

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Spring 2019.

Lecture 10/9

- Problem Set 2 was released Sunday night and is due Sunday 3/8.
- Please make sure to mark all teammates in the GradeScope submission (don't just write names on the document).
- The Midterm will be on Thursday 3/12. Will cover material **through this week.**
- Study guide/practice questions to be released soon.

Last Class:

- Continued on the frequent elements problem.
- Misra-Greis summaries (deterministic method).
- Started on Count-Min Sketch (randomized method).

Last Class:

- Continued on the frequent elements problem.
- Misra-Greis summaries (deterministic method).
- Started on Count-Min Sketch (randomized method).

This Class:

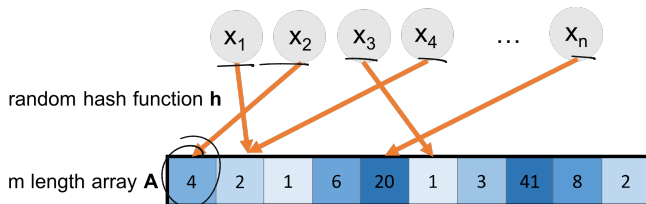
- Finish up Count-Min Sketch analysis.
- Start on randomized methods for dimensionality reduction.
- The Johnson-Lindenstrauss Lemma.

Frequent Items Problem: Identify all items with frequency $\geq n/k$ in a stream of n items. $\geq \frac{n}{k}$ $< (1-\epsilon) \frac{n}{k}$

FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Frequent Items Problem: Identify all items with frequency $\geq n/k$ in a stream of n items.

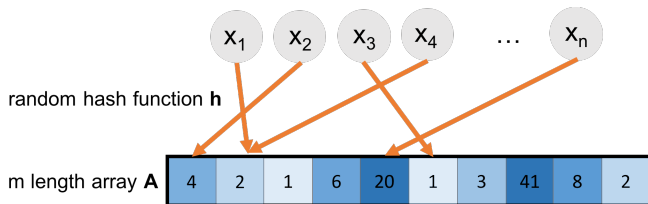
Count-Min Sketch: Bloom filter like solution using random hashing.



FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

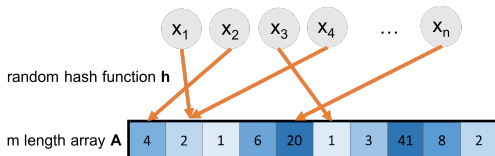
Frequent Items Problem: Identify all items with frequency $\geq n/k$ in a stream of n items.

Count-Min Sketch: Bloom filter like solution using random hashing.



Use $A[h(x)]$ to estimate $\underline{f(x)}$, the frequency of x in the stream.
I.e., $|\{x_i : x_i = x\}|$.

COUNT-MIN SKETCH ACCURACY

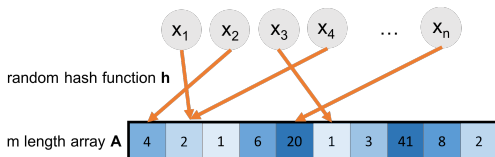


Use $A[\underline{h(x)}]$ to estimate $f(x)$

Claim 1: We always have $A[h(x)] \geq f(x)$. Why?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY



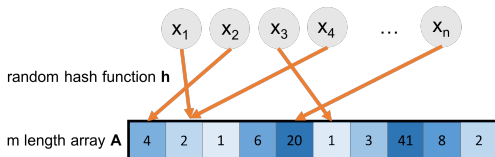
Use $A[h(x)]$ to estimate $f(x)$

Claim 1: We always have $A[h(x)] \geq f(x)$. Why?

- $A[h(x)]$ counts the number of occurrences of any y with $h(y) = h(x)$, including x itself.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY



Use $A[h(x)]$ to estimate $f(x)$

Claim 1: We always have $A[h(x)] \geq f(x)$. Why?

- $A[h(x)]$ counts the number of occurrences of any y with $h(y) = h(x)$, including x itself.

- $A[h(x)] = f(x) + \left\{ \sum_{y \neq x: h(y)=h(x)} f(y) \right\}$

$$A[h(x)] \geq f(x)$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[h(x)] = f(x) + \underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\mathbb{E} \left[\sum_{y \neq x: h(y)=h(x)} f(y) \right] = \sum_y \frac{1}{m} f(y)$$
$$\Pr(h(y) = h(x)) = \frac{1}{m}$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] = \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y)$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &= \sum_{y \neq x} \frac{1}{m} \cdot f(y) \end{aligned}$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[h(x)] = f(x) + \underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

error in frequency estimate $\mathbb{E}[\sum_{y \neq x} f(y)] = \frac{1}{m} \cdot f(y) \cdot (1 - \frac{1}{m}) \cdot 0$

$$\mathbb{E} \left[\underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error in frequency estimate}} \right] = \sum_{y \neq x} \Pr(h(y) = h(x)) \cdot f(y)$$

$$= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m}$$

$$\sum_{\substack{y \neq x \\ h(x)=h(y)}} \mathbb{E} f(y) = \sum_{\substack{y \neq x \\ h(x)=h(y)}} f(y) = \mathbb{E} \left[\sum_{y \neq x} I_y \cdot f(y) \right] = \sum_{\substack{y \neq x \\ h(x)=h(y)}} \mathbb{E} I_y \cdot f(y)$$

$$= \sum_{\substack{y \neq x \\ h(x)=h(y)}} \frac{1}{m} \cdot f(y)$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[h(x)] = f(x) + \underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: h(y)=h(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

What is a bound on probability that the error is $\geq \frac{3n}{m}$?

Markov inequality

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

What is a bound on probability that the error is $\geq \frac{3n}{m}$?

Markov's inequality: $\Pr \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \geq \frac{3n}{m} \right] \leq \frac{1}{3}$.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[h(x)] = f(x) + \underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: h(y)=h(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \left(\frac{n}{m} \right) \end{aligned}$$

2-universal

What is a bound on probability that the error is $\geq \frac{3n}{m}$?

Markov's inequality: $\Pr \left[\sum_{y \neq x: h(y)=h(x)} f(y) \geq \frac{3n}{m} \right] \leq \frac{1}{3}$.

What property of \mathbf{h} is required to show this bound?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

$$A[h(x)] = f(x) + \underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)} \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

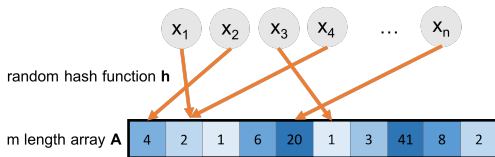
What is a bound on probability that the error is $\geq \frac{3n}{m}$?

Markov's inequality: $\Pr \left[\sum_{y \neq x: h(y)=h(x)} f(y) \geq \frac{3n}{m} \right] \leq \frac{1}{3}$.

What property of \mathbf{h} is required to show this bound? 2-universal.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY

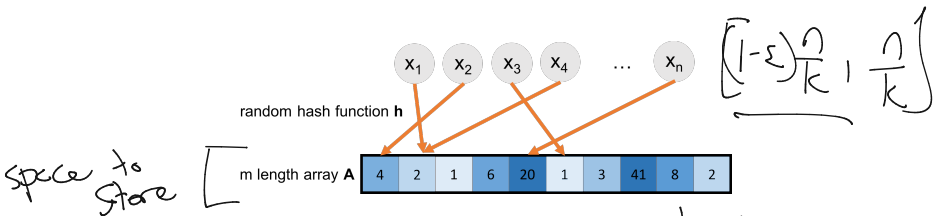


Claim: For any x , with probability at least $2/3$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{3n}{m}.$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY



Claim: For any x , with probability at least $2/3$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{3n}{m}$$

want: $\frac{\epsilon n}{k}$

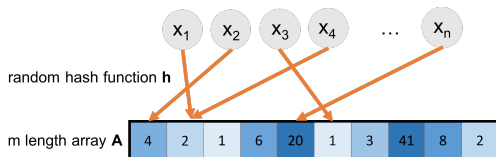
$$\frac{3n}{m} = \frac{\epsilon n}{k}$$

To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{3k}{\epsilon}$.

$$\frac{3n}{m} = \frac{3n}{3k/\epsilon} = \frac{\epsilon n}{k} \quad \text{space: } O\left(\frac{k}{\epsilon} \cdot \log n\right)$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY



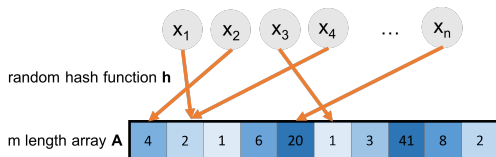
Claim: For any x , with probability at least $2/3$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{3k}{\epsilon}$.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY



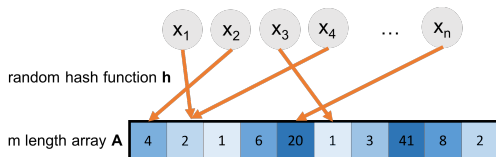
Claim: For any x , with probability at least $2/3$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{3k}{\epsilon}$.
How can we improve the success probability?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

COUNT-MIN SKETCH ACCURACY



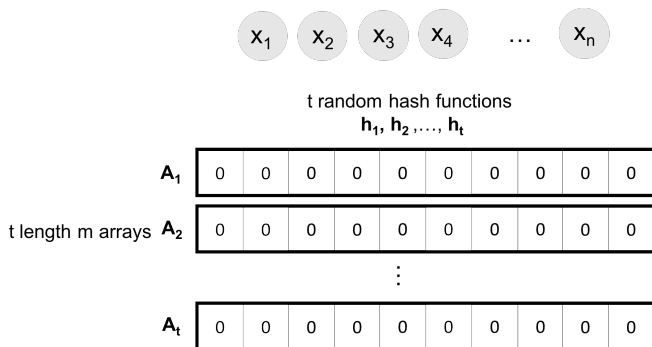
Claim: For any x , with probability at least $2/3$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

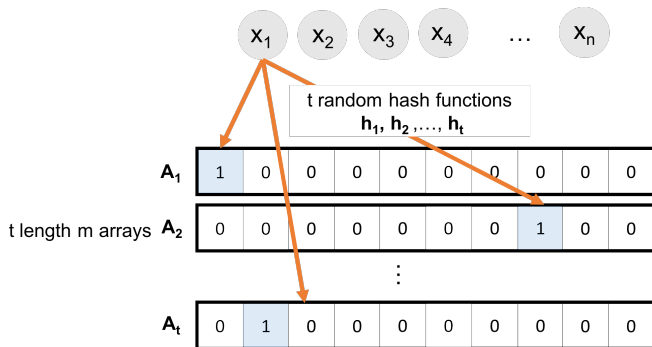
To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{3k}{\epsilon}$.
How can we improve the success probability? **Repetition.**

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of count-min sketch array.

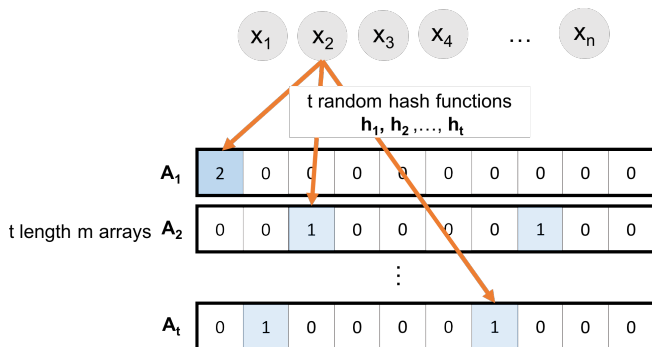
COUNT-MIN SKETCH ACCURACY



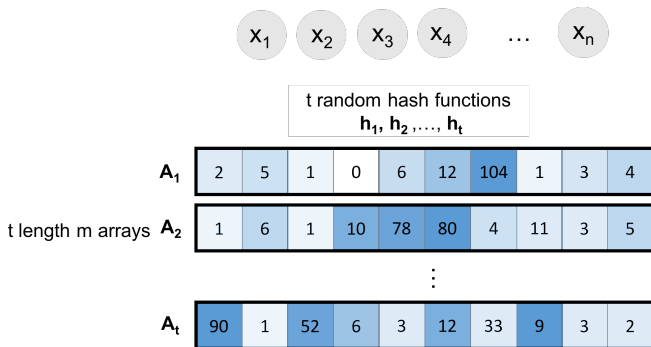
COUNT-MIN SKETCH ACCURACY



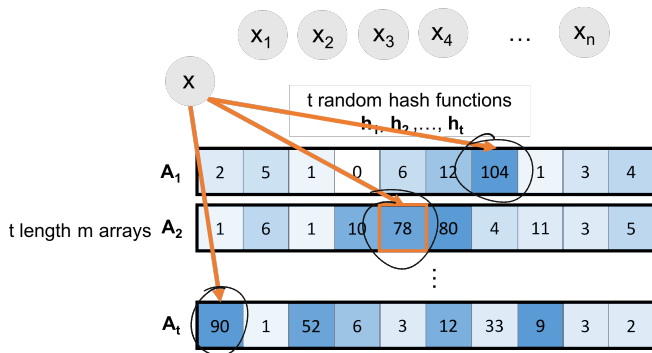
COUNT-MIN SKETCH ACCURACY



COUNT-MIN SKETCH ACCURACY

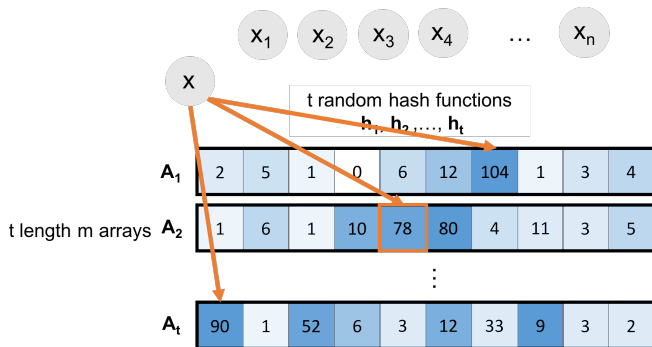


COUNT-MIN SKETCH ACCURACY



Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

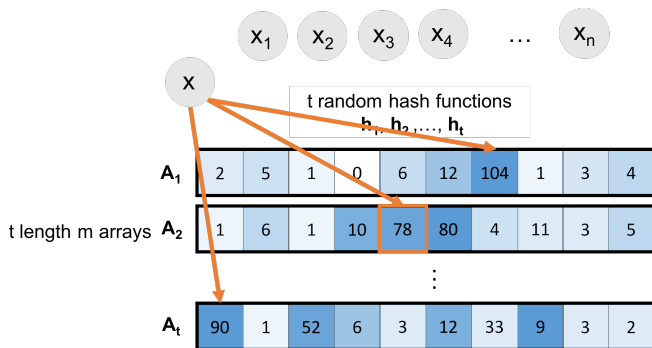
COUNT-MIN SKETCH ACCURACY



Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

Why min instead of mean or median?

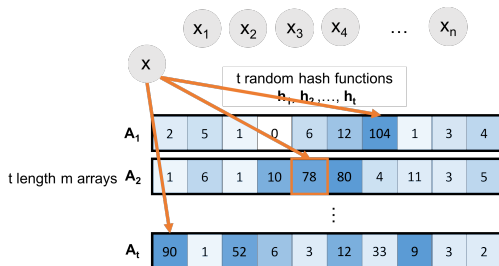
COUNT-MIN SKETCH ACCURACY



Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

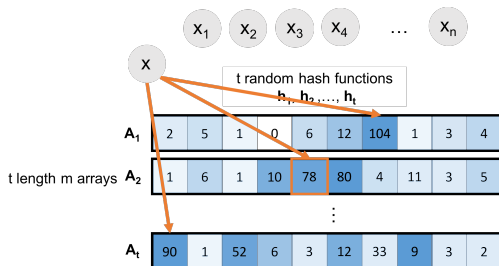
Why min instead of mean or median? The minimum estimate is always the most accurate since they are all overestimates of the true frequency!

COUNT-MIN SKETCH ANALYSIS



Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

COUNT-MIN SKETCH ANALYSIS

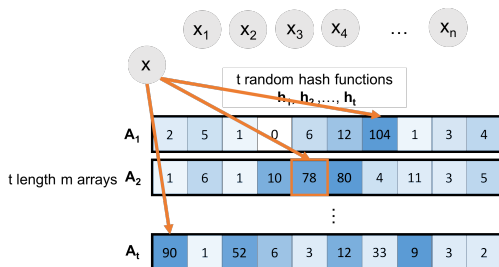


Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = O(k/\epsilon)$, with probability $\geq 2/3$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

COUNT-MIN SKETCH ANALYSIS



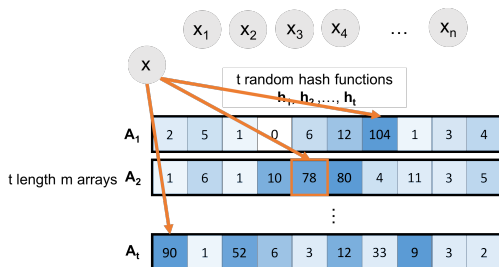
Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = O(k/\epsilon)$, with probability $\geq 2/3$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$?

COUNT-MIN SKETCH ANALYSIS



Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = O(k/\epsilon)$, with probability $\geq 2/3$:

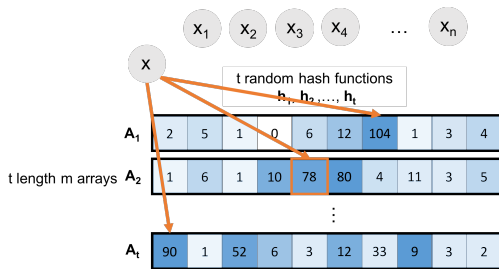
$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$? $1 - 1/3^t$.

COUNT-MIN SKETCH ANALYSIS

$$\log(1/\delta)$$

$$\frac{1}{\epsilon^2} \quad \frac{1}{\epsilon}$$



$$t \geq \frac{\log(d)}{\log(3)}$$

Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = O(k/\epsilon)$, with probability $\geq 2/3$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}$$

$$1 - 1/3^t \geq 1 - \delta$$

$$\delta \geq 1/3^t$$

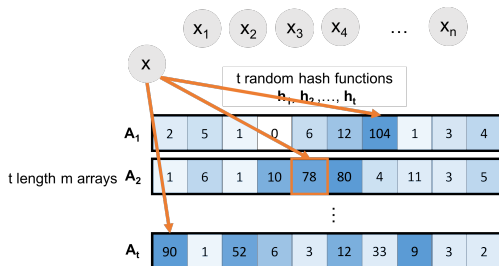
$$\log(\delta) \geq \log(1/3^t) = t \cdot \log(1/3)$$

$$\log(\delta) \geq t \cdot \log(3)$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$? $1 - 1/3^t$.

- To get a good estimate with probability $\geq 1 - \delta$, set $t = O(\log(1/\delta))$.

COUNT-MIN SKETCH ANALYSIS



Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = O(k/\epsilon)$, with probability $\geq 2/3$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$? $1 - 1/3^t$.
- To get a good estimate with probability $\geq 1 - \delta$, set $t = O(\log(1/\delta))$.

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

$\bullet \log(n)$

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

- Accurate enough to solve the (ϵ, k) -Frequent elements problem.

COUNT-MIN SKETCH

$$1000 / .01 \approx 10,000$$

[1000]

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

- Accurate enough to solve the (ϵ, k) -Frequent elements problem. $\frac{\epsilon n}{k}$
- Actually identifying the frequent elements quickly requires a little bit of further work.

One approach: Separately store a list of potential frequent elements as they come in. At step i , keep any elements whose estimated frequency is $\geq i/k$. List contains at most $O(k)$ items at any step and has all items with frequency $\geq n/k$ stored at the end of the stream.

Questions on Frequent Elements?

'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 500×500 pixels at 15 frames/second with 3 color channels is a recording of **≥ 2 billion pixel values**. Even a 500×500 pixel color image has 750,000 pixel values.

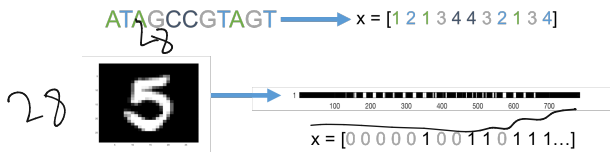
'Big Data' means not just many data points, but many measurements per data point. I.e., very **high dimensional data**.

- 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 500×500 pixels at 15 frames/second with 3 color channels is a recording of **≥ 2 billion pixel values**. Even a 500×500 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.

DATA AS VECTORS AND MATRICES

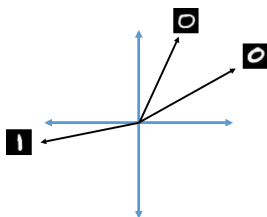
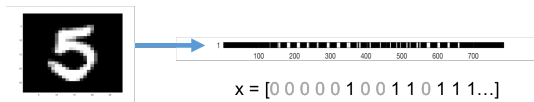
In data analysis and machine learning, data points with many attributes are often stored, processed, and interpreted as **high dimensional vectors**, with real valued entries.



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 $x = [1\ 2\ 1\ 3\ 4\ 4\ 3\ 2\ 1\ 3\ 4]$



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: $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbb{R}^d$.

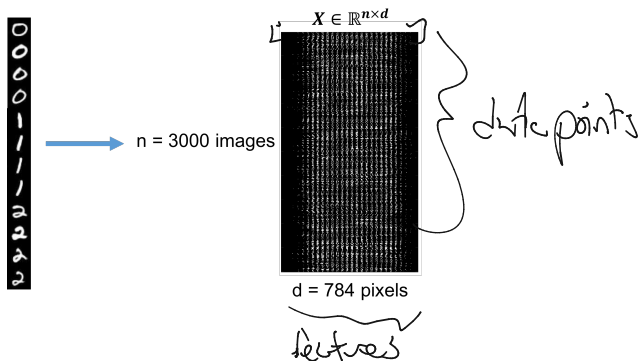
Data Set: $X \in \mathbb{R}^{n \times d}$ with i^{th} rows equal to \vec{x}_i .

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: $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbb{R}^d$.

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

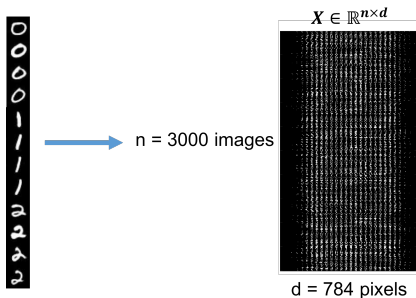


DATASETS AS VECTORS AND MATRICES

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

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

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



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.

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

5 $\rightarrow x = [00000100110111\dots] \rightarrow \tilde{x} = [-5.5 \ 4 \ 3.2 \ -1]$

DIMENSIONALITY REDUCTION

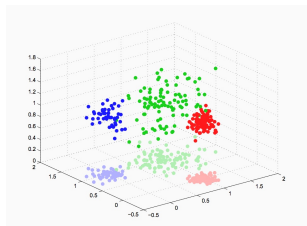
Dimensionality Reduction: Compress data points so that they lie in many fewer dimensions.

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

5

$$\rightarrow x = [00000100110111\dots] \rightarrow \tilde{x} = [-5.5 \ 4 \ 3.2 \ -1]$$

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



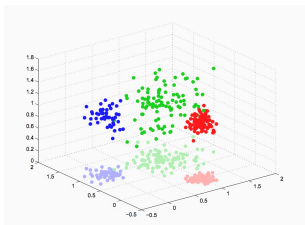
DIMENSIONALITY REDUCTION

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

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

5 $\rightarrow x = [0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 1\ 1\ \dots]$ $\rightarrow \tilde{x} = [-5.5\ 4\ 3.2\ -1]$

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



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

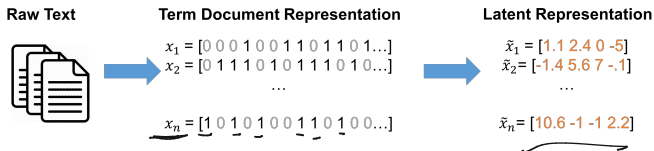
Dimensionality reduction is one of the most important techniques in data science.

Dimensionality reduction is one of the most important techniques in data science.

- Principal component analysis

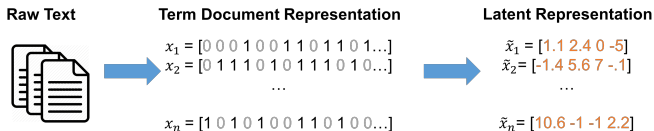
Dimensionality reduction is one of the most important techniques in data science.

- Principal component analysis
- Latent semantic analysis (LSA)



Dimensionality reduction is one of the most important techniques in data science.

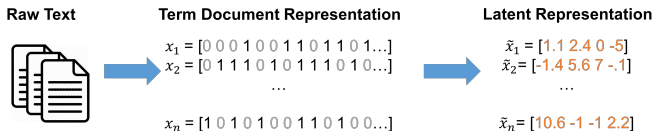
- Principal component analysis
- Latent semantic analysis (LSA)



- Linear discriminant analysis

Dimensionality reduction is one of the most important techniques in data science.

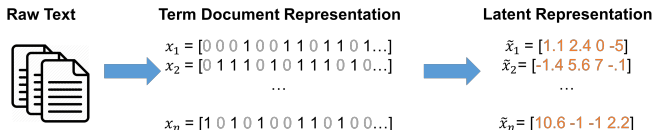
- Principal component analysis
- Latent semantic analysis (LSA)



- Linear discriminant analysis
- Autoencoders

Dimensionality reduction is one of the most important techniques in data science.

- Principal component analysis
- Latent semantic analysis (LSA)



- Linear discriminant analysis
- Autoencoders

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

Low Distortion Embedding: Given $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$, distance function D , and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (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

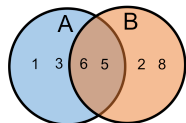
Low Distortion Embedding: Given $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$, distance function D , and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) and distance function \tilde{D} such that for all $i, j \in [n]$:

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

Have already seen one example in class: **MinHash**.

$J(A, B) = \frac{\text{# overlapping elements}}{\text{# elements}}$

$x_A = [1 0 1 0 1 1 0 0]$



$x_B = [0 1 0 0 1 1 0 1]$

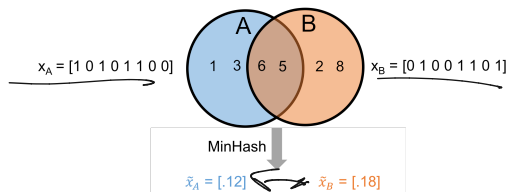


LOW DISTORTION EMBEDDING

Low Distortion Embedding: Given $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$, distance function D , and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (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**.

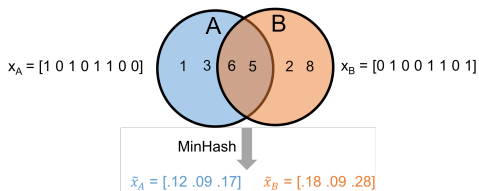


LOW DISTORTION EMBEDDING

Low Distortion Embedding: Given $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$, distance function D , and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (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**.

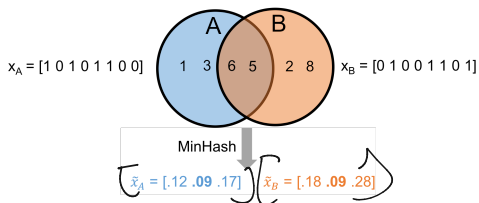


LOW DISTORTION EMBEDDING

Low Distortion Embedding: Given $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$, distance function D , and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (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**.



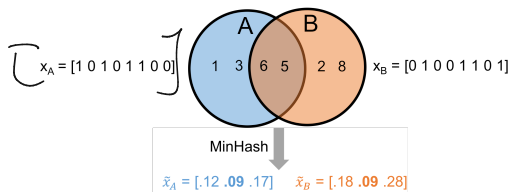
With large enough signature size r , $\frac{\# \text{ matching entries in } \tilde{x}_A, \tilde{x}_B}{r} \approx J(\vec{x}_A, \vec{x}_B)$.

LOW DISTORTION EMBEDDING

Low Distortion Embedding: Given $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$, distance function D , and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (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 , $\frac{\# \text{ matching entries in } \tilde{x}_A, \tilde{x}_B}{r} \approx J(\vec{x}_A, \vec{x}_B)$.

- Reduce dimension from $d = |U|$ to r . 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 $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ and error parameter $\epsilon \geq 0$, find $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$ (where $m \ll d$) such that for all $i, j \in [n]$:

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

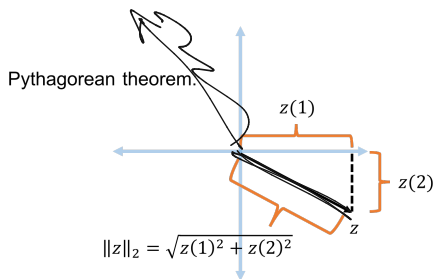
Recall that for $\vec{z} \in \mathbb{R}^n$, $\|\vec{z}\|_2 = \sqrt{\sum_{i=1}^n z(i)^2}$.

D, \tilde{D} are
euclidean distance

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

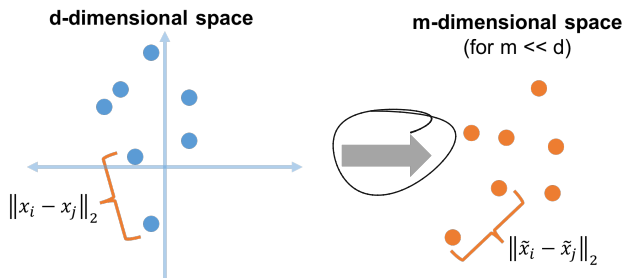
$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$

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



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

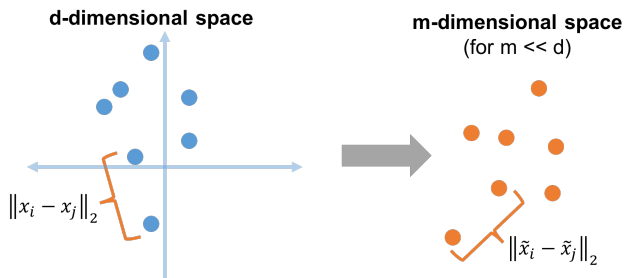
$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$



EMBEDDINGS FOR EUCLIDEAN SPACE

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

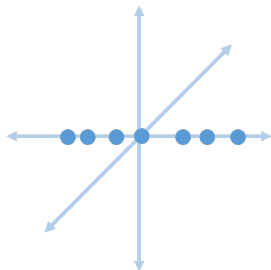
$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$



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

EMBEDDING WITH ASSUMPTIONS

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



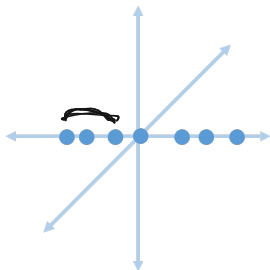
EMBEDDING WITH ASSUMPTIONS

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

$$\begin{bmatrix} x_i(1) & 0 & 0 & 0 \end{bmatrix}$$

↓

$$\tilde{x}_i = [x_i(1)]$$

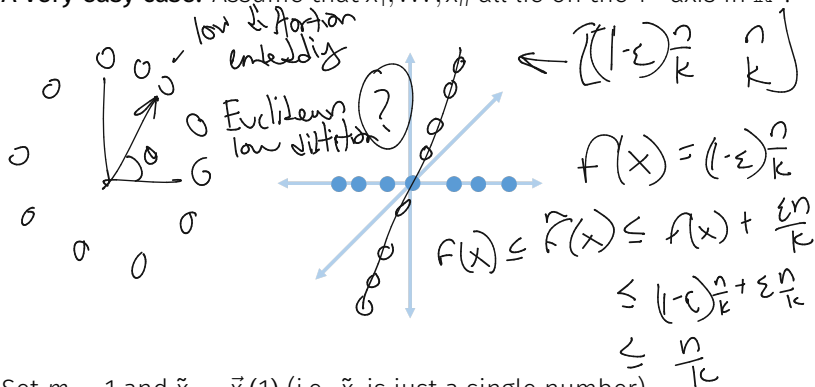


Set $m = 1$ and $\tilde{x}_i = \vec{x}_i(1)$ (i.e., \tilde{x}_i is just a single number).

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

EMBEDDING WITH ASSUMPTIONS

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



Set $m = 1$ and $\tilde{x}_i = \vec{x}_i(1)$ (i.e., \tilde{x}_i is 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 $m = 1$.