Lecture 9: Dimension Reduction

Rafael Oliveira

University of Waterloo Cheriton School of Computer Science rafael.oliveira.teaching@gmail.com

October 7, 2020

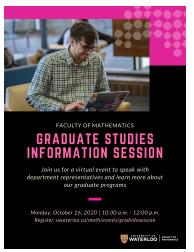
Overview

- Introduction
 - Administrivia
 - Why Reduce Dimensions?
 - Background: Continuous Probability Distributions
- Main Problem
 - Johnson-Lindenstrauss Lemma
- Acknowledgements

Grad School at UW!

Link to register:

https://uwaterloo.ca/math/events/gradinfosession



Why do we want low-dimensional objects?

When dealing with high-dimensional data, often times want to reduce dimension so that our algorithms run faster
In smaller dimension, things generally run faster, need less storage space,
easier to communicate.

Why do we want low-dimensional objects?

When dealing with high-dimensional data, often times want to reduce dimension so that our algorithms run faster In *smaller dimension*, things generally *run faster*, need *less storage space*,

easier to communicate.

- Nearest Neighbor Search
- Large Scale Regression Problems
- Minimum Enclosing Ball
- Numerical linear algebra on large matrices
- Clustering

What do we want to preserve?

distances between points

- distances between points
- angles between vectors

- distances between points
- angles between vectors
- volumes of subsets of the input

- distances between points
- angles between vectors
- volumes of subsets of the input
- optimal solutions to optimization problems

What do we want to preserve?

- distances between points
- angles between vectors
- volumes of subsets of the input
- optimal solutions to optimization problems

To preserve *distances*, need to allow some *distortion* (approximate guarantees).

What do we want to preserve?

- distances between points
- angles between vectors
- volumes of subsets of the input
- optimal solutions to optimization problems

To preserve *distances*, need to allow some *distortion* (approximate guarantees).

• Cannot compress simplex while preserving all distances.

So far we have only dealt with *discrete random variables*. Today, we will use *continuous random variables*.

So far we have only dealt with *discrete random variables*. Today, we will use *continuous random variables*.

How can we define random variables/probabilities over a *continuous* (infinite) set?

So far we have only dealt with *discrete random variables*. Today, we will use *continuous random variables*.

How can we define random variables/probabilities over a *continuous* (infinite) set?

Say we have a real-valued random variable - that is, X takes values in \mathbb{R} .

Definition (Probability Density Function)

A probability density function $f: \mathbb{R} \to \mathbb{R}_{\geq 0}$ is a function such that

- ullet f is integrable over ${\mathbb R}$
- $\int_{-\infty}^{\infty} f(x) dx = 1$

So far we have only dealt with *discrete random variables*. Today, we will use *continuous random variables*.

How can we define random variables/probabilities over a *continuous* (infinite) set?

Say we have a real-valued random variable - that is, X takes values in \mathbb{R} .

Definition (Probability Density Function)

A probability density function $f: \mathbb{R} \to \mathbb{R}_{\geq 0}$ is a function such that

- ullet f is integrable over ${\mathbb R}$
- Probability density function f(x) <u>intuitively</u> gives us <u>relative likelihood</u> that X = x.

So far we have only dealt with *discrete random variables*. Today, we will use *continuous random variables*.

How can we define random variables/probabilities over a *continuous* (infinite) set?

Say we have a real-valued random variable - that is, X takes values in \mathbb{R} .

Definition (Probability Density Function)

A probability density function $f: \mathbb{R} \to \mathbb{R}_{\geq 0}$ is a function such that

- ullet f is integrable over ${\mathbb R}$
- $\bullet \int_{-\infty}^{\infty} f(x) dx = 1$
- Probability density function f(x) <u>intuitively</u> gives us <u>relative likelihood</u> that X = x.
- •

$$\Pr[a \le X \le b] = \int_a^b f(x) dx$$

Gaussian Random Variables (Normal Random Variables)

Definition

A real-valued random variable X has the *normal distribution* with

- mean μ
- variance σ^2 ,

denoted $X \sim \mathcal{N}(\mu, \sigma^2)$, if the probability density function of X, denoted $f_X : \mathbb{R} \to \mathbb{R}_{>0}$ is:

$$f_X(x) = \frac{1}{\sigma \cdot \sqrt{2\pi}} \cdot \exp\left(-\frac{1}{2} \cdot \left(\frac{x-\mu}{\sigma}\right)^2\right)$$

Gaussian Random Variables (Normal Random Variables)

Definition

A real-valued random variable X has the *normal distribution* with

- mean μ
- variance σ^2 .

denoted $X \sim \mathcal{N}(\mu, \sigma^2)$, if the probability density function of X, denoted $f_X : \mathbb{R} \to \mathbb{R}_{>0}$ is:

$$f_X(x) = \frac{1}{\sigma \cdot \sqrt{2\pi}} \cdot \exp\left(-\frac{1}{2} \cdot \left(\frac{x-\mu}{\sigma}\right)^2\right)$$

Remark

When $\mu = 0$ and $\sigma = 1$ we say that X has standard normal distribution.

Properties of Gaussians

Proposition (Sums of Gaussians)

If
$$X \sim \mathcal{N}(\mu_X, \sigma_X^2)$$
 and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ are independent Gaussians, then

$$X + Y \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2).$$

Properties of Gaussians

Proposition (Sums of Gaussians)

If
$$X \sim \mathcal{N}(\mu_X, \sigma_X^2)$$
 and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ are independent Gaussians, then

$$X + Y \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2).$$

Proposition (Multiplication by scalar)

If $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$, then

$$\sigma \cdot X \sim \mathcal{N}(\sigma \mu_X, (\sigma \cdot \sigma_X)^2).$$

Properties of Gaussians

Proposition (Sums of Gaussians)

If $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ are independent Gaussians, then

$$X + Y \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2).$$

Proposition (Multiplication by scalar)

If $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$, then

$$\sigma \cdot X \sim \mathcal{N}(\sigma \mu_X, (\sigma \cdot \sigma_X)^2).$$

Proposition (General Linear Combinations)

If $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ are independent random Gaussians, then

$$\sum_{i=1}^{n} \alpha_i \cdot X_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \alpha_i \cdot \mu_i , \sum_{i=1}^{n} (\alpha_i \cdot \sigma_i)^2\right).$$

χ^2 Random Variables

Definition

A real-valued random variable X has the χ^2 distribution with k degrees of freedom, denoted $X \sim \chi^2(k)$, if

$$X = Z_1^2 + \ldots + Z_k^2$$

where each $Z_i \sim \mathcal{N}(0,1)$ is an independent standard normal random variable.

Lemma (Chernoff for $\chi^2(k)$)

If $Y = \sum_{i=1}^{k} X_i^2$ is a $\chi^2(k)$ random variable with k degrees of freedom (recall $X_i \sim \mathcal{N}(0,1)$), then

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \exp\left(-\frac{3}{4} \cdot d\varepsilon^2\right)$$

Lemma (Chernoff for $\chi^2(k)$)

If $Y = \sum_{i=1}^{k} X_i^2$ is a $\chi^2(k)$ random variable with k degrees of freedom (recall $X_i \sim \mathcal{N}(0,1)$), then

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \exp\left(-\frac{3}{4} \cdot d\varepsilon^2\right)$$

• Let $t \in (0, 1/2)$ be a parameter

Lemma (Chernoff for $\chi^2(k)$)

If $Y = \sum_{i=1}^{k} X_i^2$ is a $\chi^2(k)$ random variable with k degrees of freedom (recall $X_i \sim \mathcal{N}(0,1)$), then

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \exp\left(-\frac{3}{4} \cdot d\varepsilon^2\right)$$

- Let $t \in (0, 1/2)$ be a parameter
- •

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] = \Pr\left[e^{tY} > e^{t \cdot (1+\varepsilon)^2 \cdot k}\right] \leq \frac{\mathbb{E}[e^{tY}]}{e^{t \cdot (1+\varepsilon)^2 \cdot k}}$$

Lemma (Chernoff for $\chi^2(k)$)

If $Y = \sum_{i=1}^{k} X_i^2$ is a $\chi^2(k)$ random variable with k degrees of freedom (recall $X_i \sim \mathcal{N}(0,1)$), then

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \exp\left(-\frac{3}{4} \cdot d\varepsilon^2\right)$$

- Let $t \in (0, 1/2)$ be a parameter
- •

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] = \Pr\left[e^{tY} > e^{t \cdot (1+\varepsilon)^2 \cdot k}\right] \le \frac{\mathbb{E}[e^{tY}]}{e^{t \cdot (1+\varepsilon)^2 \cdot k}}$$

• By independence:

$$\mathbb{E}[e^{tY}] = \mathbb{E}\left[\exp\left(\sum_{i=1}^k t \cdot X_i^2\right)\right] = \prod_{i=1}^k \mathbb{E}[e^{tX_i^2}]$$

Need to compute $\mathbb{E}[e^{tX_i^2}]$, where $X_i \sim \mathcal{N}(0,1)$

Need to compute $\mathbb{E}[e^{tX_i^2}]$, where $X_i \sim \mathcal{N}(0,1)$

PDF of X_i:

$$f_{X_i}(x) = \frac{1}{\sqrt{2\pi}} \cdot \exp(-x^2/2)$$

Need to compute $\mathbb{E}[e^{tX_i^2}]$, where $X_i \sim \mathcal{N}(0,1)$

• PDF of *X_i*:

$$f_{X_i}(x) = \frac{1}{\sqrt{2\pi}} \cdot \exp(-x^2/2)$$

Thus we know that

$$\int_{-\infty}^{\infty} f_{X_i}(x) dx = 1$$

Need to compute $\mathbb{E}[e^{tX_i^2}]$, where $X_i \sim \mathcal{N}(0,1)$

PDF of X_i:

$$f_{X_i}(x) = \frac{1}{\sqrt{2\pi}} \cdot \exp(-x^2/2)$$

Thus we know that

$$\int_{-\infty}^{\infty} f_{X_i}(x) dx = 1$$

$$\mathbb{E}[e^{tX_i^2}] = \int_{-\infty}^{\infty} f_{X_i}(x) \cdot e^{tx^2} dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \cdot e^{-x^2/2} \cdot e^{tx^2} dx$$
"Pzobabi Lity that $X_i = x^{-1}$

$$= \frac{1}{\sqrt{2ii}} \int_{-\infty}^{\infty} \exp\left(-\left(1-2t\right) \frac{x^{2}}{2}\right) dx$$

Need to compute $\mathbb{E}[e^{tX_i^2}]$, where $X_i \sim \mathcal{N}(0,1)$

PDF of X_i:

$$f_{X_i}(x) = \frac{1}{\sqrt{2\pi}} \cdot \exp(-x^2/2)$$

Thus we know that

$$\int_{-\infty}^{\infty} f_{X_i}(x) dx = 1$$

•

$$\mathbb{E}[e^{tX_i^2}] = \int_{-\infty}^{\infty} f_{X_i}(x) \cdot e^{tx^2} dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \cdot e^{-x^2/2} \cdot e^{tx^2} dx$$

• Change of variables $z = x\sqrt{1-2t}$

$$\mathbb{E}[e^{tX_i^2}] = \frac{1}{\sqrt{2\pi} \cdot \sqrt{1 - 2t}} \cdot \int_{-\infty}^{\infty} e^{-z^2/2} dz = \frac{1}{\sqrt{1 - 2t}}$$

Putting everything together:

Putting everything together:

•

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \frac{\mathbb{E}[e^{tY}]}{e^{t \cdot (1+\varepsilon)^2 \cdot k}}$$

Putting everything together:

0

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \frac{\mathbb{E}[e^{tY}]}{e^{t \cdot (1+\varepsilon)^2 \cdot k}}$$

$$\mathbb{E}[e^{tY}] = \prod_{i=1}^k \mathbb{E}[e^{tX_i^2}] = \left(\frac{1}{\sqrt{1-2t}}\right)^k = \left(\frac{1}{\sqrt{1-2t}}\right)^k$$

Putting everything together:

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \frac{\mathbb{E}[e^{tY}]}{e^{t\cdot (1+\varepsilon)^2 \cdot k}}$$

$$\mathbb{E}[e^{tY}] = \prod_{i=1}^k \mathbb{E}[e^{tX_i^2}] = \left(\frac{1}{\sqrt{1-2t}}\right)^k$$

•

$$\Pr[\mathsf{Y} > (1+\varepsilon)^2 \cdot k] \leq e^{-t \cdot (1+\varepsilon)^2 \cdot k} \cdot (1-2t)^{-k/2} = \left[(1+\varepsilon)^2 e^{1-(1+\varepsilon)^2} \right]^{k/2}$$

• Setting
$$t = (1/2) \cdot \left(1 - \frac{1}{(1+\varepsilon)^2}\right)$$
 above
$$1 - 2t = \left(1+\epsilon\right)^2 - t\left(1+\epsilon\right)^2 = \frac{1}{2}\left(1 - \left(1+\epsilon\right)^2\right)$$

$$(1-z+)^{-k/2}=(1+\epsilon)^k$$

Concentration of χ^2 random variables

Putting everything together:

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \frac{\mathbb{E}[e^{tY}]}{e^{t\cdot (1+\varepsilon)^2 \cdot k}}$$

•

$$\mathbb{E}[e^{tY}] = \prod_{i=1}^k \mathbb{E}[e^{tX_i^2}] = \left(\frac{1}{\sqrt{1-2t}}\right)^k$$

•

$$\Pr[\mathsf{Y} > (1+\varepsilon)^2 \cdot k] \le e^{-t \cdot (1+\varepsilon)^2 \cdot k} \cdot (1-2t)^{-k/2} = \left[(1+\varepsilon)^2 e^{1-(1+\varepsilon)^2} \right]^{k/2}$$

- Setting $t = (1/2) \cdot \left(1 \frac{1}{(1+\varepsilon)^2}\right)$ above
- Use $ln(1+x) \le x x^4/4$ for $x \in [0,1]$

$$\Pr[Y > (1+\varepsilon)^2 \cdot k] \le \exp(-(3/4) \cdot k\varepsilon^2)$$

Concentration of χ^2 random variables

Putting everything together:

$$\Pr[Y > (1+arepsilon)^2 \cdot k] \leq rac{\mathbb{E}[e^{tY}]}{e^{t\cdot (1+arepsilon)^2 \cdot k}}$$

•

$$\mathbb{E}[e^{tY}] = \prod_{i=1}^k \mathbb{E}[e^{tX_i^2}] = \left(\frac{1}{\sqrt{1-2t}}\right)^k$$

 $\Pr[Y > (1+\varepsilon)^2 \cdot k] \le e^{-t \cdot (1+\varepsilon)^2 \cdot k} \cdot (1-2t)^{-k/2} = \left\lceil (1+\varepsilon)^2 e^{1-(1+\varepsilon)^2} \right\rceil^{k/2}$

•

• Setting
$$t = (1/2) \cdot \left(1 - \frac{1}{(1+\varepsilon)^2}\right)$$
 above

- - Use $ln(1+x) \le x x^4/4$ for $x \in [0,1]$ $\Pr[Y > (1+\varepsilon)^2 \cdot k] < \exp(-(3/4) \cdot k\varepsilon^2)$
 - Similar result for $\Pr[Y < (1-\varepsilon)^2 \cdot k]$ **Practice problem.**

- Introduction
 - Administrivia
 - Why Reduce Dimensions?
 - Background: Continuous Probability Distributions

- Main Problem
 - Johnson-Lindenstrauss Lemma

Acknowledgements

- Input: m points $x_1, \ldots, x_m \in \mathbb{R}^n$.
- **Output:** m points $y_1, \ldots, y_m \in \mathbb{R}^d$, where $d \ll n$ such that

- Input: m points $x_1, \ldots, x_m \in \mathbb{R}^n$.
- **Output:** m points $y_1, \ldots, y_m \in \mathbb{R}^d$, where $d \ll n$ such that

$$||y_a - y_b||_2 \approx ||x_a - x_b||_2 \quad \forall a, b \in [m]$$

- **Input:** m points $x_1, \ldots, x_m \in \mathbb{R}^n$.
- **Output:** m points $y_1, \ldots, y_m \in \mathbb{R}^d$, where $d \ll n$ such that

$$||y_a - y_b||_2 \approx ||x_a - x_b||_2 \quad \forall a, b \in [m]$$

Theorem (Johnson-Lindenstrauss Theorem)

Let $x_1, \ldots, x_m \in \mathbb{R}^n$ and $\varepsilon \in (0,1)$. For $d = O(\log(m)/\varepsilon^2)$ there exist points $y_1, \ldots, y_m \in \mathbb{R}^d$ such that:

$$(1-\varepsilon) \cdot \|x_a - x_b\|_2 \le \|y_a - y_b\|_2 \le (1+\varepsilon) \cdot \|x_a - x_b\|_2 \quad \forall a, b \in [m]$$

Moreover, the points $y_j = Lx_j$, where $L \in \mathbb{R}^{d \times n}$ is a matrix whose entries $L_{a,b} \sim \mathcal{N}(0,1)$, satisfies the above with probability $\geq 1 - 2/m$.



- **Input:** m points $x_1, \ldots, x_m \in \mathbb{R}^n$.
- **Output:** m points $y_1, \ldots, y_m \in \mathbb{R}^d$, where $d \ll n$ such that

$$||y_a - y_b||_2 \approx ||x_a - x_b||_2 \quad \forall a, b \in [m]$$

Theorem (Johnson-Lindenstrauss Theorem)

Let $x_1, \ldots, x_m \in \mathbb{R}^n$ and $\varepsilon \in (0,1)$. For $d = O(\log(m)/\varepsilon^2)$ there exist points $y_1, \ldots, y_m \in \mathbb{R}^d$ such that:

$$(1-\varepsilon)\cdot \|x_a - x_b\|_2 \le \|y_a - y_b\|_2 \le (1+\varepsilon)\cdot \|x_a - x_b\|_2 \quad \forall a, b \in [m]$$

Moreover, the points $y_j = Lx_j$, where $L \in \mathbb{R}^{d \times n}$ is a matrix whose entries $L_{a,b} \sim \mathcal{N}(0,1)$, satisfies the above with probability $\geq 1 - 2/m$.

- If one of the points is 0 then approximate norm of vectors as well!
- Independent of the original dimension n



Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[\left(1-arepsilon
ight) \ \le \ rac{\|f(v)\|_2}{\sqrt{d}} \le \ \left(1+arepsilon
ight)
ight] \ge 1-2/m^3.$$

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[\left(1-arepsilon
ight) \ \le \ rac{\|f(v)\|_2}{\sqrt{d}} \le \ \left(1+arepsilon
ight)
ight] \ge 1-2/m^3.$$

Proof of theorem given lemma:

• Define linear map $L(v) = f(v)/\sqrt{d}$

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[\left(1-arepsilon
ight) \ \le \ rac{\|f(v)\|_2}{\sqrt{d}} \le \ \left(1+arepsilon
ight)
ight] \ge 1-2/m^3.$$

Proof of theorem given lemma:

- Define linear map $L(v) = f(v)/\sqrt{d}$
- By lemma, for any $u \in \mathbb{R}^n$, we have

$$\Pr[(1-\varepsilon)\cdot \|u\|_2 \le \|L(u)\|_2 \le (1+\varepsilon)\cdot \|u\|_2] \ge 1-2/m^3$$

thus probability of failure (i.e. Large distortion)

is
$$\leq \frac{2}{m}$$
3



Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[\left(1-arepsilon
ight) \ \le \ rac{\|f(v)\|_2}{\sqrt{d}} \le \ \left(1+arepsilon
ight)
ight] \ge 1-2/m^3.$$

Proof of theorem given lemma:

- Define linear map $L(v) = f(v)/\sqrt{d}$
- By lemma, for any $u \in \mathbb{R}^n$, we have $\Pr[(1-\varepsilon) \cdot ||u||_2 \le ||L(u)||_2 \le (1+\varepsilon) \cdot ||u||_2] \ge 1 2/m^3$
- Apply this result and union bound to all vectors $x_a x_b$.
- Probability any failure on the norm $\leq m^2 \cdot 2/m^3 = 2/m$.

 JL Lemma essentially states that if we project a unit vector to a uniformly random d-dimensional subspace we can (almost) preserve the norm!

- JL Lemma essentially states that if we project a unit vector to a uniformly random d-dimensional subspace we can (almost) preserve the norm!
- One advantage of choosing random subspace is that we could *flip the* randomness: consider any d-dimensional space and take vector to be
 uniformly random unit vector

- JL Lemma essentially states that if we project a unit vector to a uniformly random d-dimensional subspace we can (almost) preserve the norm!
- One advantage of choosing random subspace is that we could *flip the* randomness: consider any d-dimensional space and take vector to be
 uniformly random unit vector
- So why not do that?

- JL Lemma essentially states that if we project a unit vector to a uniformly random d-dimensional subspace we can (almost) preserve the norm!
- One advantage of choosing random subspace is that we could flip the randomness: consider any d-dimensional space and take vector to be uniformly random unit vector
- So why not do that?
- A bit cumbersome to get random subspace (need to make L orthonormal so need to use Gram-Schmidt)

 (even though in analysis we can blip the randomness, in the algorithm we would need to use G5 to get random subspace)

- JL Lemma essentially states that if we project a unit vector to a uniformly random d-dimensional subspace we can (almost) preserve the norm!
- One advantage of choosing random subspace is that we could *flip the* randomness: consider any d-dimensional space and take vector to be
 uniformly random unit vector
- So why not do that?
- A bit cumbersome to get random subspace (need to make L orthonormal - so need to use Gram-Schmidt)
- Just taking Gaussians do the trick without Gram-Schmidt!
- More convenient algorithmically

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[(1-arepsilon) \leq \frac{\|f(v)\|_2}{\sqrt{d}} \leq (1+arepsilon)\right] \geq 1-2/m^3.$$

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[(1-arepsilon) \leq \frac{\|f(v)\|_2}{\sqrt{d}} \leq (1+arepsilon)\right] \geq 1-2/m^3.$$

Proof of upper tail: $\Pr[\|f(v)\|_2 > \sqrt{d} \cdot (1+\varepsilon)] < 1/m^3$

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[(1-arepsilon) \leq \frac{\|f(v)\|_2}{\sqrt{d}} \leq (1+arepsilon)\right] \geq 1-2/m^3.$$

Proof of upper tail: $\Pr[\|f(v)\|_2 > \sqrt{d} \cdot (1+\varepsilon)] < 1/m^3$

• Let
$$X_i = r_i^T v$$

random variable for i^{th} coordinate of f(v)

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[(1-arepsilon) \leq \frac{\|f(v)\|_2}{\sqrt{d}} \leq (1+arepsilon)\right] \geq 1-2/m^3.$$

Proof of upper tail: $\Pr[\|f(v)\|_2 > \sqrt{d} \cdot (1+\varepsilon)] < 1/m^3$

- Let $X_i = r_i^T v$ random variable for i^{th} coordinate of f(v)
- $X_i \sim \mathcal{N}(0, \sum_{i=1}^n v_i^2) = \mathcal{N}(0, 1)$ sum of Gaussians

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[\left(1-arepsilon
ight) \ \le \ rac{\|f(v)\|_2}{\sqrt{d}} \le \ \left(1+arepsilon
ight)
ight] \ge 1-2/m^3.$$

Proof of upper tail: $\Pr[\|f(v)\|_2 > \sqrt{d} \cdot (1+\varepsilon)] < 1/m^3$

- Let $X_i = r_i^T v$ random variable for i^{th} coordinate of f(v)
 - $X_i \sim \mathcal{N}(0, \sum_{i=1}^n v_i^2) = \mathcal{N}(0, 1)$ sum of Gaussians

$$||f(v)||_2^2 = \sum_{i=1}^d (r_i^T v)^2 = \sum_{i=1}^d X_i^2$$

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[(1-arepsilon) \leq \frac{\|f(v)\|_2}{\sqrt{d}} \leq (1+arepsilon)\right] \geq 1-2/m^3.$$

Proof of upper tail: $\Pr[\|f(v)\|_2 > \sqrt{d} \cdot (1+\varepsilon)] < 1/m^3$

•
$$||f(v)||_2^2 = \sum_{i=1}^d (r_i^T v)^2 = \sum_{i=1}^d X_i^2$$

 $\chi^2(d)$ random variable!

Theorem (Johnson-Lindenstrauss Lemma)

Let $v \in \mathbb{R}^n$ such that $||v||_2 = 1$, $\varepsilon \in (0,1)$ and $d = O(\log(m)/\varepsilon^2)$. Let $r_1, \ldots, r_d \in \mathbb{R}^n$ be such that $r_i \sim \mathcal{N}(0,1)$. If we let $f : \mathbb{R}^n \to \mathbb{R}^d$ s.t.

$$f(v) = (r_1^T v, r_2^T v, \dots, r_d^T v)$$

Then

$$\Pr\left[(1-\varepsilon) \le \frac{\|f(v)\|_2}{\sqrt{d}} \le (1+\varepsilon)\right] \ge 1-2/m^3.$$

Proof of upper tail: $\Pr[\|f(v)\|_2 > \sqrt{d} \cdot (1+\varepsilon)] < 1/m^3$

•
$$||f(v)||_2^2 = \sum_{i=1}^d (r_i^T v)^2 = \sum_{i=1}^d X_i^2$$

 $\chi^2(d)$ random variable!

Chernoff:

$$\Pr[\|f(v)\|_2^2 > d \cdot (1+\varepsilon)^2] < \exp(-(3/4) \cdot d\varepsilon^2) < 1/m^3$$

What if I don't like Gaussians?

- Can we even sample from a Gaussian?
- Same results also hold if pick a random matrix with entries uniformly from $\{-1,1\}$ (Rademacher random variables).
- Proof a little more involved (see Jelani's notes for a proof)

How tight is the JL lemma?

How tight is the JL lemma?

Very tight!

Theorem (Noga Alon)

Let $y_0, \ldots, y_n \in \mathbb{R}^d$ such that $1 \leq \|y_i - y_j\|_2 \leq 1 + \varepsilon$ for all $i \neq j$. Then

$$d = \Omega\left(\frac{\log n}{\varepsilon^2 \cdot \log 1/\varepsilon}\right)$$

How tight is the JL lemma?

Very tight!

Theorem (Noga Alon)

Let $y_0, \ldots, y_n \in \mathbb{R}^d$ such that $1 \leq \|y_i - y_j\|_2 \leq 1 + \varepsilon$ for all $i \neq j$. Then

$$d = \Omega\left(\frac{\log n}{\varepsilon^2 \cdot \log 1/\varepsilon}\right)$$

Can I also compress other norms?

How tight is the JL lemma?

Very tight!

Theorem (Noga Alon)

Let $y_0, \ldots, y_n \in \mathbb{R}^d$ such that $1 \leq \|y_i - y_j\|_2 \leq 1 + \varepsilon$ for all $i \neq j$. Then

$$d = \Omega\left(\frac{\log n}{\varepsilon^2 \cdot \log 1/\varepsilon}\right)$$

Can I also compress other norms?

Answer is NO in general.

How tight is the JL lemma?

Very tight!

Theorem (Noga Alon)

Let $y_0, \ldots, y_n \in \mathbb{R}^d$ such that $1 \leq \|y_i - y_j\|_2 \leq 1 + \varepsilon$ for all $i \neq j$. Then

$$d = \Omega\left(\frac{\log n}{\varepsilon^2 \cdot \log 1/\varepsilon}\right)$$

Can I also compress other norms?

- Answer is NO in general.
- [Brinkman, Charikar 2005]: For the ℓ_1 -norm, where $||x||_1 = \sum_{i=1}^n |x_i|$, if want distortion $(1 + \varepsilon)$ dimension must be $\Omega(n^{1/(1+\varepsilon)^2})$

Acknowledgement

- Lecture based largely on Jelani Nelson's and Nick Harvey's notes.
- See Jelani's notes at http://web.mit.edu/minilek/www/jl_notes.pdf
- See Nick's notes at http://www.cs.ubc.ca/~nickhar/W12/Lecture6Notes.pdf

References I



Brinkman, Bo and Charikar, Moses (2005)

On the impossibility of dimension reduction in ℓ_1 Journal of the ACM 52(5), 766–788.



William B. Johnson and Joram Lindenstrauss (1984)

Extensions of Lipschitz mappings into a Hilbert space Contemporary Mathematics, 26:189–206, 1984.