# Machine Learning for Data Science (CS4786) Lecture 6

Random Projections

Feb 09, 2015
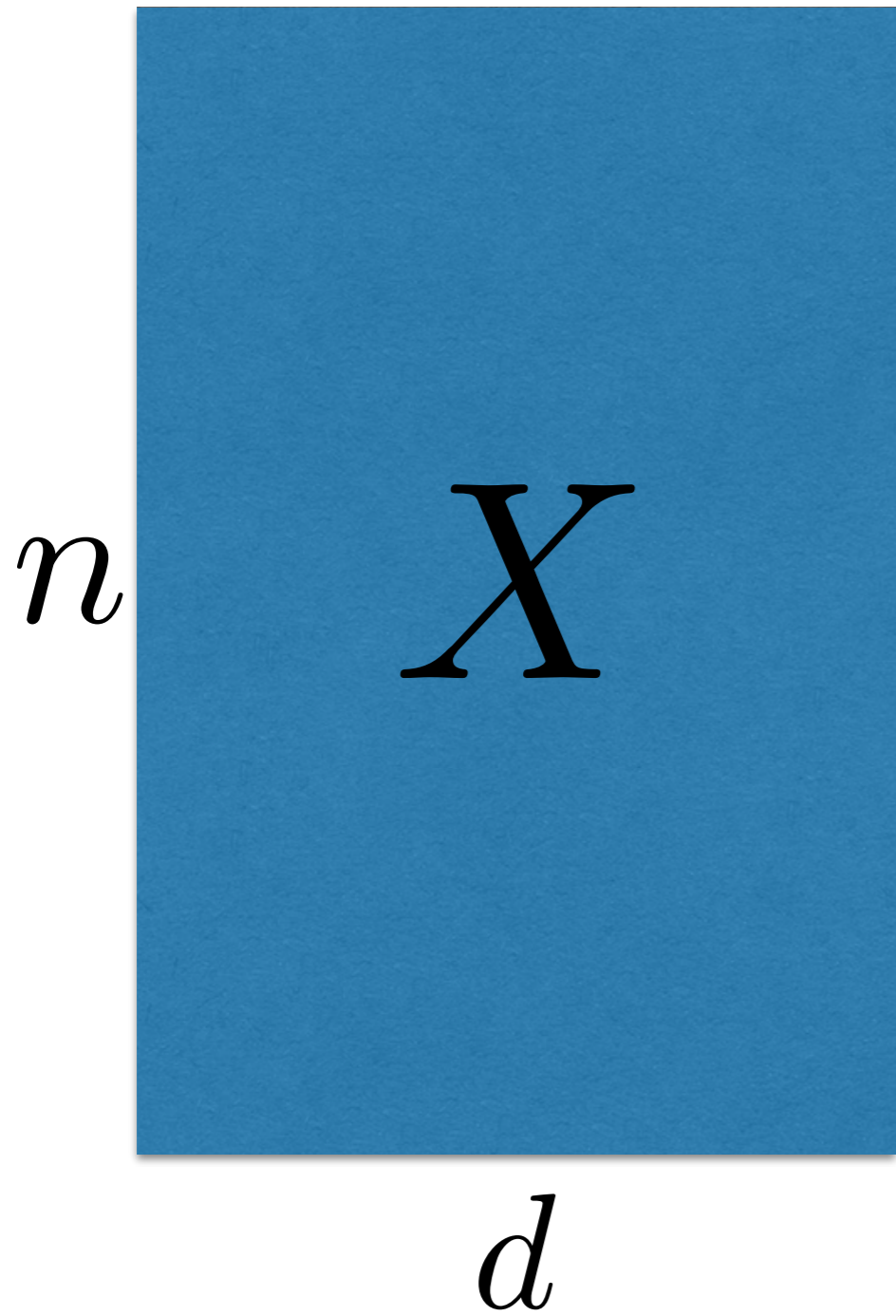
Course Webpage :
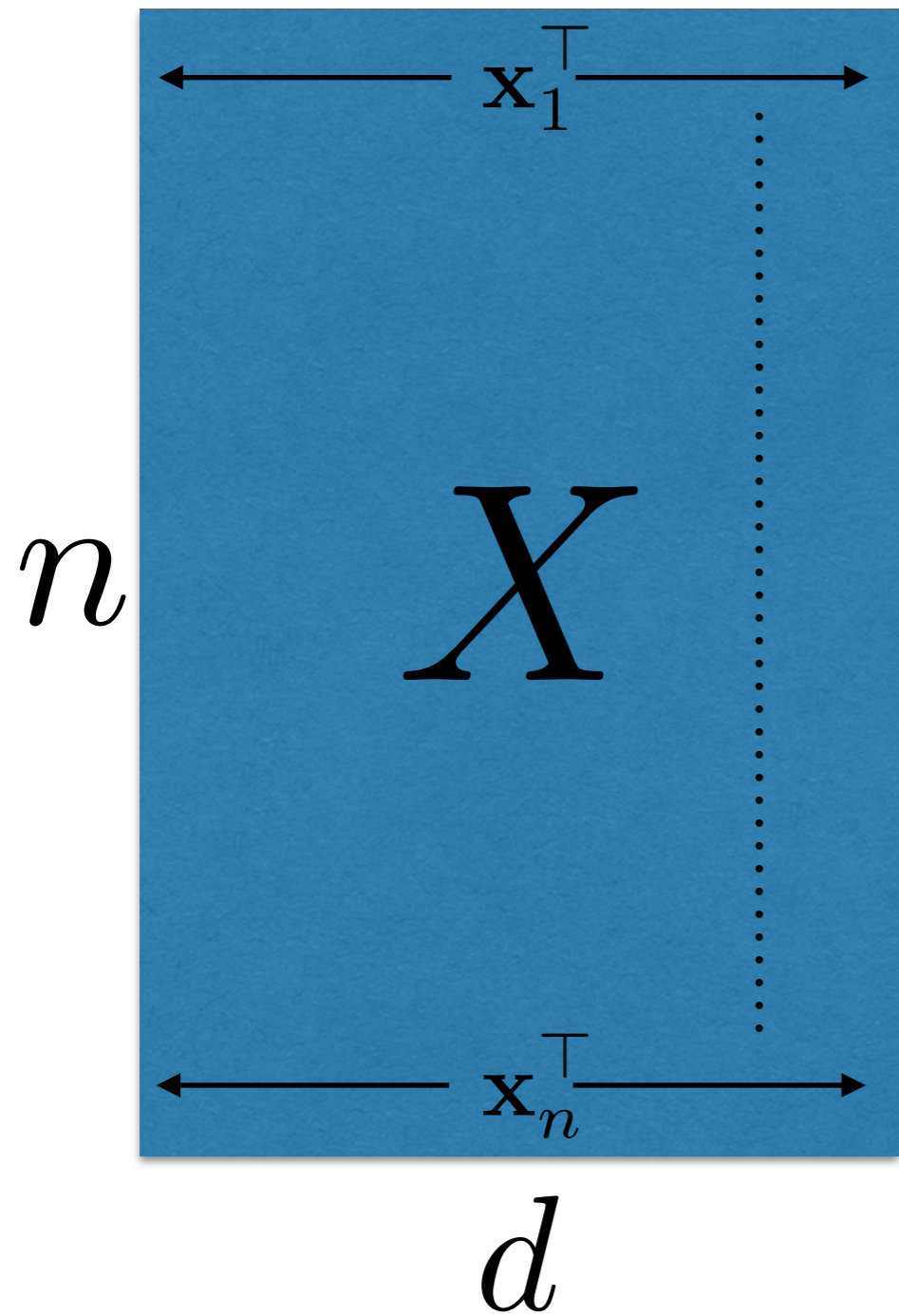http://www.cs.cornell.edu/Courses/cs4786/2015sp/

$n$ $X$

$d$

$$\mathbf{x}_1^\top$$

$$n \quad X$$

$$\mathbf{x}_n^\top$$

$$d$$

$n$ $X$ $\mathbf{x}_1^\top$ $\mathbf{x}_n^\top$ $d$

$n$ $Y$ $\mathbf{y}_1^\top$ $\mathbf{y}_n^\top$ $K$

$$n \, X \times d \, W = n \, Y$$

$\mathbf{x}_1^\top$

$\mathbf{x}_n^\top$

$d$

$K$

$\mathbf{y}_1^\top$

$\mathbf{y}_n^\top$

$K$

$$n \quad X$$

$$d$$

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

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

$$d \underset{K}{\boxed{W}} = \text{Eigs}\left(\overset{d}{\boxed{\Sigma}}, K\right)$$

$$n$$

$$X$$

$$d$$

$$X$$

- $d$ 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?

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

$$K$$

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

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

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

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

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

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

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,

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

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

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

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,

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

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,

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

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 $s, t \in \{1, \ldots, n\}$,

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

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 $s, t \in \{1, \ldots, n\}$,

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

Lets try this in Matlab …

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

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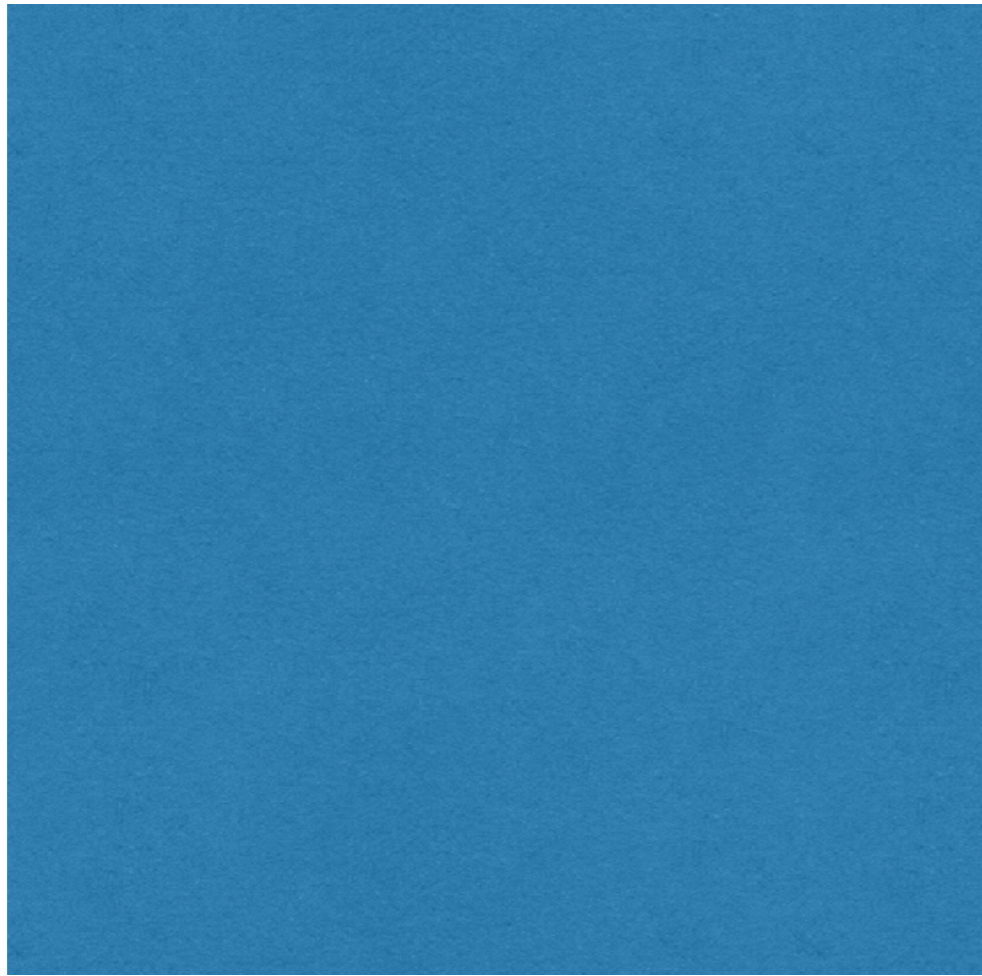
$$(1 - \epsilon) \left\| \mathbf{y}_i - \mathbf{y}_j \right\|_2 \leq \left\| \mathbf{x}_i - \mathbf{x}_j \right\|_2 \leq (1 + \epsilon) \left\| \mathbf{y}_i - \mathbf{y}_j \right\|_2$$

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 …

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 $\epsilon$