

Lecture 06: Applications I

Nikola Zlatanov

Innopolis University

Advanced Statistics

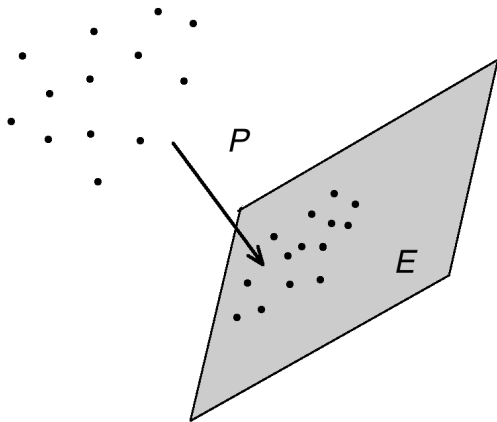
13-th of March to 20-th of March, 2023

Dimension Reduction

- We started this course by observing that high dimensions could be a major problem, due to the curse of dimensionality.
- Why don't we then fix the problem due to high-dimensions as follows: Take the high-dimensional data and transform it into a low-dimensional data.
- Let there be N high-dimensional vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, where $\mathbf{x}_j^T = [x_1, x_2, \dots, x_d]$, for $j = 1, 2, \dots, N$, where d is very high.
- We would like to make a transformation of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ into $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$, where $\mathbf{y}_j^T = [y_1, y_2, \dots, y_n]$, for $j = 1, 2, \dots, N$, where $n \ll d$, and yet the geometry of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ are preserved in $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$.

Dimension Reduction

- First we have to define what do we mean by “the geometry of the data”: By “the geometry of the data”, we mean the pairwise distances between the original data.



Is this possible? It turns out it is possible if $n = O(\ln(N)) \ll d$.

Dimension Reduction

- Thm (Johnson-Lindenstrauss Lemma): \forall vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, where $\mathbf{x}_j \in \mathbb{R}^d$, for $j = 1, 2, \dots, N$, there exists a linear map $T : \mathbf{x}_j \rightarrow \mathbf{y}_j$, where $\mathbf{y}_j \in \mathbb{R}^n$ and $n \ll d$ such that the following holds

$$\Pr \{ (1 - \delta) \|\mathbf{x}_m - \mathbf{x}_j\|_2 \leq \|\mathbf{y}_m - \mathbf{y}_j\|_2 \leq (1 + \delta) \|\mathbf{x}_m - \mathbf{x}_j\|_2 \} \geq 1 - \epsilon \quad (1)$$

for any $j \neq m$, where $j = 1, 2, \dots, N$ and $m = 1, 2, \dots, N$, and small $\delta > 0$ and $\epsilon > 0$ if

$$n > \frac{1}{c} \left(\ln(N) + \frac{1}{2} \ln \left(\frac{1}{\epsilon} \right) + \frac{1}{2} \ln(2) \right)$$

holds where

$$c = \frac{\delta(2 - \delta)}{k} \min \left\{ \frac{\delta(2 - \delta)}{k}, 1 \right\}, \quad (2)$$

where k is some constant.

Dimension Reduction

Proof:

- The vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ and $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$ are all deterministic.
- However, we will use a probabilistic method to find the mapping $T : \mathbf{x}_j \rightarrow \mathbf{y}_j$.
- Specifically, we will choose a linear map $T(\cdot)$ at random, and then we will prove that the linear map $T(\cdot)$ satisfies the properties that we seek.
- Now, a linear map is simply a multiplication of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ by a matrix \mathbf{T} , of size $n \times d$, to obtain $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$.
- Hence, if we will choose a linear map $T(\cdot)$ at random, this means that we should choose a matrix \mathbf{T} at random.
- How is it possible that a randomly chosen matrix \mathbf{T} would work?

Dimension Reduction

- Let's have a random matrix \mathbf{G} of size $n \times d$ populated by i.i.d. Gaussian entries, i.e., the (i, j) -th element of \mathbf{G} , denoted by G_{ij} , is generated i.i.d. according to the Gaussian distribution $N(0, 1)$ for $i = 1, 2, \dots, n$ and $j = 1, 2, \dots, d$.
- Let us have a fixed vector $\mathbf{z} \in \mathbb{R}^d$, with norm $\|\mathbf{z}\|_2$.
- Now, let's investigate the distribution of $\mathbf{G}\mathbf{z}$. The i -th element of $\mathbf{G}\mathbf{z}$, denoted by $(\mathbf{G}\mathbf{z})_i$ is given by

$$(\mathbf{G}\mathbf{z})_i = \sum_{j=1}^d G_{ij} z_j \sim N \left(0, \sum_{j=1}^d z_j^2 \right) = N(0, \|\mathbf{z}\|_2^2) \stackrel{(a)}{=} N(0, 1) \quad (3)$$

where (a) holds if and only if $\|\mathbf{z}\|_2 = 1$.

Dimension Reduction

On the other hand, let's fix two vectors, \mathbf{x}_m and \mathbf{x}_l . Then, the vector $\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l$, according to (3), has i.i.d. zero-mean Gaussian elements each with variance $\|\mathbf{x}_m - \mathbf{x}_l\|_2$.

Hence, if we normalize $\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l$ by $\|\mathbf{x}_m - \mathbf{x}_l\|_2$, the vector

$$\frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2}$$

will have i.i.d. zero-mean Gaussian elements each with variance one. As a result, we can use the Thin-Shell Theorem, which states that

$$\begin{aligned} & \Pr \left\{ \left| \frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \leq \delta \sqrt{n} \right\} \\ & \geq 1 - 2 \exp \left(-n \frac{\delta(2-\delta)}{k} \min \left\{ \frac{\delta(2-\delta)}{k}, 1 \right\} \right), \end{aligned} \quad (4)$$

where k is some constant.

Dimension Reduction

We will now prove the main theorem as follows. We will select a pair of vectors \mathbf{x}_m and \mathbf{x}_l . We will prove (1) for the selected \mathbf{x}_m and \mathbf{x}_l .

Then, we will use the union bound to prove that the theorem holds for all pairs satisfying the given condition.

We start with the selected pair \mathbf{x}_m and \mathbf{x}_l :

$$\begin{aligned}
 & \Pr \left\{ \left| \frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \leq \delta \sqrt{n} \right\} \\
 &= \Pr \left\{ \left| \|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2 - \|\mathbf{x}_m - \mathbf{x}_l\|_2 \sqrt{n} \right| \leq \|\mathbf{x}_m - \mathbf{x}_l\|_2 \delta \sqrt{n} \right\} \\
 &= \Pr \left\{ \|\mathbf{x}_m - \mathbf{x}_l\|_2 \sqrt{n} (1 - \delta) \leq \|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2 \leq \|\mathbf{x}_m - \mathbf{x}_l\|_2 \sqrt{n} (1 + \delta) \right\} \\
 &= \Pr \left\{ \|\mathbf{x}_m - \mathbf{x}_l\|_2 \sqrt{n} (1 - \delta) \leq \|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2 \leq \|\mathbf{x}_m - \mathbf{x}_l\|_2 \sqrt{n} (1 + \delta) \right\} \\
 &= \Pr \left\{ \|\mathbf{x}_m - \mathbf{x}_l\|_2 (1 - \delta) \leq \frac{1}{\sqrt{n}} \|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2 \leq \|\mathbf{x}_m - \mathbf{x}_l\|_2 (1 + \delta) \right\} \\
 &= \text{continuation on next page}
 \end{aligned} \tag{5}$$

Dimension Reduction

Continuation of (5)

$$\begin{aligned}
 &\stackrel{(a)}{=} \Pr \{ \|\mathbf{x}_m - \mathbf{x}_l\|_2(1 - \delta) \leq \|\mathbf{T}\mathbf{x}_m - \mathbf{T}\mathbf{x}_l\|_2 \leq \|\mathbf{x}_m - \mathbf{x}_l\|_2(1 + \delta) \} \\
 &\stackrel{(b)}{\geq} 1 - 2 \exp \left(-n \frac{\delta(2 - \delta)}{k} \min \left\{ \frac{\delta(2 - \delta)}{k}, 1 \right\} \right) \stackrel{(c)}{=} 1 - 2e^{-nc} \quad (6)
 \end{aligned}$$

where (a) comes by setting

$$\mathbf{T} = \frac{1}{\sqrt{n}} \mathbf{G},$$

and (c) comes by setting

$$c = \frac{\delta(2 - \delta)}{k} \min \left\{ \frac{\delta(2 - \delta)}{k}, 1 \right\} \quad (7)$$

Dimension Reduction

We now take the union bound over all pairs of N vectors

$$\begin{aligned}
 & \Pr \left\{ \bigcap_{m=1}^N \bigcap_{l=m+1}^N \left| \frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \leq \delta\sqrt{n} \right\} \\
 &= 1 - \Pr \left\{ \bigcup_{k=1}^N \bigcup_{l=k+1}^N \left| \frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \geq \delta\sqrt{n} \right\} \\
 &\geq 1 - \sum_{k=1}^N \sum_{l=1}^N \Pr \left\{ \left| \frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \geq \delta\sqrt{n} \right\} \\
 &= 1 - N^2 \Pr \left\{ \left| \frac{\|\mathbf{G}\mathbf{x}_m - \mathbf{G}\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \geq \delta\sqrt{n} \right\} \\
 &\geq 1 - 2N^2 e^{-nc} \\
 &= 1 - e^{\ln(2) + 2\ln(N) - nc}
 \end{aligned}$$

Dimension Reduction

Now, we want to set n , N , and c such that

$$\Pr \left\{ \left| \frac{\|G\mathbf{x}_m - G\mathbf{x}_l\|_2}{\|\mathbf{x}_m - \mathbf{x}_l\|_2} - \sqrt{n} \right| \leq \delta \sqrt{n} \right\} \geq 1 - e^{\ln(2) + 2\ln(N) - nc} \geq 1 - \epsilon$$

which occurs if

$$1 - e^{\ln(2) + 2\ln(N) - nc} \geq 1 - \epsilon$$

or equivalently if

$$\ln(2) + 2\ln(N) - nc < \ln(\epsilon)$$

or equivalently if

$$\ln(N) < nc - \frac{1}{2} \ln \left(\frac{1}{\epsilon} \right) - \frac{1}{2} \ln(2),$$

where c is given in (7) as function of δ . Q.E.D.

Note! We have lost d . Where is d ?