# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2022.

Lecture 11

$$4.1 \qquad \text{tables of size } O(m)$$
$$\text{tables of size } \quad \underline{m = O(k)}$$

- Problem Set 2 is due on Friday at 11:59pm.
- Midterm is in class next Thursday, 10/20.
- I have posted a study guide and practice questions on the course schedule.

$$\frac{n_k}{m} \qquad \frac{n - \sum_{i=1}^{k} f_i}{}$$

## Summary

### Last Class:

- Introduced the *k*-frequent elements problem – identify all elements of a stream of *n* elements that occur $\geq n/k$ times.
- Saw how to solve approximately in $O(k \log n/\epsilon)$ space using the Count-min sketch algorithm. $O(k/\epsilon)$
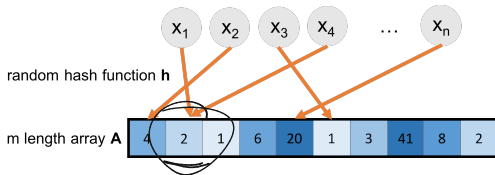- Simple analysis based on Markov's inequality and repeated random hashing.

### This Class:

- Recap and finish up Count-min sketch
- Randomized methods for dimensionality reduction.
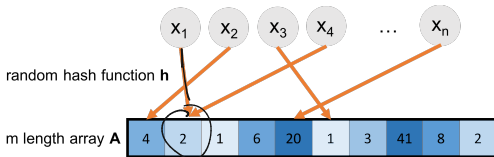- The Johnson-Lindenstrauss Lemma.

**Goal:** Return all items in a stream of $n$ elements with frequency at least $n/k$. Don't return any with frequency $\leq (1 - \epsilon) \cdot \frac{n}{k}$.

$(\epsilon, k)$ - Freq, Items problem

random hash function **h**

m length array **A** | 4 | 2 | 1 | 6 | 20 | 1 | 3 | 41 | 8 | 2 |

$F(x_1)$

$x_1$  $x_2$  $x_3$  $x_4$  $\ldots$  $x_n$

**Goal:** Return all items in a stream of $n$ elements with frequency at least $n/k$. Don't return any with frequency $\leq (1 - \epsilon) \cdot \frac{n}{k}$.



- $\underline{A[\mathbf{h}(x)]} = \underline{f(x)} + \underline{\sum_{y \neq x : \mathbf{h}(y) = \mathbf{h}(x)} f(y)}$ where
  $\underline{\mathbb{E}[\sum_{y \neq x : \mathbf{h}(y) = \mathbf{h}(x)} f(y)]} \leq \underline{\frac{n}{m}}$.

# Count-Min Sketch

**Goal:** Return all items in a stream of $n$ elements with frequency at least $n/k$. Don't return any with frequency $\leq (1 - \epsilon) \cdot \frac{n}{k}$.
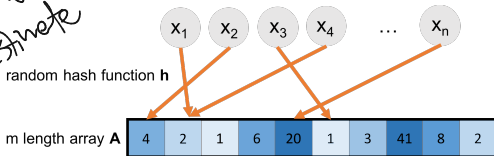
$f(x) =$ true freq. of $x$
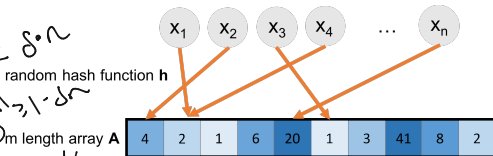
$\tilde{f}(x)$ is our estimate



random hash function **h**

m length array **A** | 4 | 2 | 1 | 6 | 20 | 1 | 3 | 41 | 8 | 2 |

- $A[\mathbf{h}(x)] = f(x) + \sum_{y \neq x : \mathbf{h}(y) = \mathbf{h}(x)} f(y)$ where $\mathbb{E}[\sum_{y \neq x : \mathbf{h}(y) = \mathbf{h}(x)} f(y)] \leq \frac{n}{m}$.

$\Pr\left(\tilde{f}(x) \geq f(x) + \frac{2n}{m}\right) \leq \frac{1}{2}$

- If we let $\tilde{f}(x)$ be the minimum of $t = \log_2(1/\delta)$ estimates, $f(x) \leq \tilde{f}(x) \leq f(x) + \frac{2n}{m}$ with probability at least $1 - \delta$.

$1 - \frac{1}{2^t} = 1 - \delta$

4

**Goal:** Return all items in a stream of $n$ elements with frequency at least $n/k$. Don't return any with frequency $\leq (1 - \epsilon) \cdot \frac{n}{k}$.

$\Pr(\tilde{f}(x) \text{ is 'bad'}) \leq \delta$

$\Pr(\tilde{f}(x) \text{ is bad for } \leq \delta \cdot n$

$\Pr(\tilde{f}(x) \text{ } x) \geq 1 - \delta \cdot n$

$\Pr(\tilde{f}(x) \text{ is good for all } x \text{ in stream}) \geq 1 - \delta'$

random hash function **h**

m length array **A**

| 4 | 2 | 1 | 6 | 20 | 1 | 3 | 41 | 8 | 2 |
|---|---|---|---|----|---|---|----|---|---|

$\delta = \delta'/n$

$1 - \delta'$

$\frac{2n}{m} \leq \frac{\epsilon n}{k}$

- $A[\mathbf{h}(x)] = f(x) + \sum_{y \neq x : \mathbf{h}(y) = \mathbf{h}(x)} f(y)$ where $\mathbb{E}[\sum_{y \neq x : \mathbf{h}(y) = \mathbf{h}(x)} f(y)] \leq \frac{n}{m}$.

- If we let $\tilde{f}(x)$ be the minimum of $t = \log_2(1/\delta)$ estimates, $f(x) \leq \tilde{f}(x) \leq f(x) + \frac{2n}{m}$ with probability at least $1 - \delta$.

- Setting $m = O(k/\epsilon)$, $\delta = O(\delta'/n)$, and applying a union bound, we have a good estimate for all $f(x)$ with probability at least $1 - \delta'$.

$m \cdot t = O\left(\frac{\log(1/\delta')}{\epsilon} \cdot k\right)$

stores

## Identifying Frequent Elements

Count-min sketch gives an accurate frequency estimate for every item in the stream. But how do we identify the frequent items without having to store/look up the estimated frequency for all elements in the stream?

## Identifying Frequent Elements

Count-min sketch gives an accurate frequency estimate for every item in the stream. But how do we identify the frequent items without having to store/look up the estimated frequency for all elements in the stream?

$\frac{1}{1-\epsilon} \cdot k$

$\left( + \left\{ \right. \right.$

[list of items]

### One approach:

· When a new item comes in at step $i$, check if its estimated frequency is $\geq i/k$ and store it if so.

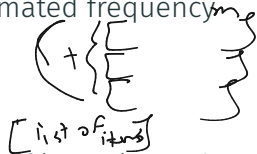At step $i$ remove any stored items whose estimated frequency drops below $i/k$.

· Store at most $O(k)$ items at once and have all items with frequency $\geq n/k$ stored at the end of the stream, no items with frequency $< (1 - \epsilon) \cdot \frac{n}{k}$.

$O\left( \log n \cdot \frac{k}{\epsilon} \right)$   $O(k)$   $O(m \cdot t + k)$

## High Dimensional Data

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

## 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.

## 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.

- 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.

# 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.

- 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.

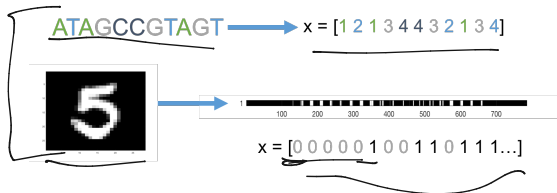- The human genome contains 3 billion+ base pairs. Genetic datasets often contain information on 100s of thousands+ mutations and genetic markers.

## Data as Vectors and Matrices

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.

## Data as Vectors and Matrices

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.

# Data as Vectors and Matrices

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.
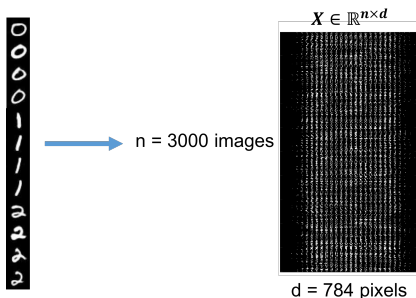
## Datasets as Vectors and Matrices

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 $i^{th}$ row equal to $\vec{x}_i$.

# Datasets as Vectors and Matrices

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 $i^{th}$ row equal to $\vec{x}_i$.



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

n = 3000 images

d = 784 pixels

# Datasets as Vectors and Matrices

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 $i^{th}$ row equal to $\vec{x}_i$.



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

n = 3000 images

d = 784 pixels

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

## Dimensionality Reduction

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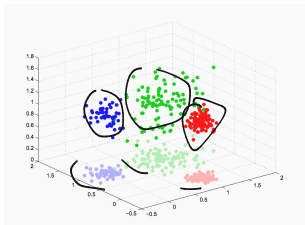
# Dimensionality Reduction

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

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

$x = [0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 1\ 1 \ldots] \longrightarrow \tilde{x} = [\text{-5.5}\ 4\ 3.2\ \text{-1}]$
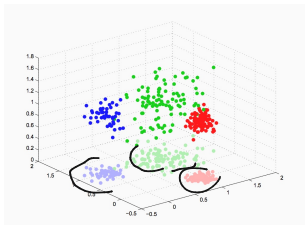
# Dimensionality Reduction

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

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

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

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

# Dimensionality Reduction

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

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

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

'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 $\vec{x}_i$.

# Dimensionality Reduction

Dimensionality reduction is one of the most important techniques in data science. What methods have you heard of?

PCA → LSA

T-SNE

Auto encoders, Nueral networks
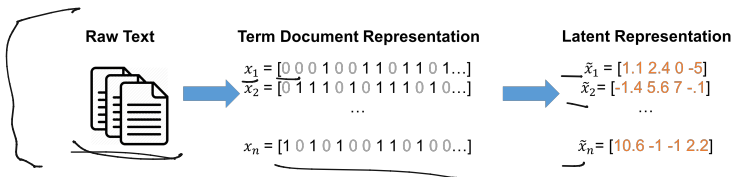
Self -organizing maps

clustering

JPEG , mp3

# Dimensionality Reduction

Dimensionality reduction is one of the most important techniques in data science. What methods have you heard of?

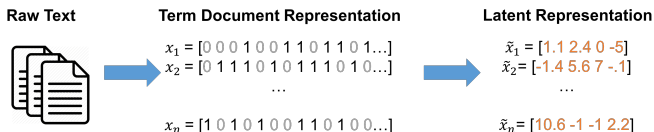- Principal component analysis
- Latent semantic analysis (LSA)



- Linear discriminant analysis
- Autoencoders

Dimensionality reduction is one of the most important techniques in data science. What methods have you heard of?

- Principal component analysis
- Latent semantic analysis (LSA)

| Raw Text | Term Document Representation | Latent Representation |
|---|---|---|
| | $x_1 = [0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 1\ 0\ 1...]$ | $\tilde{x}_1 = [1.1\ 2.4\ 0\ -5]$ |
| | $x_2 = [0\ 1\ 1\ 1\ 0\ 1\ 0\ 1\ 1\ 1\ 0\ 1\ 0...]$ | $\tilde{x}_2 = [-1.4\ 5.6\ 7\ -.1]$ |
| | ... | ... |
| | $x_n = [1\ 0\ 1\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 0\ 0...]$ | $\tilde{x}_n = [10.6\ -1\ -1\ 2.2]$ |

- Linear discriminant analysis
- Autoencoders

Compressing data makes it more efficient to work with. May also remove extraneous information/noise.
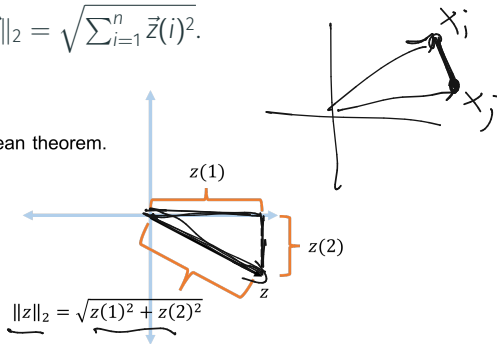
10

## Embeddings for Euclidean Space

**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 $i, j \in [n]$:

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

m-dimensel          d-dimension

# Embeddings for Euclidean Space

**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 $i, j \in [n]$:

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

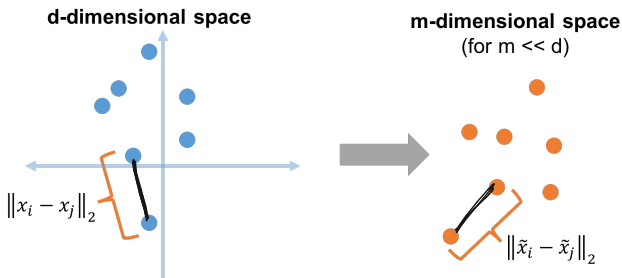Recall that for $\vec{z} \in \mathbb{R}^n$, $\|\vec{z}\|_2 = \sqrt{\sum_{i=1}^n \vec{z}(i)^2}$.

Pythagorean theorem.

$z(1)$

$z(2)$

$z$

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

$x_i$

$x_j$

# Embeddings for Euclidean Space

**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 $i, j \in [n]$:
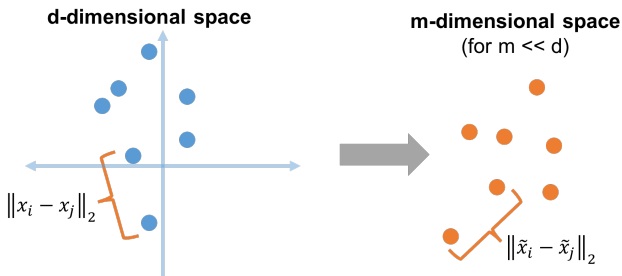
$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$

**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 $i, j \in [n]$:
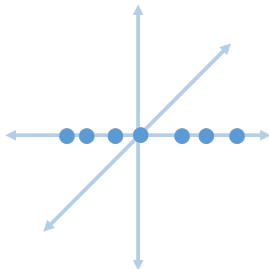
$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$



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.
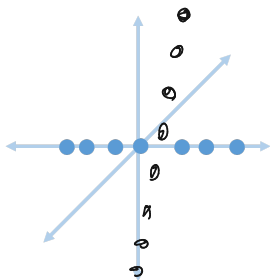
11

## Embedding with Assumptions

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

## Embedding with Assumptions

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



$$\widehat{x}_1 = \begin{bmatrix} -5 \end{bmatrix} \quad \widehat{x}_2 = \begin{bmatrix} 3 \end{bmatrix}$$

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

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

## Embedding with Assumptions

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



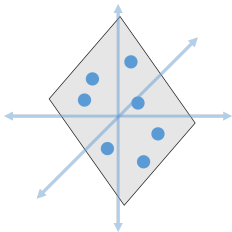Set $m = 1$ and $\tilde{x}_i = [\vec{x}_i(1)]$ (i.e., $\tilde{x}_i$ contains just a single number).

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

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

# Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



$$c_1 \vec{v}_1 + c_2 \vec{v}_2$$

$$V = \begin{bmatrix} | & | & | \\ v_1 & v_2 & v_k \\ | & | & | \end{bmatrix}$$

- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $V \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

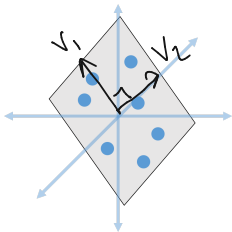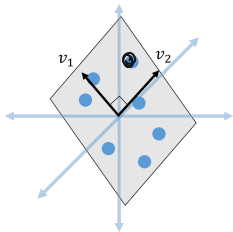## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\underbrace{\vec{x}_i - \vec{x}_j}\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle \underbrace{v_\ell}, \underbrace{\vec{x}_i - \vec{x}_j} \rangle^2} \qquad .$$

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \underbrace{\vec{x}_i - \vec{x}_j} \rangle^2} \quad .$$

# Embedding with Assumptions

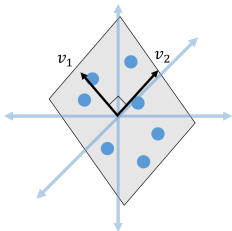Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.
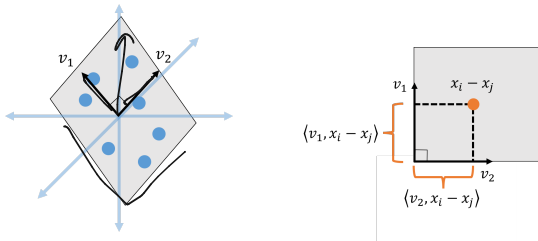


- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle \vec{v}_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

13

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.
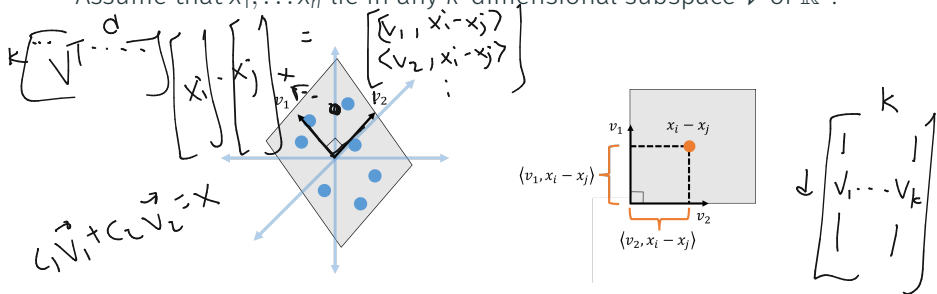
- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathsf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.

- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

- If we set $\tilde{x}_i \in \mathbb{R}^k$ to $\tilde{x}_i = \mathbf{V}^T \vec{x}_i$ we have:

$$\|\tilde{x}_i - \tilde{x}_j\|_2 = \|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2$$

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.

- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

- If we set $\tilde{x}_i \in \mathbb{R}^k$ to $\tilde{x}_i = \mathbf{V}^T\vec{x}_i$ we have:

$$\|\tilde{x}_i - \tilde{x}_j\|_2 = \|\mathbf{V}^T\vec{x}_i - \mathbf{V}^T\vec{x}_j\|_2 = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2$$

13

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.

- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

$$\widetilde{x}_1 = \begin{bmatrix} & V^T & \end{bmatrix}\begin{bmatrix} x_i \end{bmatrix} = \begin{bmatrix} \widetilde{x}_i \end{bmatrix}$$

- If we set $\tilde{x}_i \in \mathbb{R}^k$ to $\tilde{x}_i = \mathbf{V}^T \vec{x}_i$ we have:

$$k \times d \cdot d \times 1 = k \times 1$$

$$\|\tilde{x}_i - \tilde{x}_j\|_2 = \|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2 = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

$$\sqrt{\sum_{\ell=1}^{k} \mathbf{V}^T(x_i - x_j)[\ell]^2}$$

$$\langle V_\ell, x_i - x_j \rangle$$

13

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.

- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

- If we set $\tilde{x}_i \in \mathbb{R}^k$ to $\tilde{x}_i = \mathbf{V}^T \vec{x}_i$ we have:

$$\|\tilde{x}_i - \tilde{x}_j\|_2 = \|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2 = \|\mathbf{V}^T(\vec{x}_i - \vec{x}_j)\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

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

## Embedding with Assumptions

Assume that $\vec{x}_1, \ldots \vec{x}_n$ lie in any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.

- Let $\vec{v}_1, \vec{v}_2, \ldots \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and let $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- For all $i, j$ we have $\vec{x}_i - \vec{x}_j \in \mathcal{V}$ and (a good exercise!):

$$\|\vec{x}_i - \vec{x}_j\|_2 = \sqrt{\sum_{\ell=1}^{k} \langle v_\ell, \vec{x}_i - \vec{x}_j \rangle^2} = \|\mathsf{V}^T(\vec{x}_i - \vec{x}_j)\|_2.$$

- If we set $\tilde{x}_i \in \mathbb{R}^k$ to $\tilde{x}_i = \mathsf{V}^T \vec{x}_i$ we have:
$$\|\tilde{x}_i - \tilde{x}_j\|_2 = \|\mathsf{V}^T \vec{x}_i - \mathsf{V}^T \vec{x}_j\|_2 = \|\mathsf{V}^T(\vec{x}_i - \vec{x}_j)\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

- An embedding with no distortion from any $d$ into $m = k$.
- $\mathsf{V}^T : \mathbb{R}^d \to \mathbb{R}^k$ is a linear map giving our embedding.

13

## Embedding with No Assumptions

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?

## Embedding with No Assumptions

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 $\underline{\epsilon > 0}$?

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

## The Johnson-Lindenstrauss Lemma

**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 $\boldsymbol{\Pi} : \mathbb{R}^d \to \mathbb{R}^m$ such that $m = O\left(\frac{\log n}{\epsilon^2}\right)$ and letting $\tilde{x}_i = \boldsymbol{\Pi}\vec{x}_i$:

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

## The Johnson-Lindenstrauss Lemma

**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 $n = 100,000$, $m \approx 6600$.
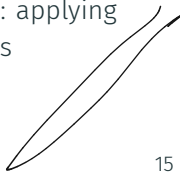
## The Johnson-Lindenstrauss Lemma

> **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 $i, j$ : $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \le \|\tilde{x}_i - \tilde{x}_j\|_2 \le (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.
>
> 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 $n = 100,000$, $m \approx 6600$.
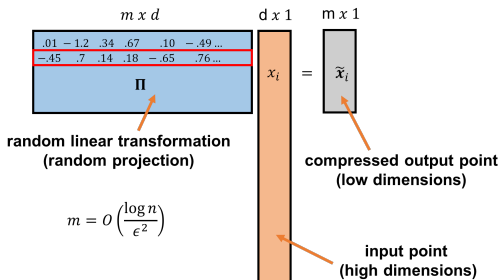
Very surprising! Powerful result with a simple construction: applying a random linear transformation to a set of points preserves distances between all those points with high probability.

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

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



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

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

$\boldsymbol{\Pi}$

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

**random linear transformation (random projection)**

**compressed output point (low dimensions)**

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

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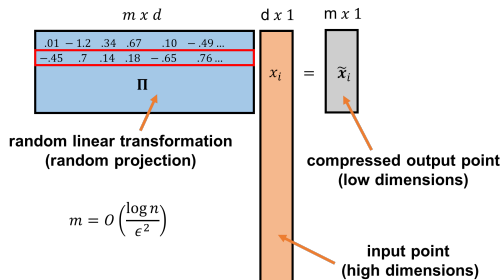
For all $i, j$: $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.



$m \times d$ — $d \times 1$ — $m \times 1$

| .01 | − 1.2 | .34 | .67 | .10 | − .49 ... |

$\mathbf{\Pi}$ — $x_i$ — $=$ — $\tilde{x}_i$

**random linear transformation (random projection)**

**compressed output point (low dimensions)**

$m = O\left(\frac{\log n}{\epsilon^2}\right)$
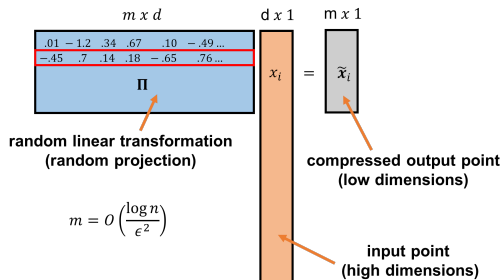
**input point (high dimensions)**

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

# Random Projection

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 $\tilde{\mathbf{x}}_i = \mathbf{\Pi}\vec{x}_i$:

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



$m \times d$    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.

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