

## LSH for $\ell_2$ distances

Lecture 15

October 15, 2020

# LSH Approach for Approximate NNS

Use **locality-sensitive hashing** to solve simplified decision problem

## Definition

A family of hash functions is  $(r, cr, p_1, p_2)$ -LSH with  $p_1 > p_2$  and  $c > 1$  if  $h$  drawn randomly from the family satisfies the following:

- $\Pr[h(x) = h(y)] \geq p_1$  when  $\text{dist}(x, y) \leq r$
- $\Pr[h(x) = h(y)] \leq p_2$  when  $\text{dist}(x, y) \geq cr$

**Key parameter:** the gap between  $p_1$  and  $p_2$  measured as  $\rho = \frac{\log p_1}{\log p_2}$  usually small.

Two-level hashing scheme:

- Amplify basic locality sensitive hash family to create better family by repetition
- Use several copies of amplified hash functions

1st level binary search based on  $r$  on top of above scheme

# LSH Approach for Approximate NNS

**Key parameter:** the gap between  $p_1$  and  $p_2$  measured as  $\rho = \frac{\log p_1}{\log p_2}$  usually small.

- $L \simeq n^\rho$  hash tables
- Storage:  $n^{1+\rho}$  (ignoring log factors)
- Query time:  $kn^\rho$  (ignoring log factors) where  $k = \log_{1/p_2} n$

# LSH for Euclidean Distances

Now  $x_1, x_2, \dots, x_n \in \mathbb{R}^d$  and  $\text{dist}(x, y) = \|x - y\|_2$

First do dimensionality reduction (JL) to reduce  $d$  (if necessary) to  $O(\log n)$  (since we are using  $c$ -approximation anyway)

# LSH for Euclidean Distances

Now  $x_1, x_2, \dots, x_n \in \mathbb{R}^d$  and  $\text{dist}(x, y) = \|x - y\|_2$

First do dimensionality reduction (JL) to reduce  $d$  (if necessary) to  $O(\log n)$  (since we are using  $c$ -approximation anyway)

What is a good basic locality-sensitive hashing scheme? That is, we want a hashing approach that makes nearby points more likely to collide than farther away points.

# LSH for Euclidean Distances

Now  $x_1, x_2, \dots, x_n \in \mathbb{R}^d$  and  $\text{dist}(x, y) = \|x - y\|_2$

First do dimensionality reduction (JL) to reduce  $d$  (if necessary) to  $O(\log n)$  (since we are using  $c$ -approximation anyway)

What is a good basic locality-sensitive hashing scheme? That is, we want a hashing approach that makes nearby points more likely to collide than farther away points.

Projections onto random lines plus bucketing

# Random unit vector

**Question:** How do we generate a random unit vector in  $\mathbb{R}^d$  (same as a uniform point on the sphere  $S^{n-1}$ )?

# Random unit vector

**Question:** How do we generate a random unit vector in  $\mathbb{R}^d$  (same as a uniform point on the sphere  $S^{n-1}$ )?

- Pick  $d$  independent rvs  $Z_1, Z_2, \dots, Z_d$  where each  $Z_i \simeq \mathcal{N}(\mathbf{0}, \mathbf{1})$  and let  $\mathbf{g} = (Z_1, Z_2, \dots, Z_d)$  (also called a random Gaussian vector)
- $\mathbf{g}$  is symmetric and hence is a random direction
- to obtain random unit vector normalize  $\mathbf{g}' = \mathbf{g} / \|\mathbf{g}\|_2$
- When  $d$  is large  $\|\mathbf{g}\|_2^2 = \sum_i Z_i^2$  is concentrated around  $d$  and hence  $\|\mathbf{g}\|_2 = (1 \pm \epsilon)\sqrt{d}$  with high probability.
- Thus  $\mathbf{g} / \sqrt{d}$  is a proxy for random unit vector and is easier to work with in many cases

# Projection onto a random gaussian vector

## Lemma

Suppose  $\mathbf{x} \in \mathbb{R}^d$  and  $\mathbf{g}$  is a random Gaussian vector. Let  $Y = \mathbf{x} \cdot \mathbf{g}$ . Then  $Y \sim \mathcal{N}(0, \|\mathbf{x}\|_2)$  and hence  $E[Y^2] = (\|\mathbf{x}\|_2)^2$ .

# Hashing scheme

- Pick a random unit Gaussian vector  $u$
- Pick a random shift  $a \in (0, r]$
- For vector  $x$  set  $h_{u,a} = \lfloor \frac{x \cdot u + a}{r} \rfloor$

# Analysis

Suppose  $x, y$  are such that  $\|x - y\|_2 \leq r$ . What is  
 $p_1 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Suppose  $x, y$  are such that  $\|x - y\|_2 \geq cr$ . What is  
 $p_2 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

# Analysis

Suppose  $x, y$  are such that  $\|x - y\|_2 \leq r$ . What is  
 $p_1 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Suppose  $x, y$  are such that  $\|x - y\|_2 \geq cr$ . What is  
 $p_2 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Let  $q = x - y$ . Let  $s = \|q\|_2$  be length of  $q$ .  
From Lemma  $q \cdot g$  is distributed as  $s\mathcal{N}(0, 1)$ .

# Analysis

Suppose  $x, y$  are such that  $\|x - y\|_2 \leq r$ . What is  $p_1 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Suppose  $x, y$  are such that  $\|x - y\|_2 \geq cr$ . What is  $p_2 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Let  $q = x - y$ . Let  $s = \|q\|_2$  be length of  $q$ . From Lemma  $q \cdot g$  is distributed as  $s\mathcal{N}(0, 1)$ .

## Observations:

- $h(x) \neq h(y)$  if  $|q \cdot g| \geq r$
- If  $|q \cdot g| < r$  then  $h(x) = h(y)$  with probability  $1 - |q \cdot g|/r$

# Analysis

Suppose  $x, y$  are such that  $\|x - y\|_2 \leq r$ . What is  
 $p_1 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Suppose  $x, y$  are such that  $\|x - y\|_2 \geq cr$ . What is  
 $p_2 = \Pr[h_{u,a}(x) = h_{u,a}(y)]$

Let  $q = x - y$ . Let  $s = \|q\|_2$  be length of  $q$ .  
From Lemma  $q \cdot g$  is distributed as  $s\mathcal{N}(0, 1)$ .

## Observations:

- $h(x) \neq h(y)$  if  $|q \cdot g| \geq r$
- If  $|q \cdot g| < r$  then  $h(x) = h(y)$  with probability  $1 - |q \cdot g|/r$

Thus collision probability depends only on  $s$

# Analysis

Let  $\mathbf{q} = \mathbf{x} - \mathbf{y}$ . Let  $s = \|\mathbf{q}\|_2$  be length of  $\mathbf{q}$ .  
From Lemma  $\mathbf{q} \cdot \mathbf{g}$  is distributed as  $s\mathcal{N}(\mathbf{0}, 1)$ .

Observations:

- $h(\mathbf{x}) \neq h(\mathbf{y})$  if  $|\mathbf{q} \cdot \mathbf{g}| \geq r$
- If  $|\mathbf{q} \cdot \mathbf{g}| < r$  then  $h(\mathbf{x}) = h(\mathbf{y})$  with probability  $1 - |\mathbf{q} \cdot \mathbf{g}|/r$

For a fixed  $s$  collision probability is

$$p(s) = \int_0^r f(t)(1 - t/r)dt$$

where  $f$  is the density function of  $|s\mathcal{N}(\mathbf{0}, 1)|$ .

Rewriting

$$p(s) = \int_0^r \frac{1}{s} f\left(\frac{t}{s}\right)(1 - t/r)dt$$

where  $f$  is the density function of the  $|\mathcal{N}(\mathbf{0}, 1)|$ .

# Analysis

$$p(s) = \int_0^r \frac{1}{s} f\left(\frac{t}{s}\right) (1 - t/r) dt$$

where  $f$  is the density function of the  $|\mathcal{N}(\mathbf{0}, \mathbf{1})|$ .

Recall  $p_1 = p(r)$  and  $p_2 = p(cr)$  and we are interested in

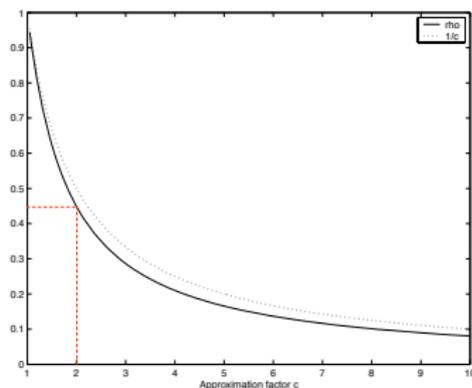
$$\rho = \frac{\log p_1}{\log p_2}.$$

# Analysis

$$\rho(s) = \int_0^r \frac{1}{s} f\left(\frac{t}{s}\right) (1 - t/r) dt$$

where  $f$  is the density function of the  $|\mathcal{N}(\mathbf{0}, \mathbf{1})|$ .

Recall  $p_1 = p(r)$  and  $p_2 = p(cr)$  and we are interested in  $\rho = \frac{\log p_1}{\log p_2}$ . Show  $\rho < 1/c$  by plot



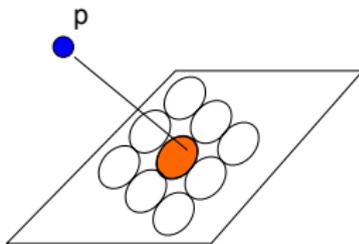
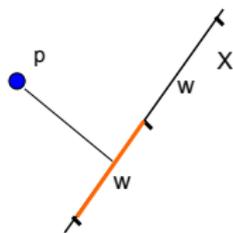
# NNS for Euclidean distances

- For any fixed  $c > 1$  use above scheme to obtain
  - Storage:  $O(n^{1+1/c} \text{polylog}(n))$
  - Query time:  $O(dn^{1/c} \text{polylog}(n))$
- Can use JL to reduce  $d$  to  $O(\log n)$ .

# Improved LSH Scheme

[Andoni-Indyk'06]

- Basic LSH scheme projects points into lines
- Better scheme: pick some small constant  $t$  and project points into  $R^t$
- Use lattice based space partitioning scheme to “bucket” instead of intervals



Figures from Piotr Indyk's slides

# Improved LSH Scheme

[Andoni-Indyk'06]

- Basic LSH scheme projects points into lines
- Better scheme: pick some small constant  $t$  and project points into  $R^t$
- Use lattice based space partitioning scheme to “bucket” instead of intervals
- Leads to  $\rho \simeq 1/c^2 + O(\log t/\sqrt{t})$  and hence tends to  $1/c^2$  for large  $t$  and fixed  $c$
- Lower bound for LSH in  $\ell_2$  says  $\rho \geq 1/c^2$

# Data dependent LSH Scheme

LSH is data oblivious. That is, the hash families are chosen before seeing the data. Can one do better by choosing hash functions based on the given set of points?

# Data dependent LSH Scheme

LSH is data oblivious. That is, the hash families are chosen before seeing the data. Can one do better by choosing hash functions based on the given set of points?

Yes.

[Andoni-Indyk-Nguyen-Razenshteyn'14, Andoni-Razenshteyn'15]

- $\rho = 1/(2c^2 - 1)$  for  $\ell_2$  improving upon  $1/c^2$  for data oblivious LSH (which is tight in worst case)
- $\rho = 1/(c^2 - 1)$  for  $\ell_1$ /Hamming cube improving upon  $1/c$  for data oblivious LSH

# LSH Summary

- A modular hashing based scheme for similarity estimation
- Main competitors are space partitioning data structures such as variants of k-d trees
- Provides speedups but uses more memory
- Does not appear to be a clear winner

# Digression: $p$ -stable distributions

For  $F_2$  estimation and JL and LSH we used important “stability” property of the Normal distribution.

## Lemma

Let  $Y_1, Y_2, \dots, Y_d$  be independent random variables with distribution  $\mathcal{N}(0, 1)$ .  $Z = \sum_i x_i Y_i$  has distribution  $\|x\|_2 \mathcal{N}(0, 1)$

Standard Gaussian is **2**-stable.

# Digression: $p$ -stable distributions

For  $F_2$  estimation and JL and LSH we used important “stability” property of the Normal distribution.

## Lemma

Let  $Y_1, Y_2, \dots, Y_d$  be independent random variables with distribution  $\mathcal{N}(0, 1)$ .  $Z = \sum_i x_i Y_i$  has distribution  $\|x\|_2 \mathcal{N}(0, 1)$

Standard Gaussian is 2-stable.

## Definition

A distribution  $\mathcal{D}$  is  $p$ -stable if  $Z = \sum_i x_i Y_i$  has distribution  $\|x\|_p \mathcal{D}$  when the  $Y_i$  are independent and each of them is distributed as  $\mathcal{D}$ .

# Digression: $p$ -stable distributions

For  $F_2$  estimation and JL and LSH we used important “stability” property of the Normal distribution.

## Lemma

Let  $Y_1, Y_2, \dots, Y_d$  be independent random variables with distribution  $\mathcal{N}(0, 1)$ .  $Z = \sum_i x_i Y_i$  has distribution  $\|x\|_2 \mathcal{N}(0, 1)$

Standard Gaussian is 2-stable.

## Definition

A distribution  $\mathcal{D}$  is  $p$ -stable if  $Z = \sum_i x_i Y_i$  has distribution  $\|x\|_p \mathcal{D}$  when the  $Y_i$  are independent and each of them is distributed as  $\mathcal{D}$ .

**Question:** Do  $p$ -stable distributions exist for  $p \neq 2$ ?

# $p$ -stable distributions

**Fact:**  $p$ -stable distributions exist for all  $p \in (0, 2]$  and do not exist for  $p > 2$ .

$p = 1$  is the Cauchy distribution which is the distribution of the ratio of two independent Gaussian random variables. Has a closed form density function  $\frac{1}{\pi(1+x^2)}$ . Mean and variance are not finite.

# $p$ -stable distributions

**Fact:**  $p$ -stable distributions exist for all  $p \in (0, 2]$  and do not exist for  $p > 2$ .

$p = 1$  is the Cauchy distribution which is the distribution of the ratio of two independent Gaussian random variables. Has a closed form density function  $\frac{1}{\pi(1+x^2)}$ . Mean and variance are not finite.

For general  $p$  no closed form formula for density but can sample from the distribution.

# $p$ -stable distributions

**Fact:**  $p$ -stable distributions exist for all  $p \in (0, 2]$  and do not exist for  $p > 2$ .

$p = 1$  is the Cauchy distribution which is the distribution of the ratio of two independent Gaussian random variables. Has a closed form density function  $\frac{1}{\pi(1+x^2)}$ . Mean and variance are not finite.

For general  $p$  no closed form formula for density but can sample from the distribution.

Streaming, sketching, LSH ideas for  $\ell_2$  generalize to  $\ell_p$  for  $p \in (0, 2]$  via  $p$ -stable distributions and additional technical work.

# Digression: Doubling dimension

**Thesis/assumption:** Real world data is high-dimensional in explicit representation but low-dimensional in "content".

Several interpretations of what it means for data to be low-dimensional

- Data lies in a low-dimensional manifold
- Data can be projected into low dimensions while preserving certain properties (JL for instance)
- Data has a latent low-dimensional description (SVD, PCA, tensor decomposition, etc)
- Data has low doubling dimension
- ...

# Intrinsic dimension

Let  $(V, \text{dist})$  be a finite metric space.

- $\text{dist}(x, y) = \text{dist}(y, x)$  for all  $x, y \in V$  (symmetry)
- $\text{dist}(x, x) = 0$  for all  $x \in V$  (reflexivity)
- $\text{dist}(x, y) + \text{dist}(y, z) \geq \text{dist}(x, z)$  for all  $x, y, z \in V$  (triangle inequality)

**Question:** Can we quantify whether  $(V, \text{dist})$  behaves like a low-dimensional Euclidean space? Does this have any benefits?

# Doubling dimension

Property of  $\mathbb{R}^d$ : A ball of radius  $r$  can be covered by  $c^d$  balls of radius  $r/2$  for some constant  $c \leq 4$ .

# Doubling dimension

Property of  $\mathbb{R}^d$ : A ball of radius  $r$  can be covered by  $c^d$  balls of radius  $r/2$  for some constant  $c \leq 4$ .

Given  $(V, d)$  let  $B(p, r)$  be the ball of radius  $r$  around  $p$  and view it as a set of points:

$$B(p, r) = \{q \mid \text{dist}(p, q) \leq r\}$$

# Doubling dimension

Property of  $\mathbb{R}^d$ : A ball of radius  $r$  can be covered by  $c^d$  balls of radius  $r/2$  for some constant  $c \leq 4$ .

Given  $(V, d)$  let  $B(p, r)$  be the ball of radius  $r$  around  $p$  and view it as a set of points:

$$B(p, r) = \{q \mid \text{dist}(p, q) \leq r\}$$

## Definition

A finite metric space  $(V, \text{dist})$  has doubling dimension  $d$  if for all  $p \in V$  and all  $r > 0$ ,  $B(p, r)$  can be covered by  $2^d$  balls of radius at most  $r/2$ .

# Doubling dimensions

## Definition

A finite metric space  $(V, \text{dist})$  has doubling dimension  $d$  if for all  $p \in V$  and all  $r > 0$ ,  $B(p, r)$  can be covered by  $2^d$  balls of radius at most  $r/2$ .

Many algorithms/data structures for  $\mathbb{R}^d$  can be extended to metric spaces with doubling dimension  $d$  with comparable running times.

Including approximate NNS.

See [Clarkson, Krauthgamer-Lee, HarPeled-Mendel]