COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Andrew McGregor

Lecture 11

'Big Data' means not just many data points, but many measurements per data point. I.e., very high dimensional data.

'Big Data' means not just many data points, but many measurements per data point. I.e., very high dimensional data.

- Twitter has 321 million active monthly users. Records (tens of) thousands of measurements per user: who they follow, who follows them, when they last visited the site, timestamps for specific interactions, how many tweets they have sent, the text of those tweets, etc.

'Big Data' means not just many data points, but many measurements per data point. I.e., very high dimensional data.

- Twitter has 321 million active monthly users. Records (tens of) thousands of measurements per user: who they follow, who follows them, when they last visited the site, timestamps for specific interactions, how many tweets they have sent, the text of those tweets, etc.

- A 3 minute Youtube clip with a resolution of $500 \times 500$ pixels at 15 frames/second with 3 color channels is a recording of $\geq 2$ billion pixel values. Even a $500 \times 500$ pixel color image has $750,000$ pixel values.

'Big Data' means not just many data points, but many measurements per data point. I.e., very high dimensional data.

- Twitter has 321 million active monthly users. Records (tens of) thousands of measurements per user: who they follow, who follows them, when they last visited the site, timestamps for specific interactions, how many tweets they have sent, the text of those tweets, etc.

- A 3 minute Youtube clip with a resolution of $500 \times 500$ pixels at 15 frames/second with 3 color channels is a recording of $\geq 2$ billion pixel values. Even a $500 \times 500$ pixel color image has $750,000$ pixel values.

- The human genome contains 3 billion+ base pairs. Genetic datasets often contain information on 100s of thousands+ mutations and genetic markers.
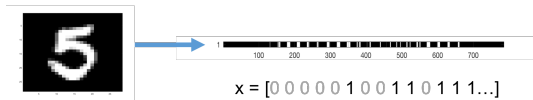
In data analysis and machine learning, data points with many attributes are often stored, processed, and interpreted as high dimensional vectors, with real valued entries.

In data analysis and machine learning, data points with many attributes are often stored, processed, and interpreted as high dimensional vectors, with real valued entries.

In data analysis and machine learning, data points with many attributes are often stored, processed, and interpreted as high dimensional vectors, with real valued entries.
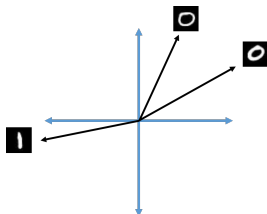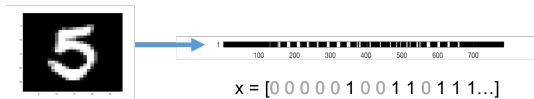
ATAGCCGTAGT $\longrightarrow$ x = [1 2 1 3 4 4 3 2 1 3 4]



x = [0 0 0 0 0 1 0 0 1 1 0 1 1 1 ...]



Similarities/distances between vectors (e.g., $\langle x, y \rangle$, $\|x - y\|_2$) have meaning for underlying data points.

Data points are interpreted as high dimensional vectors, with real valued entries. Data set is interpreted as a matrix.

**Data Points:** $\vec{x}_1, \vec{x}_2, \ldots, \vec{x}_n \in \mathbb{R}^d$.

**Data Set:** $X \in \mathbb{R}^{n \times d}$ with $i^{th}$ rows equal to $\vec{x}_i$.

Data points are interpreted as high dimensional vectors, with real valued entries. Data set is interpreted as a matrix.

**Data Points:** $\vec{x}_1, \vec{x}_2, \ldots, \vec{x}_n \in \mathbb{R}^d$.

**Data Set:** $X \in \mathbb{R}^{n \times d}$ with $i^{th}$ rows equal to $\vec{x}_i$.



$X \in \mathbb{R}^{n \times d}$

n = 3000 images

d = 784 pixels

Data points are interpreted as high dimensional vectors, with real valued entries. Data set is interpreted as a matrix.

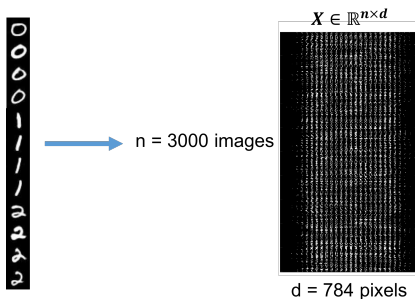**Data Points:** $\vec{x}_1, \vec{x}_2, \ldots, \vec{x}_n \in \mathbb{R}^d$.

**Data Set:** $X \in \mathbb{R}^{n \times d}$ with $i^{th}$ rows equal to $\vec{x}_i$.
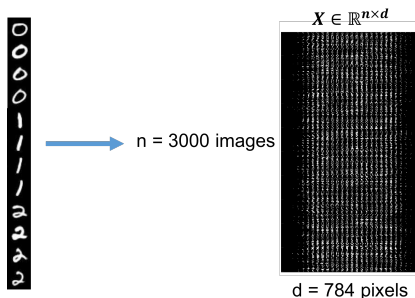


n = 3000 images

$X \in \mathbb{R}^{n \times d}$

d = 784 pixels

Many data points $n \implies$ tall. Many dimensions $d \implies$ wide.

**Dimensionality Reduction:** Compress data points so that they lie in many fewer dimensions.

**Dimensionality Reduction:** Compress data points so that they lie in many fewer dimensions.

$$\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d \to \tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m \text{ for } m \ll d.$$

 $\longrightarrow$ $x = [0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 1\ 1...]$ $\longrightarrow$ $\tilde{x} = [\text{-}5.5\ 4\ 3.2\ \text{-}1]$

**Dimensionality Reduction:** Compress data points so that they lie in many fewer dimensions.

$$\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d \rightarrow \tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m \text{ for } m \ll d.$$

$x = [0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 1\ 1...] \longrightarrow \tilde{x} = [\text{-5.5 4 3.2 -1}]$

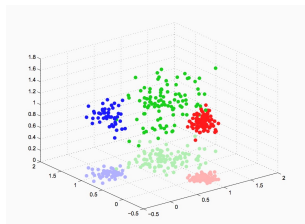'Lossy compression' that still preserves important information about the relationships between $\vec{x}_1, \ldots, \vec{x}_n$.

**Dimensionality Reduction:** Compress data points so that they lie in many fewer dimensions.

$$\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d \rightarrow \tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m \text{ for } m \ll d.$$

 ⟶ $x = [0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 1\ 1...]$ ⟶ $\tilde{x} = [\text{-5.5 4 3.2 -1}]$

'Lossy compression' that still preserves important information about the relationships between $\vec{x}_1, \ldots, \vec{x}_n$.



Generally will not consider directly how well $\tilde{x}_i$ approximates $\vec{x}_i$.

**Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$, distance function $D$, and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function $\tilde{D}$ such that for all $i, j \in [n]$:

$$(1 - \epsilon)D(\vec{x}_i, \vec{x}_j) \leq \tilde{D}(\tilde{x}_i, \tilde{x}_j) \leq (1 + \epsilon)D(\vec{x}_i, \vec{x}_j).$$

**Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$, distance function $D$, and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function $\tilde{D}$ such that for all $i, j \in [n]$:

$$(1 - \epsilon)D(\vec{x}_i, \vec{x}_j) \leq \tilde{D}(\tilde{x}_i, \tilde{x}_j) \leq (1 + \epsilon)D(\vec{x}_i, \vec{x}_j).$$

Have already seen one example in class: MinHash.



$x_A = [1\,0\,1\,0\,1\,1\,0\,0]$    A   B    $x_B = [0\,1\,0\,0\,1\,1\,0\,1]$

1   3   6   5   2   8

**Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$, distance function $D$, and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function $\tilde{D}$ such that for all $i, j \in [n]$:

$$(1 - \epsilon)D(\vec{x}_i, \vec{x}_j) \leq \tilde{D}(\tilde{x}_i, \tilde{x}_j) \leq (1 + \epsilon)D(\vec{x}_i, \vec{x}_j).$$

Have already seen one example in class: MinHash.

**Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$, distance function $D$, and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function $\tilde{D}$ such that for all $i, j \in [n]$:
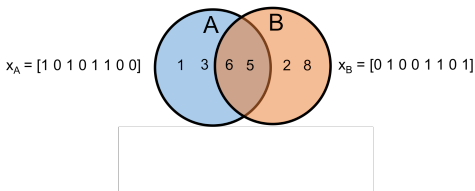
$$(1 - \epsilon)D(\vec{x}_i, \vec{x}_j) \leq \tilde{D}(\tilde{x}_i, \tilde{x}_j) \leq (1 + \epsilon)D(\vec{x}_i, \vec{x}_j).$$

Have already seen one example in class: MinHash.



$x_A = [1\ 0\ 1\ 0\ 1\ 1\ 0\ 0]$    A    B    $x_B = [0\ 1\ 0\ 0\ 1\ 1\ 0\ 1]$

1  3  6  5   2  8

MinHash

$\tilde{x}_A = [.12\ .09\ .17]$    $\tilde{x}_B = [.18\ .09\ .28]$

**Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$, distance function $D$, and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function $\tilde{D}$ such that for all $i, j \in [n]$:

$$(1 - \epsilon)D(\vec{x}_i, \vec{x}_j) \leq \tilde{D}(\tilde{x}_i, \tilde{x}_j) \leq (1 + \epsilon)D(\vec{x}_i, \vec{x}_j).$$
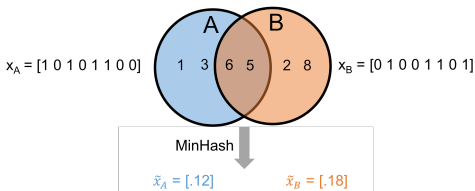
Have already seen one example in class: MinHash.



$x_A$ = [1 0 1 0 1 1 0 0]       A   B       $x_B$ = [0 1 0 0 1 1 0 1]

1   3   6   5   2   8

MinHash

$\tilde{x}_A$ = [.12 **.09** .17]     $\tilde{x}_B$ = [.18 **.09** .28]

With large enough signature size $r$, $\frac{\#\text{ matching entries in } \tilde{x}_A, \tilde{x}_B}{r} \approx J(\vec{x}_A, \vec{x}_B)$.

**Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$, distance function $D$, and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function $\tilde{D}$ such that for all $i, j \in [n]$:

$$(1 - \epsilon)D(\vec{x}_i, \vec{x}_j) \leq \tilde{D}(\tilde{x}_i, \tilde{x}_j) \leq (1 + \epsilon)D(\vec{x}_i, \vec{x}_j).$$
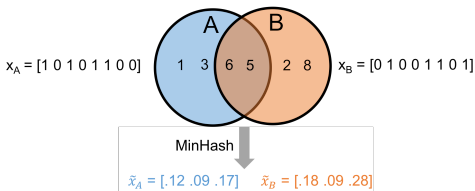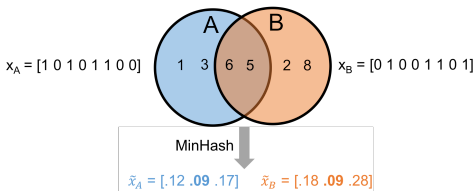
Have already seen one example in class: MinHash.



With large enough signature size $r$, $\frac{\# \text{ matching entries in } \tilde{x}_A, \tilde{x}_B}{r} \approx J(\vec{x}_A, \vec{x}_B)$.

- Note: here $J(\vec{x}_A, \vec{x}_B)$ is a similarity rather than a distance. So this is not quite low distortion embedding, but is closely related.

**Euclidean Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) such that for all $i, j \in [n]$:

$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$

Recall that for $\vec{z} \in \mathbb{R}^n$, $\|\vec{z}\|_2 = \sqrt{\sum_{i=1}^{n} \vec{z}(i)^2}$.

**Euclidean Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) such that for all $i, j \in [n]$:

$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$

Recall that for $\vec{z} \in \mathbb{R}^n$, $\|\vec{z}\|_2 = \sqrt{\sum_{i=1}^n \vec{z}(i)^2}$.



Pythagorean theorem.

$z(1)$

$z(2)$

$z$

$\|z\|_2 = \sqrt{z(1)^2 + z(2)^2}$

**Euclidean Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) such that for all $i, j \in [n]$:

$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$



**d-dimensional space**

$\|x_i - x_j\|_2$

**m-dimensional space**
(for m << d)

$\|\tilde{x}_i - \tilde{x}_j\|_2$

**Euclidean Low Distortion Embedding:** Given $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \ldots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) such that for all $i, j \in [n]$:

$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$
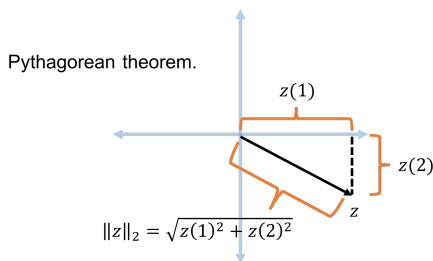


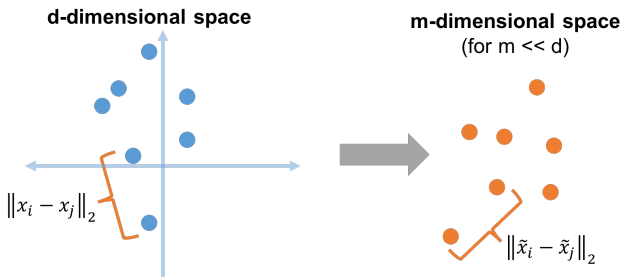Can use $\tilde{x}_1, \ldots, \tilde{x}_n$ in place of $\vec{x}_1, \ldots, \vec{x}_n$ in clustering, SVM, linear classification, near neighbor search, etc.

**A very easy case:** Assume that $\vec{x}_1, \ldots, \vec{x}_n$ all lie on the $1^{st}$ axis in $\mathbb{R}^d$.

## EMBEDDING WITH ASSUMPTIONS

**A very easy case:** Assume that $\vec{x}_1, \ldots, \vec{x}_n$ all lie on the $1^{st}$ axis in $\mathbb{R}^d$.



Set $m = 1$ and $\tilde{x}_i = [\vec{x}_i(1)]$ (i.e., $\tilde{x}_i$ contains just a single number).

• $\|\tilde{x}_i - \tilde{x}_j\|_2 = \sqrt{[\vec{x}_i(1) - \vec{x}_j(1)]^2} = |\vec{x}_i(1) - \vec{x}_j(1)| = \|\vec{x}_i - \vec{x}_j\|_2$.

**A very easy case:** Assume that $\vec{x}_1, \ldots, \vec{x}_n$ all lie on the $1^{st}$ axis in $\mathbb{R}^d$.



Set $m = 1$ and $\tilde{x}_i = [\vec{x}_i(1)]$ (i.e., $\tilde{x}_i$ contains just a single number).

- $\|\tilde{x}_i - \tilde{x}_j\|_2 = \sqrt{[\vec{x}_i(1) - \vec{x}_j(1)]^2} = |\vec{x}_i(1) - \vec{x}_j(1)| = \|\vec{x}_i - \vec{x}_j\|_2$.
- An embedding with no distortion from any $d$ into $m = 1$.

# EMBEDDING WITH ASSUMPTIONS

**A very easy case:** Assume that $\vec{x}_1, \ldots, \vec{x}_n$ all lie on the $1^{st}$ axis in $\mathbb{R}^d$.



Set $m = 1$ and $\tilde{x}_i = [\vec{x}_i(1)]$ (i.e., $\tilde{x}_i$ contains just a single number).

- $\|\tilde{x}_i - \tilde{x}_j\|_2 = \sqrt{[\vec{x}_i(1) - \vec{x}_j(1)]^2} = |\vec{x}_i(1) - \vec{x}_j(1)| = \|\vec{x}_i - \vec{x}_j\|_2$.

- An embedding with no distortion from any $d$ into $m = 1$.

- More generally. there's a no distortion embedding into $m = D$ dimensions if all the points lie is a $D$ dimensional space.

What about when we don't make any assumptions on $\vec{x}_1, \ldots, \vec{x}_n$.
I.e., they can be scattered arbitrarily around $d$-dimensional space?

- Can we find a no-distortion embedding into $m \ll d$ dimensions?

What about when we don't make any assumptions on $\vec{x}_1, \ldots, \vec{x}_n$.
I.e., they can be scattered arbitrarily around $d$-dimensional space?

- Can we find a no-distortion embedding into $m \ll d$ dimensions?
  No. Require $m = d$.

What about when we don't make any assumptions on $\vec{x}_1, \ldots, \vec{x}_n$.
I.e., they can be scattered arbitrarily around $d$-dimensional space?

- Can we find a no-distortion embedding into $m \ll d$ dimensions?
  No. Require $m = d$.

- Can we find an $\epsilon$-distortion embedding into $m \ll d$ dimensions
  for $\epsilon > 0$?

  For all $i, j : (1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.

What about when we don't make any assumptions on $\vec{x}_1, \ldots, \vec{x}_n$.
I.e., they can be scattered arbitrarily around $d$-dimensional space?

- Can we find a no-distortion embedding into $m \ll d$ dimensions?
  No. Require $m = d$.

- Can we find an $\epsilon$-distortion embedding into $m \ll d$ dimensions
  for $\epsilon > 0$? Yes! Always, with $m$ depending on $\epsilon$.

  For all $i, j : (1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$

**Johnson-Lindenstrauss Lemma:** For any set of points $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and $\epsilon > 0$ there exists a linear map $\mathbf{\Pi} : \mathbb{R}^d \to \mathbb{R}^m$ such that $m = O\left(\frac{\log n}{\epsilon^2}\right)$ and letting $\tilde{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all $i, j : (1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.

Further, if $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ has each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, it satisfies the guarantee with high probability.

**Johnson-Lindenstrauss Lemma:** For any set of points $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and $\epsilon > 0$ there exists a linear map $\mathbf{\Pi} : \mathbb{R}^d \to \mathbb{R}^m$ such that $m = O\left(\frac{\log n}{\epsilon^2}\right)$ and letting $\tilde{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all $i, j : \ (1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$

Further, if $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ has each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, it satisfies the guarantee with high probability.

For $d = 1$ trillion, $\epsilon = .05$, and $n = 100,000$, $m \approx 6600$.

**Johnson-Lindenstrauss Lemma:** For any set of points $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ and $\epsilon > 0$ there exists a linear map $\mathbf{\Pi} : \mathbb{R}^d \to \mathbb{R}^m$ such that $m = O\left(\frac{\log n}{\epsilon^2}\right)$ and letting $\tilde{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all $i, j$ : $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.

Further, if $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ has each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, it satisfies the guarantee with high probability.

For $d = 1$ trillion, $\epsilon = .05$, and $n = 100,000$, $m \approx 6600$.

Very surprising! Powerful result with a simple construction: applying a random linear transformation to a set of points preserves distances between all those points with high probability.

For any $\vec{x}_1, \ldots, \vec{x}_n$ and $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ with each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, with high probability, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all $i, j$: $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.



$m \times d$

| .01 | −1.2 | .34 | .67 | .10 | −.49 ... |
| −.45 | .7 | .14 | .18 | −.65 | .76 ... |

$\mathbf{\Pi}$

d x 1  $x_i$

m x 1  $\tilde{\mathbf{x}}_i$

$=$

**random linear transformation
(random projection)**

**compressed output point
(low dimensions)**

$$m = O\left(\frac{\log n}{\epsilon^2}\right)$$

**input point
(high dimensions)**

For any $\vec{x}_1, \ldots, \vec{x}_n$ and $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ with each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, with high probability, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all $i, j$ : $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.



$m \times d$     d x 1     m x 1

| .01 | − 1.2 | .34 | .67 | .10 | − .49 | ... |
| −.45 | .7 | .14 | .18 | − .65 | .76 | ... |

$\mathbf{\Pi}$

$x_i$ = $\tilde{x}_i$

**random linear transformation**
**(random projection)**

**compressed output point**
**(low dimensions)**

$$m = O\left(\frac{\log n}{\epsilon^2}\right)$$

**input point**
**(high dimensions)**

- $\mathbf{\Pi}$ is known as a random projection. It is a random linear function, mapping length $d$ vectors to length $m$ vectors.

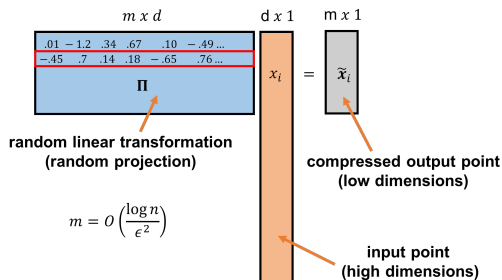For any $\vec{x}_1, \ldots, \vec{x}_n$ and $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ with each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, with high probability, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all $i, j$ : $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.



$m \times d$     d x 1    m x 1

| .01 | − 1.2 | .34 | .67 | .10 | − .49 ... |
| −.45 | .7 | .14 | .18 | − .65 | .76 ... |

$\mathbf{\Pi}$

$x_i \quad = \quad \widetilde{x}_i$

**random linear transformation**
**(random projection)**

**compressed output point**
**(low dimensions)**

$$m = O\left(\frac{\log n}{\epsilon^2}\right)$$

**input point**
**(high dimensions)**

- $\mathbf{\Pi}$ is known as a random projection. It is a random linear function, mapping length $d$ vectors to length $m$ vectors.

- $\mathbf{\Pi}$ is data oblivious. Stark contrast to methods like PCA.

- Many alternative constructions: $\pm 1$ entries, sparse (most entries 0), Fourier structured, etc. $\implies$ more efficient computation of $\mathbf{x}_i = \mathbf{\Pi} \vec{x}_i$.

- Many alternative constructions: $\pm 1$ entries, sparse (most entries 0), Fourier structured, etc. $\implies$ more efficient computation of $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$.

- Data oblivious property means that once $\mathbf{\Pi}$ is chosen, $\mathbf{x}_1, \ldots, \mathbf{x}_n$ can be computed in a stream with little memory.

- Storage is just $O(nm)$ rather than $O(nd)$.

- Many alternative constructions: $\pm 1$ entries, sparse (most entries 0), Fourier structured, etc. $\implies$ more efficient computation of $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$.
- Data oblivious property means that once $\mathbf{\Pi}$ is chosen, $\mathbf{x}_1, \ldots, \mathbf{x}_n$ can be computed in a stream with little memory.
- Storage is just $O(nm)$ rather than $O(nd)$.
- Compression can be performed in parallel on different servers.

- Many alternative constructions: $\pm 1$ entries, sparse (most entries 0), Fourier structured, etc. $\implies$ more efficient computation of $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$.

- Data oblivious property means that once $\mathbf{\Pi}$ is chosen, $\mathbf{x}_1, \ldots, \mathbf{x}_n$ can be computed in a stream with little memory.

- Storage is just $O(nm)$ rather than $O(nd)$.

- Compression can be performed in parallel on different servers.

- When new data points are added, can be easily compressed, without updating existing points.

The Johnson-Lindenstrauss Lemma is a direct consequence of a closely related lemma:

**Distributional JL Lemma:** Let $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$. If we set $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then for any $\vec{y} \in \mathbb{R}^d$, with probability $\geq 1 - \delta$

$$(1 - \epsilon)\|\vec{y}\|_2 \leq \|\mathbf{\Pi}\vec{y}\|_2 \leq (1 + \epsilon)\|\vec{y}\|_2$$

$\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

The Johnson-Lindenstrauss Lemma is a direct consequence of a closely related lemma:

**Distributional JL Lemma:** Let $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$. If we set $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then for any $\vec{y} \in \mathbb{R}^d$, with probability $\geq 1 - \delta$

$$(1 - \epsilon)\|\vec{y}\|_2 \leq \|\mathbf{\Pi}\vec{y}\|_2 \leq (1 + \epsilon)\|\vec{y}\|_2$$

Applying a random matrix $\mathbf{\Pi}$ to any vector $\vec{y}$ preserves $\vec{y}$'s norm with high probability.

$\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

The Johnson-Lindenstrauss Lemma is a direct consequence of a closely related lemma:

> **Distributional JL Lemma:** Let $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$. If we set $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then for any $\vec{y} \in \mathbb{R}^d$, with probability $\geq 1 - \delta$
> $$(1 - \epsilon)\|\vec{y}\|_2 \leq \|\mathbf{\Pi}\vec{y}\|_2 \leq (1 + \epsilon)\|\vec{y}\|_2$$

Applying a random matrix $\mathbf{\Pi}$ to any vector $\vec{y}$ preserves $\vec{y}$'s norm with high probability.

- Like a low-distortion embedding, but for the length of a compressed vector rather than distances between vectors.

$\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Distributional JL Lemma $\implies$ JL Lemma:** Distributional JL show that a random projection $\mathbf{\Pi}$ preserves the norm of any $y$. The main JL Lemma says that $\mathbf{\Pi}$ preserves distances between vectors.

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Distributional JL Lemma $\implies$ JL Lemma:** Distributional JL show that a random projection $\mathbf{\Pi}$ preserves the norm of any $y$. The main JL Lemma says that $\mathbf{\Pi}$ preserves distances between vectors.

Since $\mathbf{\Pi}$ is linear these are the same thing!

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Distributional JL Lemma $\implies$ JL Lemma:** Distributional JL show that a random projection $\Pi$ preserves the norm of any $y$. The main JL Lemma says that $\Pi$ preserves distances between vectors.

Since $\Pi$ is linear these are the same thing!

**Proof:** Given $\vec{x}_1, \ldots, \vec{x}_n$, define $\binom{n}{2}$ vectors $\vec{y}_{ij}$ where $\vec{y}_{ij} = \vec{x}_i - \vec{x}_j$.
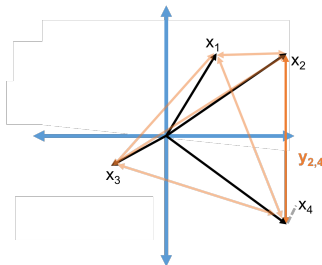
---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\Pi \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Distributional JL Lemma $\implies$ JL Lemma:** Distributional JL show that a random projection $\mathbf{\Pi}$ preserves the norm of any $y$. The main JL Lemma says that $\mathbf{\Pi}$ preserves distances between vectors.

Since $\mathbf{\Pi}$ is linear these are the same thing!

**Proof:** Given $\vec{x}_1, \ldots, \vec{x}_n$, define $\binom{n}{2}$ vectors $\vec{y}_{ij}$ where $\vec{y}_{ij} = \vec{x}_i - \vec{x}_j$.



$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

13

**Distributional JL Lemma** $\implies$ **JL Lemma:** Distributional JL show that a random projection $\mathbf{\Pi}$ preserves the norm of any $y$. The main JL Lemma says that $\mathbf{\Pi}$ preserves distances between vectors.

Since $\mathbf{\Pi}$ is linear these are the same thing!

**Proof:** Given $\vec{x}_1, \ldots, \vec{x}_n$, define $\binom{n}{2}$ vectors $\vec{y}_{ij}$ where $\vec{y}_{ij} = \vec{x}_i - \vec{x}_j$.

- If we choose $\mathbf{\Pi}$ with $m = O\left(\frac{\log 1/\delta'}{\epsilon^2}\right)$, for each $\vec{y}_{ij}$ with probability $\geq 1 - \delta'$ we have:

$$(1 - \epsilon)\|\vec{y}_{ij}\|_2 \leq \|\mathbf{\Pi}\vec{y}_{ij}\|_2 \leq (1 + \epsilon)\|\vec{y}_{ij}\|_2$$

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Distributional JL Lemma** $\implies$ **JL Lemma:** Distributional JL show that a random projection $\mathbf{\Pi}$ preserves the norm of any $y$. The main JL Lemma says that $\mathbf{\Pi}$ preserves distances between vectors.

Since $\mathbf{\Pi}$ is linear these are the same thing!

**Proof:** Given $\vec{x}_1, \ldots, \vec{x}_n$, define $\binom{n}{2}$ vectors $\vec{y}_{ij}$ where $\vec{y}_{ij} = \vec{x}_i - \vec{x}_j$.

- If we choose $\mathbf{\Pi}$ with $m = O\left(\frac{\log 1/\delta'}{\epsilon^2}\right)$, for each $\vec{y}_{ij}$ with probability $\geq 1 - \delta'$ we have:

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

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Distributional JL Lemma $\implies$ JL Lemma:** Distributional JL show that a random projection $\boldsymbol{\Pi}$ preserves the norm of any $y$. The main JL Lemma says that $\boldsymbol{\Pi}$ preserves distances between vectors.

Since $\boldsymbol{\Pi}$ is linear these are the same thing!

**Proof:** Given $\vec{x}_1, \ldots, \vec{x}_n$, define $\binom{n}{2}$ vectors $\vec{y}_{ij}$ where $\vec{y}_{ij} = \vec{x}_i - \vec{x}_j$.

- If we choose $\boldsymbol{\Pi}$ with $m = O\left(\frac{\log 1/\delta'}{\epsilon^2}\right)$, for each $\vec{y}_{ij}$ with probability $\geq 1 - \delta'$ we have:

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

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x}_1, \ldots, \mathbf{x}_n$: compressed points, $\boldsymbol{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x}_1, \ldots, \mathbf{x}_n$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x}_1, \ldots, \mathbf{x}_n$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$.

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x}_1, \ldots, \mathbf{x}_n$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$. $\implies$ for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$. $\implies$ for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

$$m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$$

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x}_1, \ldots, \mathbf{x}_n$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$. $\implies$ for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

$$m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right) = O\left(\frac{\log(\binom{n}{2}/\delta)}{\epsilon^2}\right)$$

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$. $\implies$ for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

$$m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right) = O\left(\frac{\log(\binom{n}{2}/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n^2/\delta)}{\epsilon^2}\right)$$

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x_1}, \ldots, \mathbf{x_n}$: compressed points, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\Pi$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \Pi \vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$. $\implies$ for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

$$m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right) = O\left(\frac{\log(\binom{n}{2}/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n^2/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n/\delta)}{\epsilon^2}\right)$$

---

$\vec{x}_1, \ldots, \vec{x}_n$: original points, $\mathbf{x}_1, \ldots, \mathbf{x}_n$: compressed points, $\Pi \in \mathbb{R}^{m \times d}$: random projection matrix. $d$: original dimension. $m$: compressed dimension, $\epsilon$: embedding error, $\delta$: embedding failure prob.

**Claim:** If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\mathbf{x}_i = \mathbf{\Pi}\vec{x}_i$, for each pair $\vec{x}_i, \vec{x}_j$ with probability $\geq 1 - \delta'$ we have:

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

With what probability are all pairwise distances preserved?

**Union bound:** With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta / \binom{n}{2}$. $\implies$ for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

$$m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right) = O\left(\frac{\log(\binom{n}{2}/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n^2/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n/\delta)}{\epsilon^2}\right)$$

Yields the JL lemma.