

Today

- 1) Metric embeddings
- 2) Random Projections
- 3) Johnson-Lindenstrauss Lemma
- 4) Applications of (3)

Recall

The Metric Framework

- i) Find a metric in your problem
- ii) Find structure in your metric
- iii) Use metric structure to solve problem

The Concentration Framework

- a) Show (X) true if all RVs near \mathbb{E}
- b) Concentration: each RV at $\mathbb{E} (\pm \log n)$ whp
- c) Union bound: all RVs

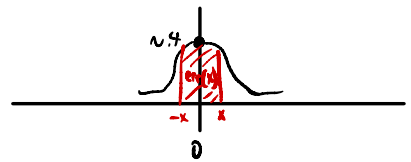
Fact: $\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$

Gaussians RV

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

$$Z \sim N(0, 1) \rightarrow \text{"standard Gaussian"}$$

\mathbb{E} \uparrow var



Amazing Fact 1: $\int_{-\infty}^{\infty} \varphi(x) dx = 1$ \rightarrow More generally, $\int_{-\infty}^{\infty} e^{-ax^2} dx = \sqrt{\frac{\pi}{a}}$

Amazing Fact 2: Rotational symmetry: $(N(0,1), N(0,1), \dots)$ in a uniformly random direction

Dfn. $N(\mu, \sigma^2) = \mu + \sigma Z$ where $Z \sim N(0,1)$ is a "non-standard Gaussian"

Amazing Fact 3: Sum of Gaussians is a Gaussian

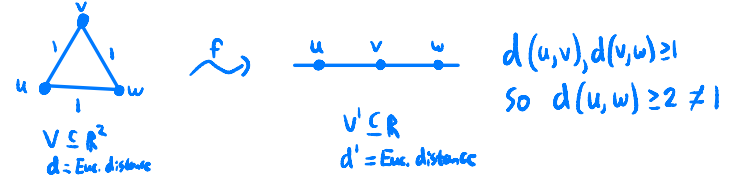
If $X \sim N(0, 1)$, $Y \sim N(0, 1)$, X, Y independent, $a, b \in \mathbb{R}$
then $aX + bY \sim N(0, a^2 + b^2)$

Metric Embeddings: how can we approximate a complex metric w/ a simple one ((ii) today)

An embedding of metric space (V, d) into metric space (V', d') is a function $f: V \rightarrow V'$

(V, d) and (V', d') are isometric iff \exists embedding $f: V \rightarrow V'$ s.t. $d(u, v) = d'(f(u), f(v)) \forall u, v \in V$
 And vice versa

Isometry not always possible, e.g.



f has distortion α if $d(u, v) \leq d'(f(u), f(v)) \leq \alpha \cdot d(u, v) \forall u, v \in V$

E.g. $\alpha = 2$ above

Ideally: α small + (V', d') simpler/more structured than (V, d)

Today: linear embedding of (V, d) into (V', d') for $V \subseteq \mathbb{R}^n, V' \subseteq \mathbb{R}^k, k \ll n$
 $|V| = m$ \uparrow \uparrow
 Euc. distance

w/ distortion $\alpha = \sqrt{\frac{1+\epsilon}{1-\epsilon}}$ for small $\epsilon > 0$

Reduction: for goal, suffices to give linear f s.t. $\|f(w)\|^2 \in [1-\epsilon, 1+\epsilon]$ for $\|w\|^2 \leq m^2$ unit w

g as in goal

$$\|u-v\| \leq \|g(u)-g(v)\| \leq \alpha \|u-v\| \quad \forall u, v \in V$$

$$\|u-v\| \leq \|g(u-v)\| \leq \alpha \|u-v\|$$

$$\|z\| \leq \|g(z)\| \leq \alpha \|z\| \quad \forall z \in V-V$$

$$1 \leq \|g(\frac{z}{\|z\|})\| \leq \alpha \quad \forall z \in V-V$$

$$1 \leq \|g(w)\|^2 \leq \alpha^2 \quad \forall w = \frac{z}{\|z\|} \quad \forall z \in V-V$$

Let $f := \sqrt{1-\epsilon} \cdot g$ and $\alpha = \sqrt{\frac{1+\epsilon}{1-\epsilon}}$

goal

$$1-\epsilon \leq (1-\epsilon) \|g(w)\|^2 \leq 1+\epsilon \quad \forall w = \frac{z}{\|z\|} \quad \forall z \in V-V$$

$$1-\epsilon \leq \|\sqrt{1-\epsilon} g(w)\|^2 \leq 1+\epsilon \quad \forall w = \frac{z}{\|z\|} \quad \forall z \in V-V$$

$$\|f(w)\|^2 \in [1-\epsilon, 1+\epsilon] \quad \forall w = \frac{z}{\|z\|} \quad \forall z \in V-V$$

χ^2 Distributions + Concentration

X is a chi-squared RV w/ k degrees of freedom if

$$X = \sum_{i=1}^k z_i^2$$

where each $z_i \sim N(0,1)$ independently

Notated $X \sim \chi_k^2$

Claim: $\mathbb{E}[X] = k$ for $X \sim \chi_k^2$

Have $\mathbb{E}[z_i^2] = \text{Var}(z_i) + \mathbb{E}[z_i]^2 = 1 + 0$

Claim follows by LoE

Can't apply Chernoff for concentration b/c not $\in \{0,1\}$ (or even bounded)

Nonetheless, similar proof works

Can get $-O(\epsilon)$ in exponent instead of $-O(\epsilon^2)$ $\forall \epsilon \in (0,1)$

Skip, come back if time

Claim: $\Pr(|X - \mathbb{E}[X]| \geq \epsilon \cdot k) \leq 2 \cdot \exp(-\epsilon^2 k / 8)$ for $X \sim \chi_k^2$ ($\epsilon \in (0,1)$)

Will prove upper tail; lower tail symmetric; Claim follows by union bound

Let $t = \epsilon/4$ so $t \in (0, 1/4)$ and $3t + \epsilon t = \frac{3}{4}\epsilon + \frac{1}{4}\epsilon^2 \geq \epsilon^2/8 \forall \epsilon \in (0,1)$

$$\Pr(X - \mathbb{E}[X] \geq \epsilon k) = \Pr(X \geq (1+\epsilon) \cdot k) \leq \Pr(e^{tX} \geq e^{(1+\epsilon) \cdot tk}) \stackrel{\text{Markov}}{\leq} \frac{\mathbb{E}[e^{tX}]}{e^{(1+\epsilon) \cdot tk}} \stackrel{z_i \text{ ind.}}{=} \prod_{i=1}^k \frac{\mathbb{E}[e^{t z_i^2}]}{e^{(1+\epsilon) \cdot tk}}$$

$$\text{Now calculate } \mathbb{E}[e^{t z_i^2}] = \frac{1}{\sqrt{2\pi}} \int_a e^{ta^2} \cdot e^{-a^2/2} da = \frac{1}{\sqrt{2\pi}} \int_a e^{-a^2(\frac{1}{2}-t)} da \stackrel{\text{Gaussian integral}}{=} \frac{1}{\sqrt{2(\frac{1}{2}-t)}}$$

$$\text{Plugging } \mathbb{E}[e^{t z_i^2}] \text{ in, get } \Pr(X \geq (1+\epsilon) \cdot k) \leq \left(\frac{1}{e^{(1+\epsilon) \cdot t} \sqrt{1-2t}} \right)^k$$

$$\text{But } \frac{1}{e^{(1+\epsilon)t} \sqrt{1-2t}} = \exp(-t - \frac{1}{2} \ln(1-2t) - \epsilon t) \stackrel{\uparrow}{=} \exp(-3t - \epsilon t) \leq \exp(-\epsilon^2/8)$$

Combining above gives result

$e^{-x} \leq 1 - sx$ for $x \in (0,1)$
so $-1 \ln(1-2t) \leq -4t$

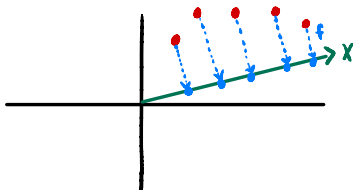
above choice of t

Random Projection

Simplest case: $k=1$

Let $X = (X_1, X_2, \dots, X_m)$ where $X_i \sim N(0,1) \forall i$ be a vector in a uniformly random direction

Let $f: \mathbb{R}^m \rightarrow \mathbb{R}$ be $f(w) = \langle w, X \rangle$



Pick random direction X
 w mapped to how far along X $\text{Proj}(w \rightarrow X)$ is

Claim: $\|f(w)\|^2 \sim \chi_1^2 \forall w \text{ s.t. } \|w\|=1$

$$\|w\|=1 \rightarrow \|w\|^2=1 \rightarrow \sum_i w_i^2=1$$

$$\text{So } \|f(w)\| = \sum_i w_i X_i$$

$\sim N(0, w_1^2 + w_2^2 + \dots + w_m^2)$ by \sum_i Gaussians = Gaussian

$$= N(0,1) \text{ b/c } \|w\|=1$$

Corollaries:

1) $\mathbb{E}[\|f(w)\|^2] = 1 \quad \forall w \text{ s.t. } \|w\|=1 \rightarrow$ "Random Projection Preserves unit Vector length in \mathbb{E} "

2) $\Pr(\|f(w)\|^2 - 1 \geq \epsilon) \leq 2 \cdot \exp(-\epsilon^2/8) \rightarrow$ Vector length even concentrates under random projection

Problem: Upper tail useful but lower tail not $\rightarrow \pm \log n$ from \mathbb{E}
= bad distortion

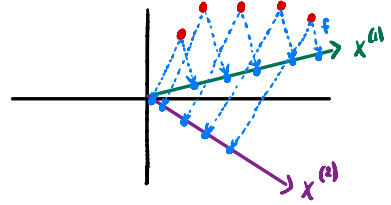
Solution: bring \mathbb{E} up to $\approx \log n$ by repeated trials

Johnson-Lindenstrauss Lemma

Let $X^{(i)} \sim (X_1^{(i)}, X_2^{(i)}, \dots, X_n^{(i)})$ where each $X_j^{(i)} \sim N(0,1)$

Let $f_i(w) := \langle X^{(i)}, w \rangle$

Let $F(w) := (f_1(w), f_2(w), \dots, f_k(w)) \rightarrow$ Equivalent to Aw where $A_{ij} \sim N(0,1)$
So F linear



Claim: $\|F(w)\|^2 \sim \chi_k^2 \quad \forall w \text{ s.t. } \|w\|=1$

$\|F(w)\|^2 = \sum_i (f_i(w))^2 = \sum_i \|f_i(w)\|^2$ and $\|f_i(w)\|^2 \sim \chi_1^2$ by previous claim

Corollaries:

1) $E[\|F(w)\|^2] = k \quad \forall w \text{ s.t. } \|w\|=1 \rightarrow$ So, $k \approx \log n \rightarrow$ good concentration

2) $\Pr(\|F(w)\|^2 - k \geq \epsilon k) \leq 2 \cdot \exp(-\epsilon^2 k / 8) \quad \forall w \text{ s.t. } \|w\|=1$

Let $\tilde{F}(w) := \frac{1}{\sqrt{k}} \cdot F(w) \rightarrow$ Scale down so unit vector mapped to unit vector
 \rightarrow Trivially Poly-time computable (even w/o knowing points)

JL Lemma: For any m^2 unit vectors W and $k = 24 \cdot \frac{\ln m}{\epsilon^2}$ $\tilde{F}: \mathbb{R}^m \rightarrow \mathbb{R}^k$ satisfies

$$\|\tilde{F}(w)\|^2 \in [1-\epsilon, 1+\epsilon] \quad \forall w \in W$$

except w/ $\Pr \leq 1 - \frac{2}{n}$

Note $E[\|\tilde{F}(w)\|^2] = \frac{1}{k} E[\|F(w)\|^2] = 1 \quad (a)$

$\|\tilde{F}(w)\|^2 \notin [1-\epsilon, 1+\epsilon]$ only if $\|F(w)\|^2 \notin [(1-\epsilon)k, (1+\epsilon)k]$ only if $|\|F(w)\|^2 - k| \geq \epsilon k$

But $\Pr(\|F(w)\|^2 - k \geq \epsilon k) \leq 2 \cdot \exp(-3 \ln m) = \frac{2}{n^3} \quad (b)$

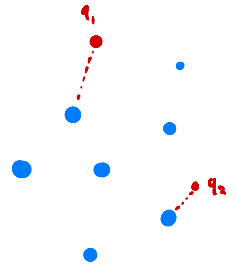
By union bound, $\tilde{F}(w) \in [1-\epsilon, 1+\epsilon] \quad \forall n^2 w$ except w/ $\Pr \leq \frac{2}{n} \quad (c)$

Application of JL

$\sqrt{\frac{1+\epsilon}{1-\epsilon}}$ -NN-Search: initially given $X \subseteq \mathbb{R}^n$, w/ $|X|=m$ (i)

repeatedly given queries q_1, q_2, \dots

$\forall q_i$, return $x \in X$ s.t. $d(q_i, x) \leq \sqrt{\frac{1+\epsilon}{1-\epsilon}} \cdot d(q_i, X)$
as fast as possible



Naive solution: $O(n \cdot m)$ time per query

JL solution: $O\left(\frac{\log n}{\epsilon^2} \cdot m\right)$ per query

let $\hat{F}: \mathbb{R}^n \rightarrow \mathbb{R}^k$ be linear embedding w/ distortion $\sqrt{\frac{1+\epsilon}{1-\epsilon}}$ from JL
for $k = O\left(\frac{\log n}{\epsilon^2}\right)$ so

let $Y := \{\hat{F}(x) : x \in X\}$ (ii)

For query q_i ,

return $x \in X$ s.t. $\hat{F}(x) = y$ where $y = \arg \min_{y \in Y} d(y, q_i)$ (iii)

Easy to verify result is (approximately) correct

Takes $O\left(\frac{\log n}{\epsilon^2} \cdot m\right)$ per query

\hookrightarrow Can even reduce to $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ w/ a little more work!