CS-GY 9223 I: Lecture 4
Near neighbor search + locality sensitive hashing

NYU Tandon School of Engineering, Prof. Christopher Musco

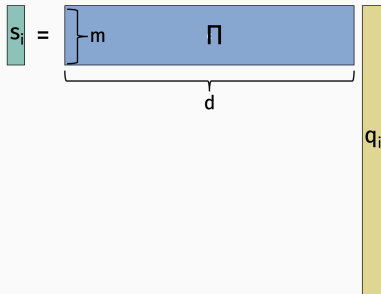- Problem set 1.
- Reading group.

Lemma (Johnson-Lindenstrauss, 1984)

*For any set of n data points $\mathbf{q}_1, \ldots, \mathbf{q}_n \in \mathbb{R}^d$ there exists a <u>linear map</u> $\Pi : \mathbb{R}^d \to \mathbb{R}^m$ where $m = O\left(\frac{\log n}{\epsilon^2}\right)$ such that <u>for all i, j</u>,*

$$(1 - \epsilon)\|\mathbf{q}_i - \mathbf{q}_j\|_2 \le \|\Pi\mathbf{q}_i - \Pi\mathbf{q}_j\|_2 \le (1 + \epsilon)\|\mathbf{q}_i - \mathbf{q}_j\|_2.$$



3

$\Pi \in \mathbb{R}^{k \times d}$ be chosen so that each entry equals $\frac{1}{\sqrt{m}} \mathcal{N}(0, 1)$.

… or each entry equals $\frac{1}{\sqrt{m}} \pm 1$ with equal probability.



```
>> Pi = randn(m,d);
>> s = (1/sqrt(m))*Pi*q;
```



```
>> Pi = 2*randi(2,m,d)-3;
>> s = (1/sqrt(m))*Pi*q;
```

Often called "random projections". Why?

4

$k$-**means objective**: Find clusters $C_1, \ldots, C_k \subseteq \{1, \ldots, n\}$ to minimize:

$$Cost(C_1, \ldots, C_k) = \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|\mathsf{X}_u - \mathsf{X}_v\|_2^2.$$

**Approximation algorithm**: Find optimal clusters $\tilde{C}_1^*, \ldots, \tilde{C}_k^*$ for the $k$ dimension data set $\Pi X_1, \ldots, \Pi X_n$.

$$Cost(C_1, \ldots, C_k) = \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|\mathsf{x}_u - \mathsf{x}_v\|_2^2.$$

$$\widetilde{Cost}(C_1, \ldots, C_k) = \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|\Pi\mathsf{x}_u - \Pi\mathsf{x}_v\|_2^2.$$

For any $C_1, \ldots, C_k$,

$$(1 - \epsilon)Cost(C_1, \ldots, C_k) \leq \widetilde{Cost}(C_1, \ldots, C_k) \leq (1 + \epsilon)Cost(C_1, \ldots, C_k)$$

Let $C_1^*, \ldots, C_k^* = \arg\min Cost(C_1, \ldots, C_k)$ and
$\tilde{C}_1^*, \ldots, \tilde{C}_k^* = \arg\min \widetilde{Cost}(C_1, \ldots, C_k)$

**Want to prove:** $Cost(\tilde{C}_1^*, \ldots, \tilde{C}_k^*) \leq (1 + O(\epsilon))Cost(C_1^*, \ldots, C_k^*)$

input data

high dimensional vector representation

| 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |

| .45 | .68 | .10 | .92 |

sketched representation

**Goal:** Given input vectors q and y, $C(q)$ and $C(y)$ should be similar if q and y are similar.

Other Example: Binary valued vectors.

Definition (Jaccard Similarity)

$$J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|} = \frac{\text{\# of non-zero entries in common}}{\text{total \# of non-zero entries}}$$

$0 \leq J(\mathbf{q}, \mathbf{y}) \leq 1$.

- Choose $k$ random hash functions
  $h_1, \ldots, h_k : \{1, \ldots, n\} \to [0, 1]$.
- For $i \in 1, \ldots, k$, let $c_i = \min_{j, \mathbf{q}_j = 1} h_i(j)$.
- $C(\mathbf{q}) = [c_1, \ldots, c_k]$.

**Example 1:** Binary valued vectors.

If $J(\mathbf{q}, \mathbf{y}) = v$ then the expected number of common entries between $C(\mathbf{q})$ and $C(\mathbf{y})$ is $v$.

| $C(\mathbf{q})$ | .12 | .24 | .76 | .35 |
|---|---|---|---|---|

| $C(\mathbf{y})$ | .12 | .98 | .76 | .11 |
|---|---|---|---|---|

Using a Chernoff bound, we proved that if $C$ maps to dimension $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, we can approximate the Jaccard similarity between any two binary vectors to accuracy $\epsilon$ with probability $1 - \delta$.

**Common goal:** Find all fingerprints in database $q_1, \ldots, q_n \in \mathbb{R}^d$ that are close to some input finger print $y \in \mathbb{R}^d$.

- Audio + video search.
- Finding duplicate or near duplicate documents.
- Seismic applications (here they want all pairs of close fingerprints).

Does similarity sketching help in these applications?

New goal: <u>Sublinear</u> $o(n)$ time to find near neighbors.

Let $h : \mathbb{R}^d \to \{1, \ldots, m\}$ be a random hash function.

We call $h$ <u>locality sensitive</u> if $\Pr[h(\mathbf{q}) == h(\mathbf{y})]$ is:

- Higher when $\mathbf{q}$ and $\mathbf{y}$ are more similar.
- Lower when $\mathbf{q}$ and $\mathbf{y}$ differ substantially.

LSH for Jaccard similarity:

- Let $c : \{0,1\}^d \to [0,1]$ be a single instantiation of MinHash.
- Let $g : [0,1] \to \{1,\ldots,m\}$ be a fully random hash function.
- Let $h(\mathbf{x}) = g(c(\mathbf{x}))$.

LSH for Jaccard similarity:

- Let $c : \{0, 1\}^d \rightarrow [0, 1]$ be a single instantiation of MinHash.
- Let $g : [0, 1] \rightarrow \{1, \ldots, m\}$ be a fully random hash function.
- Let $h(\mathbf{x}) = g(c(\mathbf{x}))$.

If $J(\mathbf{q}, \mathbf{y}) = v$,

$$\Pr\left[h(\mathbf{q}) == h(\mathbf{y})\right] =$$

Basic approach for near neighbor search in a database.

Basic approach for near neighbor search in a database.

## Pre-processing:

- Select random LSH function $h : \{0,1\}^d \to 1, \ldots, m$.
- Create table $T$ with $m$ slots.
- For $i = 1, \ldots, n$, insert $\mathbf{q}_i$ into $T(h(\mathbf{q}_i))$.
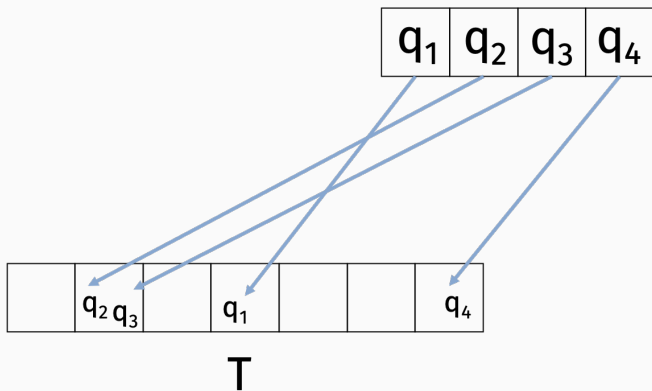
Basic approach for near neighbor search in a database.

### Pre-processing:

- Select random LSH function $h : \{0, 1\}^d \to 1, \dots, m$.
- Create table $T$ with $m$ slots.
- For $i = 1, \dots, n$, insert $\mathbf{q}_i$ into $T(h(\mathbf{q}_i))$.

### Query:

- Want to find near neighbors of input $\mathbf{y} \in \{0, 1\}^d$.
- Linear scan through all vectors in $T(h(\mathbf{y}))$.

T

Two main considerations:
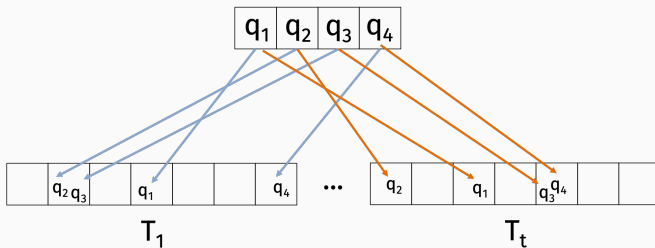
- **False Negative Rate**: What's the probability we do not find a vector that <u>is close</u> to $y$?
- **False Positive Rate**: What's the probability we need to scan over vectors that <u>are not close</u> to $y$?

Suppose the nearest database point $\mathbf{q}$ has $J(\mathbf{y}, \mathbf{q}) = .4$.

What's the probability we do not find q?

Pre-processing:

- Select $t$ independent LSH's $h_1, \ldots, h_t : \{0,1\}^d \to 1, \ldots, m$.
- Create tables $T_1, \ldots, T_t$, each with $m$ slots.
- For $i = 1, \ldots, n, j = 1, \ldots, t$, insert $\mathbf{q}_i$ into $T_j(h_j(\mathbf{q}_i))$.

Query:

- Want to find near neighbors of input $y \in \{0, 1\}^d$.
- Linear scan through all vectors in $T_1(h_1(y)), T_2(h_2(y)), \ldots, T_t(h_t(y))$.

**Query:**

- Want to find near neighbors of input $\mathbf{y} \in \{0, 1\}^d$.
- Linear scan through all vectors in
  $T_1(h_1(\mathbf{y})), T_2(h_2(\mathbf{y})), \ldots, T_t(h_t(\mathbf{y}))$.

Suppose the nearest database point $\mathbf{q}$ has $J(\mathbf{y}, \mathbf{q}) = .4$.

What's the probability we find q?

Suppose there is some other database point $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .2$? What is the probability we will consider that point in our original scheme?

In the new scheme?

Change our locality sensitive hash function.

If $J(\mathbf{q}, \mathbf{y}) = v$,

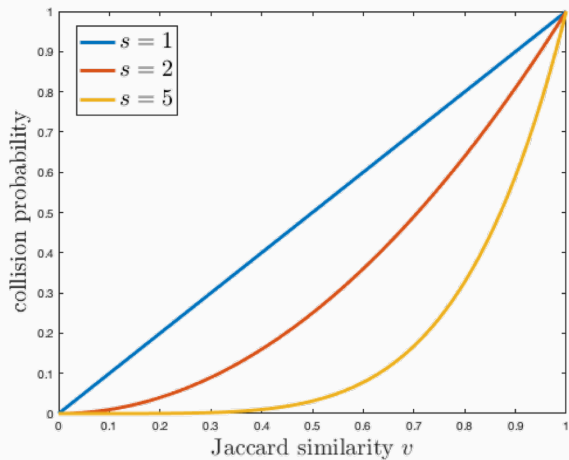$$\Pr[h(\mathbf{q}) == h(\mathbf{y})] =$$

Change our locality sensitive hash function.

Tunable LSH for Jaccard similarity:

- Choose parameter $s \in \mathbb{Z}^+$.
- Let $c_1, \ldots, c_s : \{0, 1\}^d \to [0, 1]$ be random MinHashs.
- Let $g : [0, 1]^s \to \{1, \ldots, m\}$ be a fully random hash function.
- Let $h(\mathbf{x}) = g(c_1(\mathbf{x}), \ldots, c_s(\mathbf{x}))$.

If $J(\mathbf{q}, \mathbf{y}) = v$,

$$\Pr[h(\mathbf{q}) == h(\mathbf{y})] =$$

Parameter: $S = 1$.

Chance we find $\mathbf{q}_i$ with $J(\mathbf{y}, \mathbf{q}_i) = .8$:

Chance we need to scan $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .4$:

Parameter: $S = 2$.

Chance we find $\mathbf{q}_i$ with $J(\mathbf{y}, \mathbf{q}_i) = .8$:

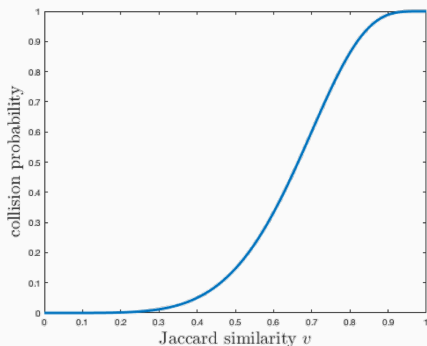Chance we need to scan $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .4$:

Parameter: $S = 5$.

Chance we find $\mathbf{q}_i$ with $J(\mathbf{y}, \mathbf{q}_i) = .8$:

Chance we need to scan $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .4$:

Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$1 - (1 - v^s)^t$$



$s = 5, t = 5$

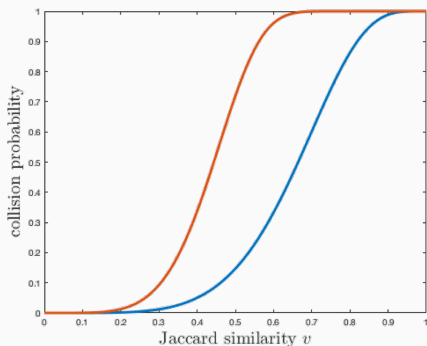Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$1 - (1 - v^s)^t$$
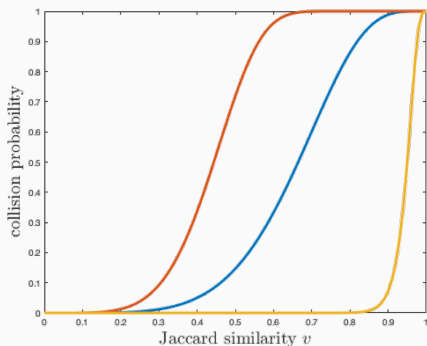


$$s = 5, t = 40$$

Probability we see $\mathbf{q}$ when querying $\mathbf{y}$ if $J(\mathbf{q}, \mathbf{y}) = v$:

$$\approx 1 - (1 - v^s)^t$$



$s = 40, t = 5$
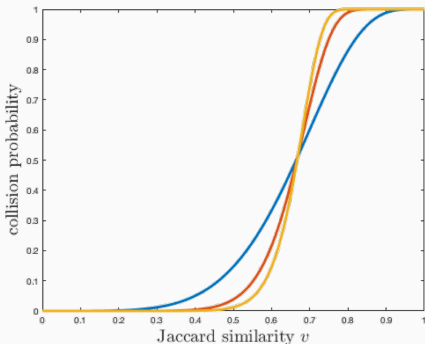
Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$1 - (1 - v^s)^t$$



Increasing both *s* and *t* gives a steeper curve.

Better for search, but worse space complexity.

33

**Use Case 1:** Fixed threshold.

- Shazam wants to find match to audio clip **y** in a database of 10 million clips.
- There are 10 <u>true matches</u> with $J(\mathbf{y}, \mathbf{q}) > .9$.
- There are 10,0000 <u>near matches</u> with $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$.

With $s = 25$ and $t = 40$,

- Hit probability for $J(\mathbf{y}, \mathbf{q}) > .9$ is $\gtrsim 1 - (1 - .9^{25})^{40} = .95$
- Hit probability for $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$ is $\lesssim 1 - (1 - .9^{25})^{40} = .95$
- Hit probability for $J(\mathbf{y}, \mathbf{q}) < .7$ is $\lesssim 1 - (1 - .7^{25})^{40} = .005$

**Total number of items scanned:**

$.95 \cdot 10 + .95 \cdot 10,000 + .005 \cdot 9,989,990 \approx 60,000 \ll 10,000,000$.

Space complexity: 40 hash tables $\approx 40 \cdot O(n)$.

Directly trade space for fast search.

Concrete worst case result:

### Theorem (Indyk, Motwani, 1998)

*If there exists some q with $\|q - y\|_0 \leq R$, return a vector $\tilde{q}$ with $\|\tilde{q} - y\|_0 \leq C \cdot R$ in:*

- *Time: $O\left(n^{1/C}\right)$.*
- *Space: $O\left(n^{1+1/C}\right)$.*

$\|q - y\|_0 = $ "hamming distance" = number of elements that differ between q and y.

Theorem (Indyk, Motwani, 1998)

*Let $q$ be the closest database vector to $\mathbf{y}$. Return a vector $\tilde{\mathbf{q}}$ with $\|\tilde{\mathbf{q}} - \mathbf{y}\|_0 \leq C \cdot \|\mathbf{q} - \mathbf{y}\|_0$ in:*

- *Time: $\tilde{O}\left(n^{1/C}\right)$.*
- *Space: $\tilde{O}\left(n^{1+1/C}\right)$.*
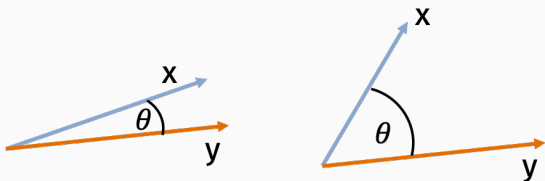
Any ideas for how this is done?

Good locality sensitive hash functions exists for many other similarity measures.

Good locality sensitive hash functions exists for many other similarity measures.

**Cosine similarity** $\cos\left(\theta(\mathbf{x}, \mathbf{y})\right) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$:



$$-1 \leq \cos\left(\theta(\mathbf{x}, \mathbf{y})\right) \leq 1.$$

Cosine similarity is natural "inverse" for Euclidean distance.
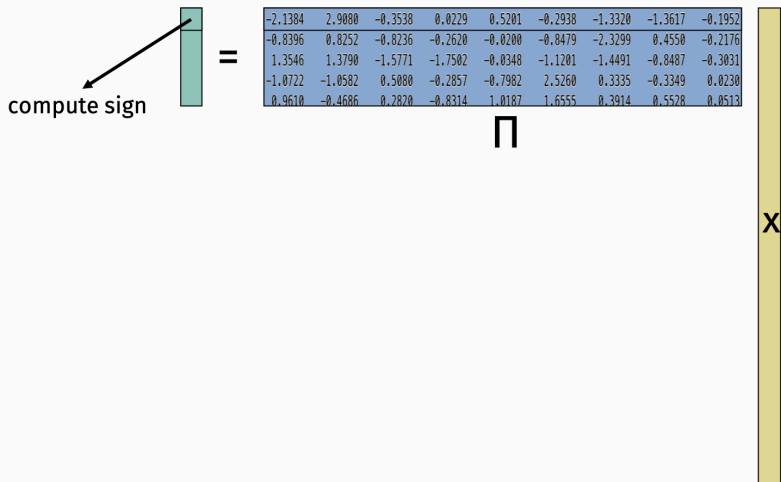
**Euclidean distance** $\|x - y\|_2^2$**:**

- Suppose for simplicity that $\|x\|_2^2 = \|y\|_2^2 = 1$.

Locality sensitive hash for cosine similarity:
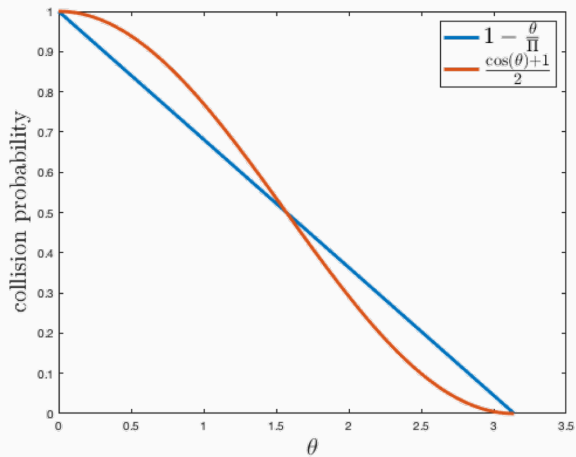
- Let $\mathbf{g} \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0,1)$.
- $h : \mathbb{R}^d \to \{-1, 1\}$ is definied $h(\mathbf{x}) = \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)$.

  If $\cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $\Pr[h(\mathbf{x} == h(\mathbf{y})]$?

Inspired by Johnson-Lindenstrauss sketching



| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| −2.1384 | 2.9080 | −0.3538 | 0.0229 | 0.5201 | −0.2938 | −1.3320 | −1.3617 | −0.1952 |
| −0.8396 | 0.8252 | −0.8236 | −0.2620 | −0.0200 | −0.8479 | −2.3299 | 0.4550 | −0.2176 |
| 1.3546 | 1.3790 | −1.5771 | −1.7502 | −0.0348 | −1.1201 | −1.4491 | −0.8487 | −0.3031 |
| −1.0722 | −1.0582 | 0.5080 | −0.2857 | −0.7982 | 2.5260 | 0.3335 | −0.3349 | 0.0230 |
| 0.9610 | −0.4686 | 0.2820 | −0.8314 | 1.0187 | 1.6555 | 0.3914 | 0.5528 | 0.0513 |

**=**

**Π**

compute sign

**X**

Locality sensitive hash for cosine similarity:

- Let $\mathbf{g} \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- $h : \mathbb{R}^d \to \{-1, 1\}$ is definied $h(\mathbf{x}) = \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)$.

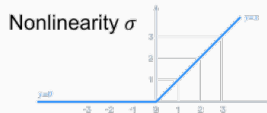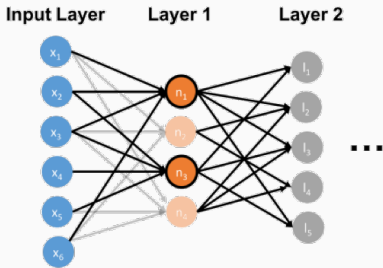If $\cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $\Pr[h(\mathbf{x} == h(\mathbf{y})]$?

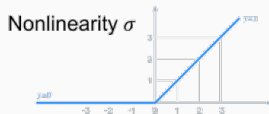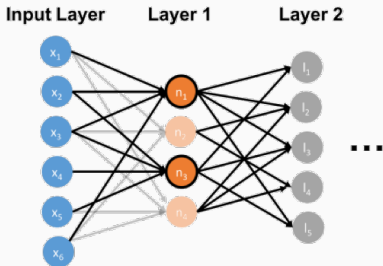Work of Anshumali Shrivastava at Rice University and coauthors.



$$n_i = \sigma\left(\sum_{j=1}^{m} w(x_j, n_i) \cdot x_j\right) = \sigma(\langle w_i, x\rangle)$$

- Number of multiplications to evaluate $\mathcal{N}(\mathbf{x})$:
  $|\mathbf{x}| \cdot |\text{layer 1}| + |\text{layer 1}| \cdot |\text{layer 2}| + |\text{layer 2}| \cdot |\text{layer 3}| + \ldots$
- For an approximate solution, only consider neurons on each each with high activation.

45

Work of Anshumali Shrivastava at Rice University and coauthors.



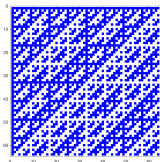$$n_i = \sigma\left(\sum_{j=1}^{m} w(x_j, n_i) \cdot x_j\right) = \sigma(\langle w_i, x \rangle)$$

- High activation = large value of $\sigma(\langle w_i, x \rangle)$.
- Typically $\sigma(\langle w_i, x \rangle)$ increases as $\langle w_i, x \rangle$ increases.
- Use LSH/SimHash to quickly find all $w_i$ for which $\langle w_i, x \rangle$ is large and only include these terms in the sum.
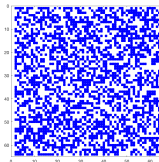
46

Why can't we just sample entries from vectors?

[Ailon, Chazelle, 2009 – The Fast Johnson-Lindenstrauss Transform]

[Ailon, Chazelle, 2009 – The Fast Johnson-Lindenstrauss Transform]

Deterministic Hadamard matrix.

Randomized Hadamard $HD$.

Fully random sign matrix.