

CS-GY 6763: Lecture 5

Near neighbor search in high dimensions

NYU Tandon School of Engineering, Prof. Christopher Musco

- Sign-up to present or lead discussion for 1 reading group slot. We need presenters for next Friday!

Lemma (Distributional JL Lemma)

Let $\Pi \in \mathbb{R}^{k \times d}$ be a random Gaussian/sign matrix. For any two real-valued vectors $\mathbf{q}, \mathbf{y} \in \mathbb{R}^d$, then with probability $1 - \delta$,

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

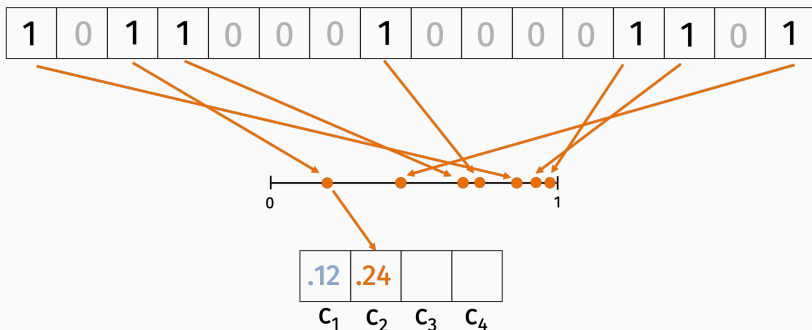
as long as $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$.

LAST CLASS: MINHASH SKETCHES FOR BINARY VECTORS

- Choose k random hash functions

$$h_1, \dots, h_k : \{1, \dots, n\} \rightarrow [0, 1].$$

- For $i \in 1, \dots, k$,
 - Let $c_i = \min_{j, q_j=1} h_i(j)$.
- $C(\mathbf{q}) = [c_1, \dots, c_k]$.



LAST CLASS: MINHASH SKETCHES FOR BINARY VECTORS

Let $\tilde{J}(C(\mathbf{q}), C(\mathbf{y})) = \frac{1}{k} \sum_{i=1}^k \mathbb{1}[C(\mathbf{q})_i = C(\mathbf{y})_i]$.

Lemma (Distributional JL Lemma)

For any two binary vectors $\mathbf{q}, \mathbf{y} \in \mathbb{R}^d$, with probability $1 - \delta$,

$$J(\mathbf{q}, \mathbf{y}) - \epsilon \leq \tilde{J}(C(\mathbf{q}), C(\mathbf{y})) \leq J(\mathbf{q}, \mathbf{y}) + \epsilon,$$

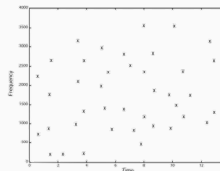
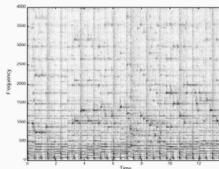
as long as $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$.

Recall that $J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|}$.

SIMILARITY SKETCHING



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

NEAR NEIGHBOR SEARCH

Common goal: Find all vectors in database $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^d$ that are close to some input query vector $\mathbf{y} \in \mathbb{R}^d$. I.e. find all of \mathbf{y} 's “nearest neighbors” in the database.

- The Shazam problem.
- Audio + video search.
- Finding duplicate or near duplicate documents.
- Detecting seismic events.

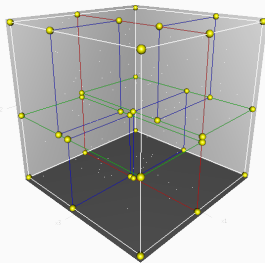
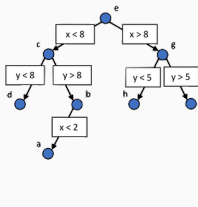
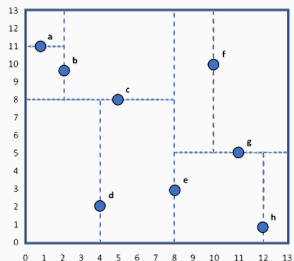
How does similarity sketching help in these applications?

- Improves runtime of “linear scan” from $O(nd)$ to $O(nk)$.
- Improves space complexity from $O(nd)$ to $O(nk)$. This can be super important – e.g. if it means the linear scan only accesses vectors in fast memory.

New goal: Sublinear $o(n)$ time to find near neighbors.

BEYOND A LINEAR SCAN

This problem can already be solved for a small number of dimensions using space partitioning approaches (e.g. kd-tree).



Runtime is roughly $O(d \cdot \min(n, 2^d))$, which is only sublinear for $d = o(\log n)$.

Only been attacked much more recently:

- **Locality-sensitive hashing [Indyk, Motwani, 1998]**
- Spectral hashing [Weiss, Torralba, and Fergus, 2008]
- Vector quantization [Jégou, Douze, Schmid, 2009]

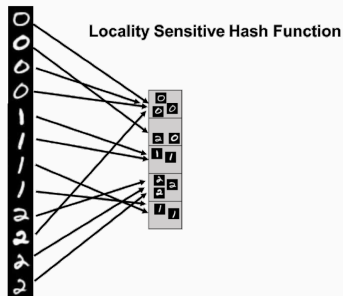
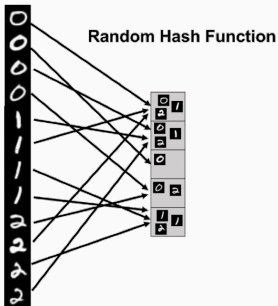
Key Insight of LSH: Trade worse space-complexity for better time-complexity. I.e. typically use more than $O(n)$ space.

LOCALITY SENSITIVE HASH FUNCTIONS

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

We call h locality sensitive for similarity function $s(\mathbf{q}, \mathbf{y})$ if $\Pr[h(\mathbf{q}) == h(\mathbf{y})]$ is:

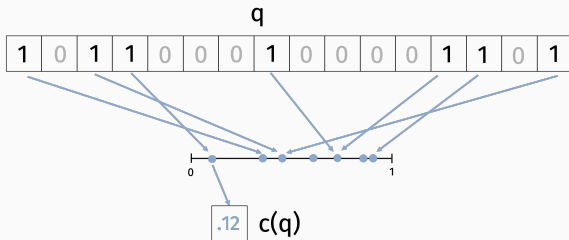
- Higher when \mathbf{q} and \mathbf{y} are more similar, i.e. $s(\mathbf{q}, \mathbf{y})$ is higher.
- Lower when \mathbf{q} and \mathbf{y} are more dissimilar, i.e. $s(\mathbf{q}, \mathbf{y})$ is lower.



LOCALITY SENSITIVE HASH FUNCTIONS

LSH for $s(\mathbf{q}, \mathbf{y})$ equal to Jaccard similarity:

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



LSH for Jaccard similarity:

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

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

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

NEAR NEIGHBOR SEARCH

Basic approach for near neighbor search in a database.

Pre-processing:

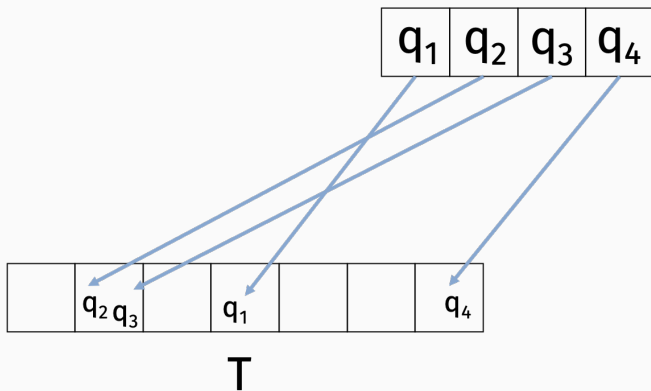
- Select random LSH function $h : \{0, 1\}^d \rightarrow 1, \dots, m$.
- Create table T with $m = O(n)$ slots.¹
- For $i = 1, \dots, n$, insert \mathbf{q}_i into $T(h(\mathbf{q}_i))$.

Query:

- Want to find near neighbors of input $\mathbf{y} \in \{0, 1\}^d$.
- Linear scan through all vectors $\mathbf{q} \in T(h(\mathbf{y}))$ and return any that are close to \mathbf{y} . Time required is $O(d \cdot |T(h(\mathbf{y}))|)$.

¹Enough to make the $O(1/m)$ term negligible.

NEAR NEIGHBOR SEARCH



Two main considerations:

- **False Negative Rate:** What's the probability we do not find a vector that is close to \mathbf{y} ?
- **False Positive Rate:** What's the probability that a vector in $T(h(\mathbf{y}))$ is not close to \mathbf{y} ?

A higher false negative rate means we miss near neighbors.

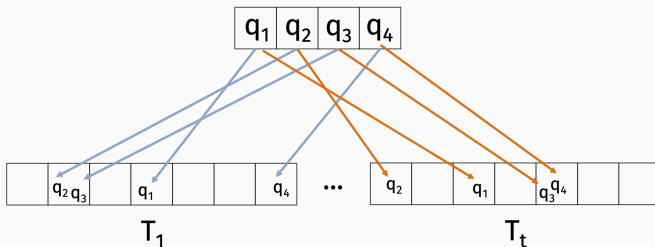
A higher false positive rate means increased runtime – we need to compute $J(\mathbf{q}, \mathbf{y})$ for every $\mathbf{q} \in T(h(\mathbf{y}))$ to check if it's actually close to \mathbf{y} .

Note: The meaning of “close” and “not close” is application dependent. E.g. we might specify that we want to find anything with Jaccard similarity $> .4$, but not with Jaccard similarity $< .2$.

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

What's the probability we do not find \mathbf{q} ?

REDUCING FALSE NEGATIVE RATE



Pre-processing:

- Select t independent LSH's $h_1, \dots, h_t : \{0, 1\}^d \rightarrow 1, \dots, m$.
- Create tables T_1, \dots, T_t , each with m slots.
- For $i = 1, \dots, n, j = 1, \dots, t$,
 - Insert q_i into $T_j(h_j(q_i))$.

Query:

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

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

What's the probability we find \mathbf{q} ?

WHAT HAPPENS TO FALSE POSITIVES?

Suppose there is some other database point \mathbf{z} with $J(\mathbf{y}, \mathbf{z}) = .2$.

What is the probability we will need to compute $J(\mathbf{z}, \mathbf{y})$ in our hashing scheme with one table? I.e. the probability that \mathbf{y} hashes into at least one bucket containing \mathbf{z} .

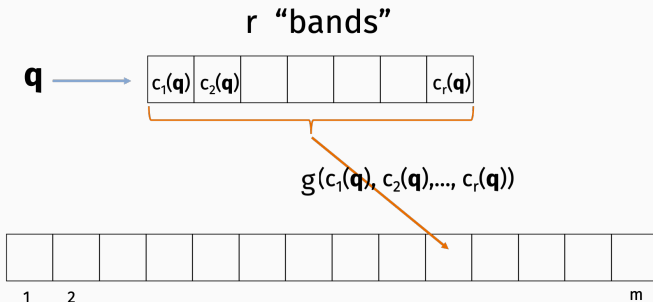
In the new scheme with $t = 10$ tables?

REDUCING FALSE POSITIVES

Change our locality sensitive hash function.

Tunable LSH for Jaccard similarity:

- Choose parameter $r \in \mathbb{Z}^+$.
- Let $c_1, \dots, c_r : \{0, 1\}^d \rightarrow [0, 1]$ be random MinHash.
- Let $g : [0, 1]^r \rightarrow \{1, \dots, m\}$ be a uniform random hash function.
- Let $h(x) = g(c_1(x), \dots, c_r(x))$.

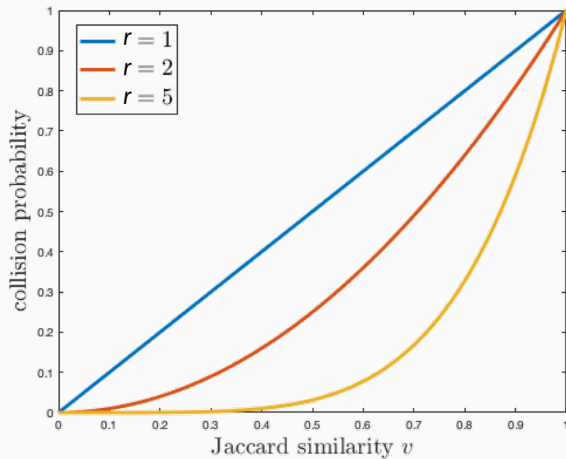


REDUCING FALSE POSITIVES

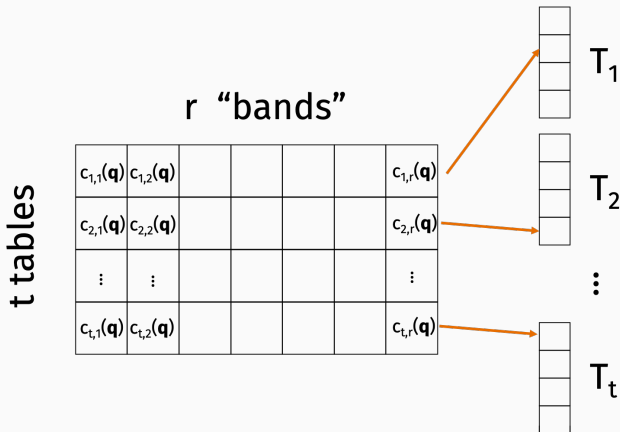
Tunable LSH for Jaccard similarity:

- Choose parameter $r \in \mathbb{Z}^+$.
- Let $c_1, \dots, c_r : \{0, 1\}^d \rightarrow [0, 1]$ be random MinHash.
- Let $g : [0, 1]^r \rightarrow \{1, \dots, m\}$ be a uniform random hash function.
- Let $h(\mathbf{x}) = g(c_1(\mathbf{x}), \dots, c_r(\mathbf{x}))$.

If $J(\mathbf{q}, \mathbf{y}) = v$, then $\Pr[h(\mathbf{q}) = h(\mathbf{y})] =$



Full LSH scheme has two parameters to tune:



Effect of **increasing number of tables t** on:

False Negatives

False Positives

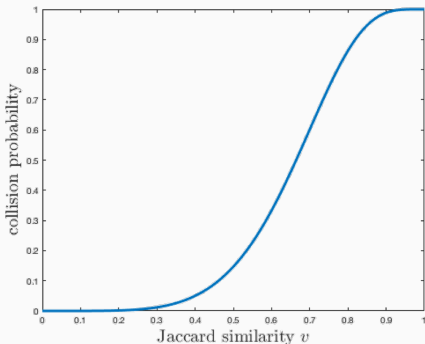
Effect of **increasing number of bands r** on:

False Negatives

False Positives

S-CURVE TUNING

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

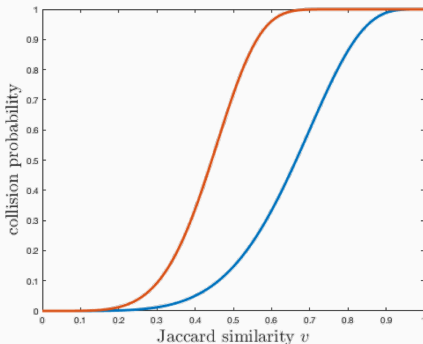


$$r = 5, t = 5$$

S-CURVE TUNING

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

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

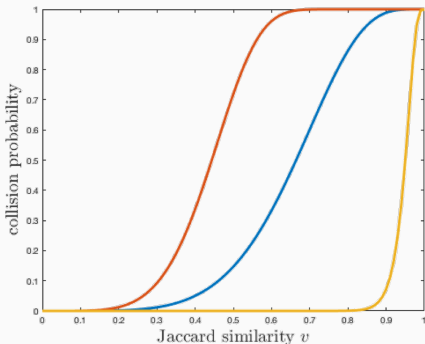


$$r = 5, t = 40$$

S-CURVE TUNING

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

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

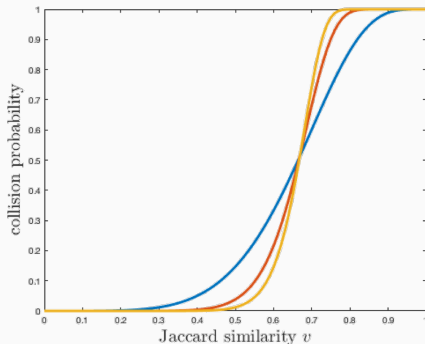


$$r = 40, t = 5$$

S-CURVE TUNING

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

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



Increasing both r and t gives a steeper curve.

Better for search, but worse space complexity.

FIXED THRESHOLD

Use Case 1: Fixed threshold.

- Shazam wants to find match to audio clip \mathbf{y} in a database of 10 million clips.
- There are 10 true matches with $J(\mathbf{y}, \mathbf{q}) > .9$.
- There are 10,000 near matches with $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$.
- All other items have $J(\mathbf{y}, \mathbf{q}) < .7$.

With $r = 25$ 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 $\lesssim 1 - (1 - .7^{25})^{40} = .005$

Upper bound on total number of items checked:

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

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

Directly trade space for fast search.

Near Neighbor Search Problem

Concrete worst case result:

Theorem (Indyk, Motwani, 1998)

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

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

$\|\mathbf{q} - \mathbf{y}\|_0$ = “hamming distance” = number of elements that differ between \mathbf{q} and \mathbf{y} .

Theorem (Indyk, Motwani, 1998)

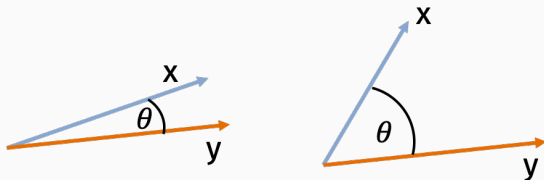
Let q be the closest database vector to \mathbf{y} . Return a vector $\tilde{\mathbf{q}}$ with $\|\tilde{\mathbf{q}} - \mathbf{y}\|_0 \leq C \cdot \|\mathbf{q} - \mathbf{y}\|_0$ in:

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

OTHER LSH FUNCTIONS

Good locality sensitive hash functions exist for other similarity measures.

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



$$-1 \leq \cos(\theta(x, y)) \leq 1.$$

Cosine similarity is natural “inverse” for Euclidean distance.

Euclidean distance $\|\mathbf{x} - \mathbf{y}\|_2^2$:

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

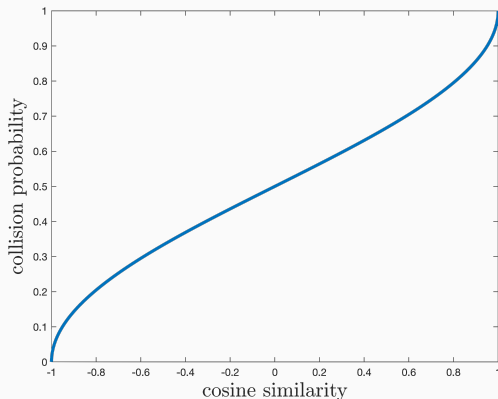
Locality sensitive hash for **cosine similarity**:

- Let $\mathbf{g} \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $f: \{-1, 1\} \rightarrow \{1, \dots, m\}$ be a uniformly random hash function.
- $h: \mathbb{R}^d \rightarrow \{1, \dots, m\}$ is defined $h(\mathbf{x}) = f(\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle))$.

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

Theorem (to be prove): If $\cos(\theta(x, y)) = v$, then

$$\Pr[h(x) == h(y)] = 1 - \frac{\theta}{\pi} + \frac{1}{m} = 1 - \frac{\cos^{-1}(v)}{\pi} + \frac{1}{m}$$



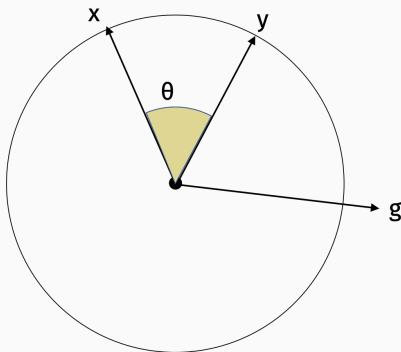
SimHash can be tuned, just like our MinHash based LSH function for Jaccard similarity:

- Let $\mathbf{g}_1, \dots, \mathbf{g}_r \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $f: \{-1, 1\}^r \rightarrow \{1, \dots, m\}$ be a uniformly random hash function.
- $h: \mathbb{R}^d \rightarrow \{1, \dots, m\}$ is defined
 $h(\mathbf{x}) = f([\text{sign}(\langle \mathbf{g}_1, \mathbf{x} \rangle), \dots, \text{sign}(\langle \mathbf{g}_r, \mathbf{x} \rangle)])$.

$$\Pr[h(\mathbf{x}) == h(\mathbf{y})] = \left(1 - \frac{\theta}{n}\right)^r$$

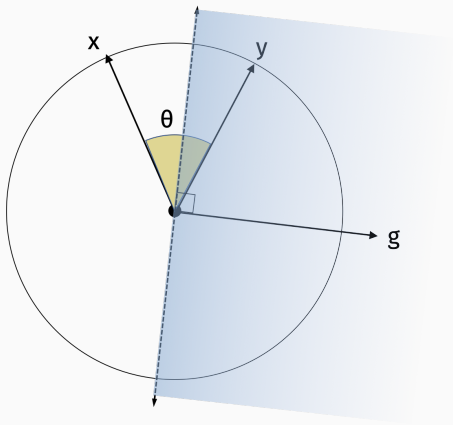
SIMHASH ANALYSIS IN 2D

To prove: $\Pr[h(x) == h(y)] = 1 - \frac{\theta}{\pi}$, where $h(x) = f(\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle))$ and f is uniformly random hash function.



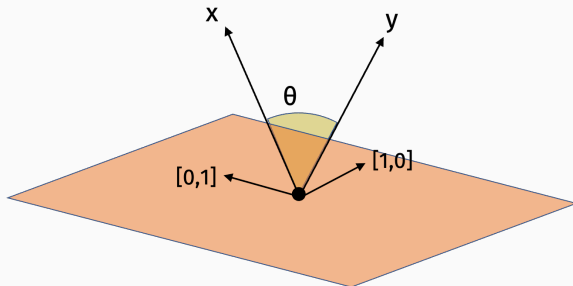
$$\Pr[h(x) == h(y)] = z + \frac{1 - v}{m} \approx z.$$

where $z = \Pr[\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle) == \text{sign}(\langle \mathbf{g}, \mathbf{y} \rangle)]$

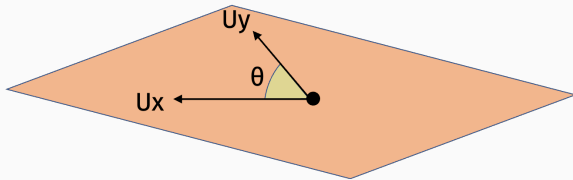


$\Pr[h(\mathbf{x}) == h(\mathbf{y})] \approx$ probability \mathbf{x} and \mathbf{y} are on the same side of hyperplane orthogonal to \mathbf{g} .

SIMHASH ANALYSIS HIGHER DIMENSIONS



There is always some rotation matrix \mathbf{U} such that $\mathbf{U}\mathbf{x}$, $\mathbf{U}\mathbf{y}$ are spanned by the first two-standard basis vectors and have the same cosine similarity as \mathbf{x} and \mathbf{y} .



There is always some rotation matrix \mathbf{U} such that \mathbf{x}, \mathbf{y} are spanned by the first two-standard basis vectors.

Claim:

$$\begin{aligned} 1 - \frac{\theta}{\pi} &= \Pr[\text{sign}(\langle \mathbf{g}, \mathbf{Ux} \rangle) == \text{sign}(\langle \mathbf{g}, \mathbf{Uy} \rangle)] \\ &= \Pr[\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle) == \text{sign}(\langle \mathbf{g}, \mathbf{y} \rangle)] \end{aligned}$$

Why?

BREAK

Have some function $f: \mathbb{R}^d \rightarrow \mathbb{R}$. Want to find \mathbf{x}^* such that:

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}).$$

Or at least $\hat{\mathbf{x}}$ which is close to a minimum. E.g.

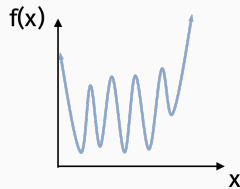
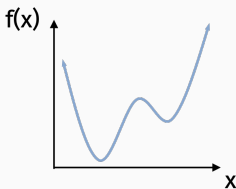
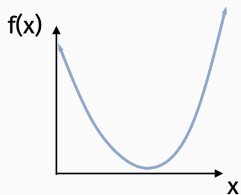
$$f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x}) + \epsilon$$

Often we have some additional constraints:

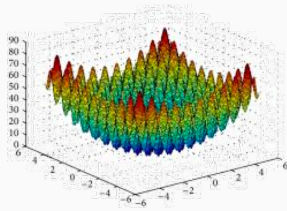
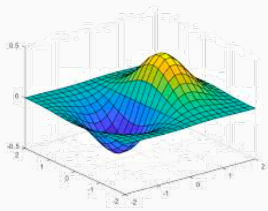
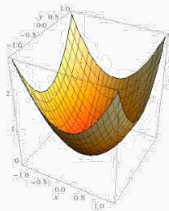
- $\mathbf{x} > 0$.
- $\|\mathbf{x}\|_2 \leq R, \|\mathbf{x}\|_1 \leq R$.
- $\mathbf{a}^T \mathbf{x} > c$.

CONTINUOUS OPTIMIZATION

Dimension $d = 1$:



Dimension $d = 2$:



Continuous optimization is the foundation of modern machine learning.

Supervised learning: Want to learn a model that maps inputs

- numerical data vectors
- images, video
- text documents

to predictions

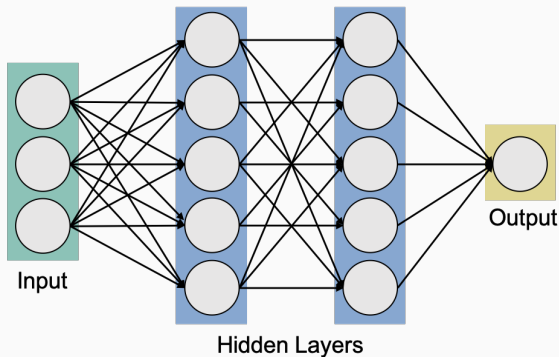
- numerical value (probability stock price increases)
- label (is the image a cat? does the image contain a car?)
- decision (turn car left, rotate robotic arm)

Let $M_{\mathbf{x}}$ be a model with parameters $\mathbf{x} = \{x_1, \dots, x_k\}$, which takes as input a data vector \mathbf{a} and outputs a prediction.

Example:

$$M_{\mathbf{x}}(\mathbf{a}) = \text{sign}(\mathbf{a}^T \mathbf{x})$$

Example:



$\mathbf{x} \in \mathbb{R}^{(\# \text{ of connections})}$ is the parameter vector containing all the network weights.

Classic approach in supervised learning: Find a model that works well on data that you already have the answer for (labels, values, classes, etc.).

- Model $M_{\mathbf{x}}$ parameterized by a vector of numbers \mathbf{x} .
- Dataset $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(n)}$ with outputs $y^{(1)}, \dots, y^{(n)}$.

Want to find $\hat{\mathbf{x}}$ so that $M_{\hat{\mathbf{x}}}(\mathbf{a}^{(i)}) \approx y^{(i)}$ for $i \in 1, \dots, n$.

How do we turn this into a function minimization problem?

Loss function $L(M_x(\mathbf{a}), y)$: Some measure of distance between prediction $M_x(\mathbf{a})$ and target output y . Increases if they are further apart.

- Squared (ℓ_2) loss: $|M_x(\mathbf{a}) - y|^2$
- Absolute deviation (ℓ_1) loss: $|M_x(\mathbf{a}) - y|$
- Hinge loss: $1 - y \cdot M_x(\mathbf{a})$
- Cross-entropy loss (log loss).
- Etc.

Empirical risk minimization:

$$f(\mathbf{x}) = \sum_{i=1}^n L\left(M_{\mathbf{x}}(\mathbf{a}^{(i)}), y^{(i)}\right)$$

Solve the optimization problem $\min_{\mathbf{x}} f(\mathbf{x})$.

EXAMPLE: LINEAR REGRESSION

- $M_{\mathbf{x}}(\mathbf{a}) = \mathbf{x}^T \mathbf{a}$. \mathbf{x} contains the regression coefficients.
- $L(z, y) = |z - y|^2$.
- $f(\mathbf{x}) = \sum_{i=1}^n |\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)}|^2$

$$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|_2^2$$

where \mathbf{A} is a matrix with $\mathbf{a}^{(i)}$ as its i^{th} row and \mathbf{y} is a vector with $y^{(i)}$ as its i^{th} entry.

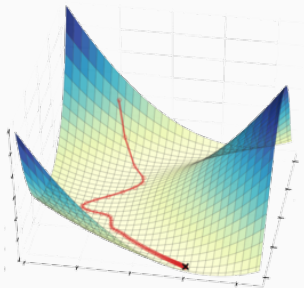
The choice of algorithm to minimize $f(\mathbf{x})$ will depend on:

- The form of $f(\mathbf{x})$ (is it linear, is it quadratic, does it have finite sum structure, etc.)
- If there are any additional constraints imposed on \mathbf{x} . E.g. $\|\mathbf{x}\|_2 \leq c$.

What are some example algorithms for continuous optimization?

FIRST TOPIC: GRADIENT DESCENT + VARIANTS

Gradient descent: A greedy algorithm for minimizing functions of multiple variables that often works amazingly well.



Runtime generally scales linearly with the dimension of x (although this is a bit of an over-simplification).

- Cutting plane methods (e.g. center-of-gravity, ellipsoid)
- Interior point methods

Fast and more accurate in low-dimensions, slower in very high dimensions. Generally runtime scales polynomially with the dimension of \mathbf{x} .

For $i = 1, \dots, d$, let x_i be the i^{th} entry of \mathbf{x} . Let $\mathbf{e}^{(i)}$ be the i^{th} standard basis vector.

Partial derivative:

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{e}^{(i)}) - f(\mathbf{x})}{t}$$

Directional derivative:

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}$$

Gradient:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \frac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_d}(\mathbf{x}) \end{bmatrix}$$

Directional derivative:

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^T \mathbf{v}.$$

Given a function f to minimize, assume we have:

- **Function oracle:** Evaluate $f(\mathbf{x})$ for any \mathbf{x} .
- **Gradient oracle:** Evaluate $\nabla f(\mathbf{x})$ for any \mathbf{x} .

We view the implementation of these oracles as black-boxes, but they can often require a fair bit of computation.

EXAMPLE GRADIENT EVALUATION

Linear least-squares regression:

- Given $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(n)} \in \mathbb{R}^d, y^{(1)}, \dots, y^{(n)} \in \mathbb{R}$.
- Want to minimize:

$$f(\mathbf{x}) = \sum_{i=1}^n \left(\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)} \right)^2 = \|\mathbf{Ax} - \mathbf{y}\|_2^2.$$

$$\frac{\partial f}{\partial x_j} = \sum_{i=1}^n 2 \left(\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)} \right) \cdot a_j^{(i)} = (2\mathbf{Ax} - \mathbf{y})^T \boldsymbol{\alpha}^{(j)}$$

where $\boldsymbol{\alpha}^{(j)}$ is the j^{th} column of \mathbf{A} .

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{y})$$

What is the time complexity of a gradient oracle for $\nabla f(\mathbf{x})$?

Greedy approach: Given a starting point \mathbf{x} , make a small adjustment that decreases $f(\mathbf{x})$. In particular, $\mathbf{x} \leftarrow \mathbf{x} + \eta \mathbf{v}$ and $f(\mathbf{x}) \leftarrow f(\mathbf{x} + \eta \mathbf{v})$.

What property do I want in \mathbf{v} ?

Leading question: When η is small, what's an approximation for $f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x})$?

$$f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x}) \approx$$

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^T \mathbf{v}.$$

So:

$$f(\mathbf{x} + \eta\mathbf{v}) - f(\mathbf{x}) \approx$$

How should we choose \mathbf{v} so that $f(\mathbf{x} + \eta\mathbf{v}) < f(\mathbf{x})$?

Prototype algorithm:

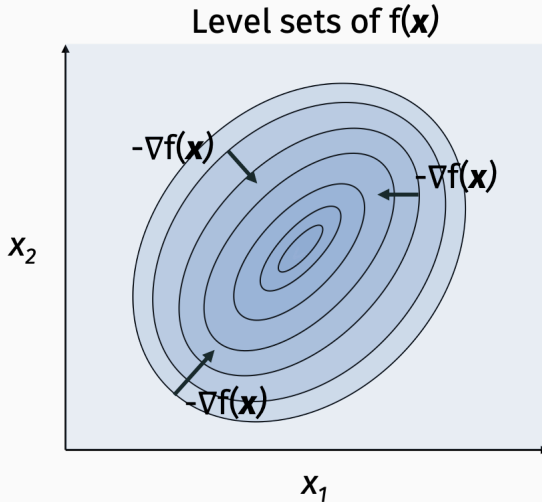
- Choose starting point $\mathbf{x}^{(0)}$.
- For $i = 0, \dots, T$:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\mathbf{x}^{(T)}$.

η is a step-size parameter, which is often adapted on the go.
For now, assume it is fixed ahead of time.

1 dimensional example:

GRADIENT DESCENT INTUITION

2 dimensional example:



KEY RESULTS

For a convex function $f(\mathbf{x})$: For sufficiently small η and a sufficiently large number of iterations T , gradient descent will converge to a **near global minimum**:

$$f(\mathbf{x}^{(T)}) \leq f(\mathbf{x}^*) + \epsilon.$$

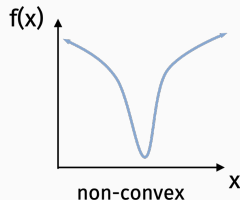
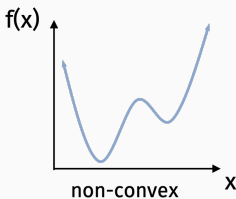
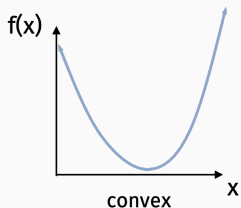
Examples: least squares regression, logistic regression, kernel regression, SVMs.

For a non-convex function $f(\mathbf{x})$: For sufficiently small η and a sufficiently large number of iterations T , gradient descent will converge to a **near stationary point**:

$$\|\nabla f(\mathbf{x}^{(T)})\|_2 \leq \epsilon.$$

Examples: neural networks, matrix completion problems, mixture models.

CONVEX VS. NON-CONVEX



One issue with non-convex functions is that they can have **local minima**. Even when they don't, convergence analysis requires different assumptions than convex functions.

APPROACH FOR THIS UNIT

We care about how fast gradient descent and related methods converge, not just that they do converge.

- Bounding iteration complexity requires placing some assumptions on $f(\mathbf{x})$.
- Stronger assumptions lead to better bounds on the convergence.

Understanding these assumptions can help us design faster variants of gradient descent (there are many!).