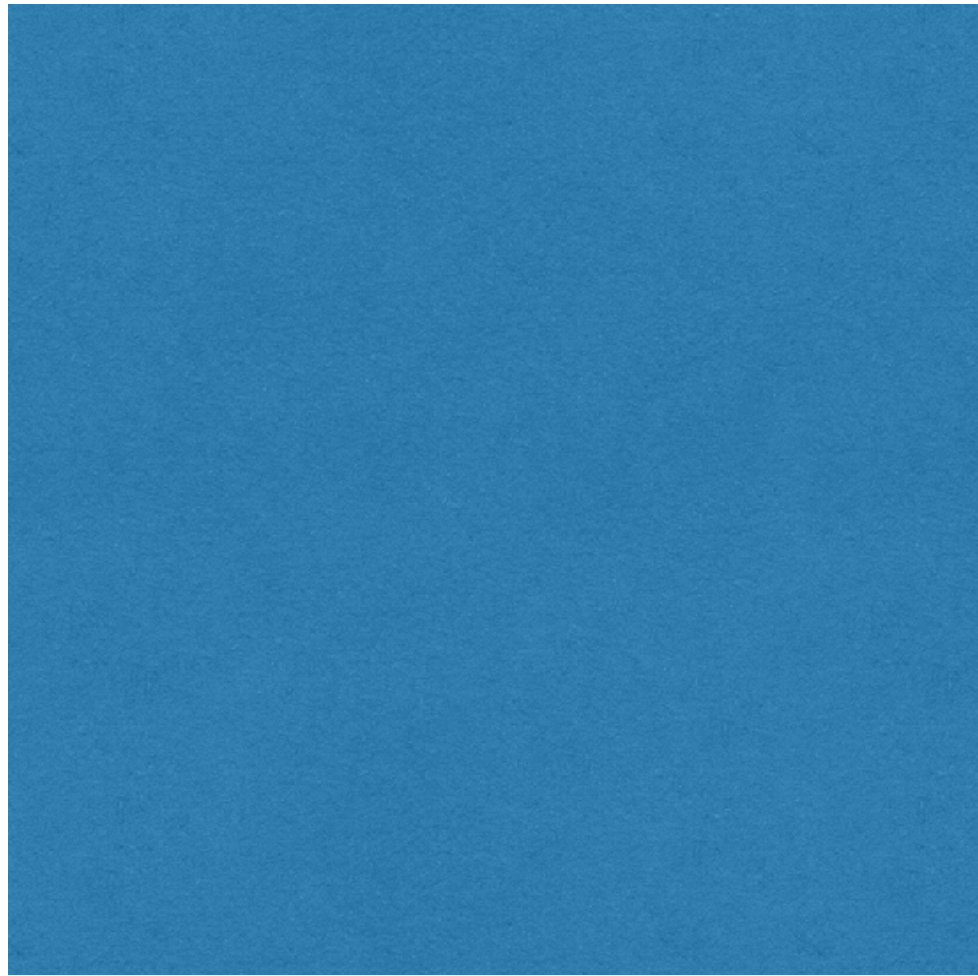
$$(1 - \epsilon)\left\|\mathbf{y}_i - \mathbf{y}_j\right\|_2 \leq \left\|\mathbf{x}_i - \mathbf{x}_j\right\|_2 \leq (1 + \epsilon)\left\|\mathbf{y}_i - \mathbf{y}_j\right\|_2$$

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
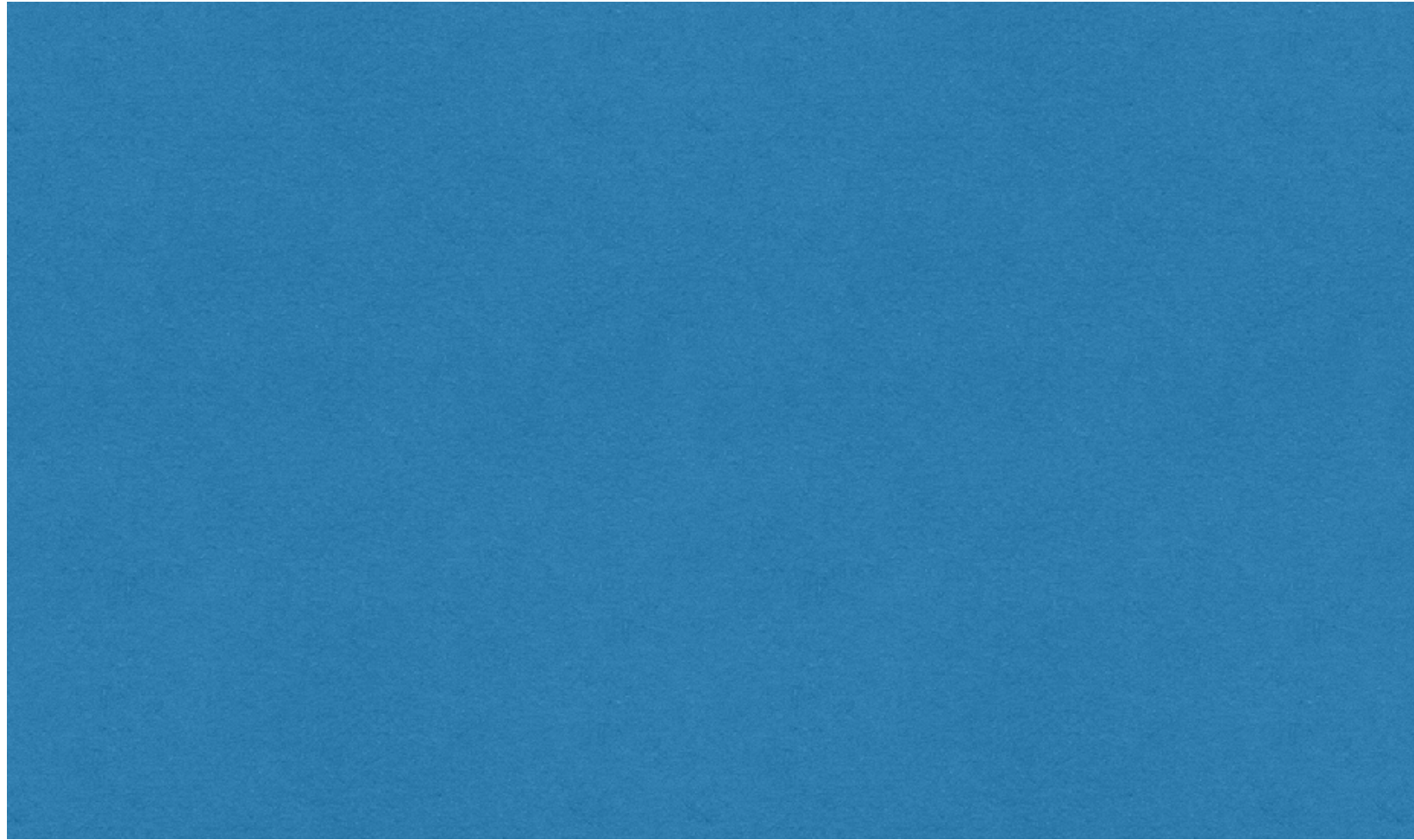


d = 10000

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 $\epsilon$