

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2022.

Lecture 12

- Problem Set 2 is due Friday, 11:59pm.
- No quiz this week.
- The exam will be held next Thursday in class.
- We will do some midterm review in class on Tuesday. I will also hold additional office hours for midterm prep, TBD.

Last Class: The Johnson-Lindenstrauss Lemma

- Intro to dimensionality reduction and low-distortion embeddings.
- Statement of the JL Lemma: we can obtain low-distortion embeddings for **any set of points** via random projection.

This Class:

- Reduction of the JL Lemma to the ‘distributional JL Lemma’.
- Proof of the distributional JL lemma.
- Example application to clustering.

The Johnson-Lindenstrauss Lemma

Johnson-Lindenstrauss Lemma: For any set of points $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ and $\epsilon > 0$ there exists a linear map $\mathbf{\Pi} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ such that $m = O\left(\frac{\log n}{\epsilon^2}\right)$ and letting $\tilde{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all i, j : $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.

Further, if $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ has each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, it satisfies the guarantee with high probability.

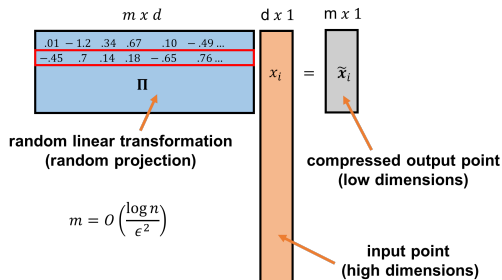
For $d = 1$ trillion, $\epsilon = .05$, and $n = 100,000$, $m \approx 6600$.

Very surprising! Powerful result with a simple construction: applying a random linear transformation to a set of points preserves distances between all those points with high probability.

Random Projection

For any $\vec{x}_1, \dots, \vec{x}_n$ and $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ with each entry chosen i.i.d. from $\mathcal{N}(0, 1/m)$, with high probability, letting $\tilde{x}_i = \mathbf{\Pi}\vec{x}_i$:

For all i, j : $(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{x}_i - \tilde{x}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2$.



- $\mathbf{\Pi}$ is known as a **random projection**. It is a random linear function, mapping length d vectors to length m vectors.
- $\mathbf{\Pi}$ is **data oblivious**. Stark contrast to methods like PCA.

Algorithmic Considerations

- Many alternative constructions: ± 1 entries, sparse (most entries 0), Fourier structured, etc. \implies more efficient computation of $\tilde{\mathbf{x}}_j = \mathbf{\Pi}\vec{x}_j$.
- Data oblivious property means that once $\mathbf{\Pi}$ is chosen, $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ can be computed in a stream with little memory.
- Memory needed is just $O(d + nm)$ vs. $O(nd)$ to store the full data set.
- Compression can also be easily performed in parallel on different servers.
- When new data points are added, can be easily compressed, without updating existing points.

Distributional JL

The Johnson-Lindenstrauss Lemma is a direct consequence of a closely related lemma:

Distributional JL Lemma: Let $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$. If we set $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then **for any** $\vec{y} \in \mathbb{R}^d$, with probability $\geq 1 - \delta$

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

Applying a random matrix $\mathbf{\Pi}$ to any vector \vec{y} preserves \vec{y} 's norm with high probability.

- Like a low-distortion embedding, but for the length of a compressed vector rather than distances between vectors.
- Can be proven from first principles.

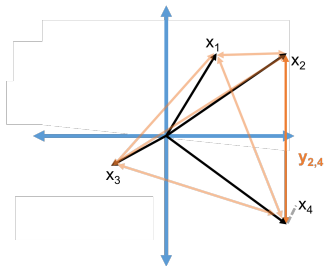
$\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection matrix. d : original dimension. m : compressed dimension, ϵ : embedding error, δ : embedding failure prob.

Distributional JL \implies JL

Distributional JL Lemma \implies JL Lemma: Distributional JL show that a random projection Π preserves the **norm** of any y . The main JL Lemma says that Π preserves **distances** between vectors.

Since Π is **linear** these are the same thing!

Proof: Given $\vec{x}_1, \dots, \vec{x}_n$, define $\binom{n}{2}$ vectors \vec{y}_{ij} where $\vec{y}_{ij} = \vec{x}_i - \vec{x}_j$.



- If we choose Π with $m = O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$, for each \vec{y}_{ij} with probability $\geq 1 - \delta$ we have:

$$(1 - \epsilon)\|\vec{v}_i - \vec{v}_j\|_2 \leq \|\Pi\vec{v}_i - \Pi\vec{v}_j\|_2 \leq (1 + \epsilon)\|\vec{v}_i - \vec{v}_j\|_2$$

Distributional JL \implies JL

Claim: If we choose $\mathbf{\Pi}$ with i.i.d. $\mathcal{N}(0, 1/m)$ entries and $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, letting $\tilde{\mathbf{x}}_i = \mathbf{\Pi}\vec{x}_i$, for each pair \vec{x}_i, \vec{x}_j with probability $\geq 1 - \delta'$ we have:

$$(1 - \epsilon)\|\vec{x}_i - \vec{x}_j\|_2 \leq \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2 \leq (1 + \epsilon)\|\vec{x}_i - \vec{x}_j\|_2.$$

With what probability are all pairwise distances preserved?

Union bound: With probability $\geq 1 - \binom{n}{2} \cdot \delta'$ all pairwise distances are preserved.

Apply the claim with $\delta' = \delta/\binom{n}{2}$. \implies for $m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$, all pairwise distances are preserved with probability $\geq 1 - \delta$.

$$m = O\left(\frac{\log(1/\delta')}{\epsilon^2}\right) = O\left(\frac{\log(\binom{n}{2}/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n^2/\delta)}{\epsilon^2}\right) = O\left(\frac{\log(n/\delta)}{\epsilon^2}\right)$$

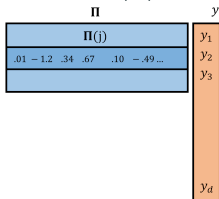
Yields the JL lemma.

Distributional JL Proof

Distributional JL Lemma: Let $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$. If we set $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then for any $\vec{y} \in \mathbb{R}^d$, with probability $\geq 1 - \delta$

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

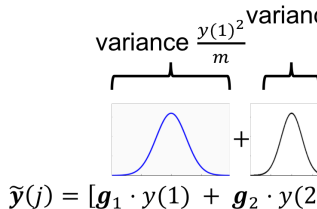
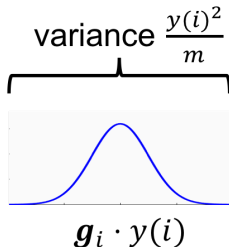
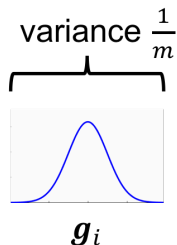
- Let $\tilde{\mathbf{y}}$ denote $\mathbf{\Pi}\vec{y}$ and let $\mathbf{\Pi}(j)$ denote the j^{th} row of $\mathbf{\Pi}$.
- For any j , $\tilde{\mathbf{y}}(j) = \langle \mathbf{\Pi}(j), \vec{y} \rangle = \sum_{i=1}^d \mathbf{g}_i \cdot \vec{y}(i)$ where $\mathbf{g}_i \sim \mathcal{N}(0, 1/m)$.



$\vec{y} \in \mathbb{R}^d$: arbitrary vector, $\tilde{\mathbf{y}} \in \mathbb{R}^m$: compressed vector, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection. d : original dim. m : compressed dim, ϵ : error, δ : failure prob.

Distributional JL Proof

- Let $\tilde{\mathbf{y}}$ denote $\mathbf{\Pi}\vec{\mathbf{y}}$ and let $\mathbf{\Pi}(j)$ denote the j^{th} row of $\mathbf{\Pi}$.
- For any j , $\tilde{\mathbf{y}}(j) = \langle \mathbf{\Pi}(j), \vec{\mathbf{y}} \rangle = \sum_{i=1}^d \mathbf{g}_i \cdot \vec{\mathbf{y}}(i)$ where $\mathbf{g}_i \sim \mathcal{N}(0, 1/m)$.
- $\mathbf{g}_i \cdot \vec{\mathbf{y}}(i) \sim \mathcal{N}(0, \frac{\vec{\mathbf{y}}(i)^2}{m})$: normally distributed with variance $\frac{\vec{\mathbf{y}}(i)^2}{m}$.



What is the distribution of $\tilde{\mathbf{y}}(j)$? Also Gaussian!

$\vec{\mathbf{y}} \in \mathbb{R}^d$: arbitrary vector, $\tilde{\mathbf{y}} \in \mathbb{R}^m$: compressed vector, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection mapping $\vec{\mathbf{y}} \rightarrow \tilde{\mathbf{y}}$. $\mathbf{\Pi}(j)$: j^{th} row of $\mathbf{\Pi}$, d : original dimension. m : compressed dimension, \mathbf{g}_i : normally distributed random variable.

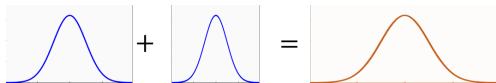
Distributional JL Proof

Letting $\tilde{\mathbf{y}} = \mathbf{\Pi}\vec{y}$, we have $\tilde{y}(j) = \langle \mathbf{\Pi}(j), \vec{y} \rangle$ and:

$$\tilde{y}(j) = \sum_{i=1}^d \mathbf{g}_i \cdot \vec{y}(i) \text{ where } \mathbf{g}_i \cdot \vec{y}(i) \sim \mathcal{N}\left(0, \frac{\vec{y}(i)^2}{m}\right).$$

Stability of Gaussian Random Variables. For independent $a \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $b \sim \mathcal{N}(\mu_2, \sigma_2^2)$ we have:

$$a + b \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$



Thus, $\tilde{y}(j) \sim \mathcal{N}\left(0, \frac{\vec{y}(1)^2}{m} + \frac{\vec{y}(2)^2}{m} + \dots + \frac{\vec{y}(d)^2}{m} \frac{\|\vec{y}\|_2^2}{m}\right)$ I.e., $\tilde{\mathbf{y}}$ itself is a random Gaussian vector. **Rotational invariance of the Gaussian distribution.**

$\vec{y} \in \mathbb{R}^d$: arbitrary vector, $\tilde{\mathbf{y}} \in \mathbb{R}^m$: compressed vector, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection mapping $\vec{y} \mapsto \tilde{\mathbf{y}}$ $\mathbf{\Pi}(i)$: i th row of $\mathbf{\Pi}$ d : original dimension, m : com-

Distributional JL Proof

So far: Letting $\mathbf{\Pi} \in \mathbb{R}^{d \times m}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$, for any $\vec{y} \in \mathbb{R}^d$, letting $\tilde{\mathbf{y}} = \mathbf{\Pi}\vec{y}$:

$$\tilde{y}(j) \sim \mathcal{N}(0, \|\vec{y}\|_2^2/m).$$

What is $\mathbb{E}[\|\tilde{\mathbf{y}}\|_2^2]$?

$$\begin{aligned}\mathbb{E}[\|\tilde{\mathbf{y}}\|_2^2] &= \mathbb{E}\left[\sum_{j=1}^m \tilde{y}(j)^2\right] = \sum_{j=1}^m \mathbb{E}[\tilde{y}(j)^2] \\ &= \sum_{j=1}^m \frac{\|\vec{y}\|_2^2}{m} = \|\vec{y}\|_2^2\end{aligned}$$

So $\tilde{\mathbf{y}}$ has the right norm in expectation.

How is $\|\tilde{\mathbf{y}}\|_2^2$ distributed? Does it concentrate?

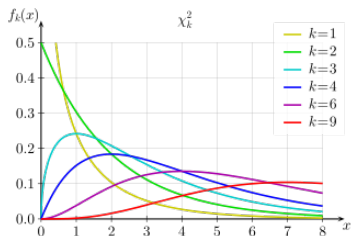
$\vec{y} \in \mathbb{R}^d$: arbitrary vector, $\tilde{\mathbf{y}} \in \mathbb{R}^m$: compressed vector, $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$: random projection mapping $\vec{y} \rightarrow \tilde{\mathbf{y}}$. $\mathbf{\Pi}(j)$: j^{th} row of $\mathbf{\Pi}$, d : original dimension. m : compressed dimension, \mathbf{g}_j : normally distributed random variable

Distributional JL Proof

So far: Letting $\mathbf{\Pi} \in \mathbb{R}^{d \times m}$ have each entry chosen i.i.d. as $\mathcal{N}(0, 1/m)$, for any $\vec{y} \in \mathbb{R}^d$, letting $\tilde{\mathbf{y}} = \mathbf{\Pi}\vec{y}$:

$$\tilde{y}(j) \sim \mathcal{N}(0, \|\vec{y}\|_2^2/m) \text{ and } \mathbb{E}[\|\tilde{\mathbf{y}}\|_2^2] = \|\vec{y}\|_2^2$$

$\|\tilde{\mathbf{y}}\|_2^2 = \sum_{j=1}^m \tilde{y}(j)^2$ a **Chi-Squared random variable with m degrees of freedom** (a sum of m squared independent Gaussians)

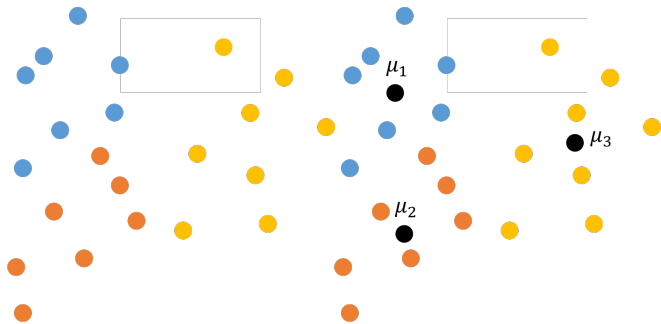


Lemma: (Chi-Squared Concentration) Letting Z be a Chi-Squared random variable with m degrees of freedom,

$$\Pr[|Z - \mathbb{E}Z| > \epsilon \mathbb{E}Z] < 2e^{-m\epsilon^2/8}.$$

Example Application: k -means clustering

Goal: Separate n points in d dimensional space into k groups.



k-means Objective: $Cost(\mathcal{C}_1, \dots, \mathcal{C}_k) = \min_{\mathcal{C}_1, \dots, \mathcal{C}_k} \sum_{j=1}^k \sum_{\vec{x} \in \mathcal{C}_k} \|\vec{x} - \mu_j\|_2^2.$

Write in terms of distances:

$$Cost(\mathcal{C}_1, \dots, \mathcal{C}_k) = \min_{\mathcal{C}_1, \dots, \mathcal{C}_k} \sum_{j=1}^k \sum_{\vec{x}_1, \vec{x}_2 \in \mathcal{C}_k} \|\vec{x}_1 - \vec{x}_2\|_2^2$$

Example Application: k -means clustering

k-means Objective: $Cost(\mathcal{C}_1, \dots, \mathcal{C}_k) = \min_{\mathcal{C}_1, \dots, \mathcal{C}_k} \sum_{j=1}^k \sum_{\vec{x}_1, \vec{x}_2 \in \mathcal{C}_k} \|\vec{x}_1 - \vec{x}_2\|_2^2$

If we randomly project to $m = O\left(\frac{\log n}{\epsilon^2}\right)$ dimensions, for all pairs \vec{x}_1, \vec{x}_2 ,

$$(1 - \epsilon)\|\vec{x}_1 - \vec{x}_2\|_2^2 \leq \|\tilde{\mathbf{x}}_1 - \tilde{\mathbf{x}}_2\|_2^2 \leq (1 + \epsilon)\|\vec{x}_1 - \vec{x}_2\|_2^2 \implies$$

Letting $\overline{Cost}(\mathcal{C}_1, \dots, \mathcal{C}_k) = \min_{\mathcal{C}_1, \dots, \mathcal{C}_k} \sum_{j=1}^k \sum_{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2 \in \mathcal{C}_k} \|\tilde{\mathbf{x}}_1 - \tilde{\mathbf{x}}_2\|_2^2$

$$(1 - \epsilon)Cost(\mathcal{C}_1, \dots, \mathcal{C}_k) \leq \overline{Cost}(\mathcal{C}_1, \dots, \mathcal{C}_k) \leq (1 + \epsilon)Cost(\mathcal{C}_1, \dots, \mathcal{C}_k).$$

Upshot: Can cluster in m dimensional space (much more efficiently) and minimize $\overline{Cost}(\mathcal{C}_1, \dots, \mathcal{C}_k)$. The optimal set of clusters will have true cost within $1 + c\epsilon$ times the true optimal. **Good exercise to prove this.**