

Lecture 01: Curse of Dimensionality: Problem and Solutions

Nikola Zlatanov

Innopolis University
Advanced Statistics

30-th of January, 2023

Motivation - Big Data

- Under the 'umbrella' term Big Data, the following is understood to hold:
 - a) Very big number of observations (or data points)
 - b) Very big number of dimensions (or parameters) associated with each data point
- Example for a)
 - The monthly incomes of the residents of Tatarstan
 - Hence, each month we have 3.8 million observations (or data points), which is big.
 - In one year, we will have 45.8 million observations, which is big.
 - BUT NOTE, the dimension is one, since for each person we record only one number.
- Example for b)
 - The HD image of the residents of Tatarstan
 - Each image is an observation and each observation is comprised of $1440 \times 1080 = 1,555,200$ pixels, and each pixel represents one dimension. Hence, 1,555,200 dimensions per observation – big num.
 - Other examples: text, sound, video, human genome, etc

Difficulty - Big Data

- Empirically (i.e., by experience), it is exponentially harder to deal with large number of dimensions than to deal with large number of observations
- Large number of observations is usually not a problem. This has been dealt with tools from classical statistics.
- In fact, the more observations we have, the more we can understand the statistics (i.e., of what is going on) of the underlying random process.
- But large number of dimensions makes the problem of understanding the statistics very difficult.
- In fact, adding more dimensions to an observation makes the problem of understanding the statistics exponentially more difficult.

Volume of Big Data - Curse of Dimensional

- Example of volume of a cube in high dimensions:
 - The volume of a 2-dimensional cube, where each size is a , is $a \times a = a^2$. For $a = 2$, the volume is 2^2
 - The volume of a 3-dimensional cube, where each size is a , is $a \times a \times a = a^3$. For $a = 2$, the volume is 2^3
 - The volume of a d -dimensional cube, where each size is a , is $a \times a \times \dots \times a = a^d$. For $a = 2$, the volume is 2^d .
 - Hence, the volume grows exponentially with the dimension d . This means the for higher dimensions, there is exponentially more space than in lower dimensions.
- This phenomenon is known as the “curse of high-dimensionality” (CD), and is the main problem in DS and ML.
- We have to find methods that can go around the “curse of high-dimensionality”, if possible.
- This course is aimed to help you understand how to do that.

Example of ACT: Covering Number of Balls

- Assume we have an object, described as a set \mathcal{T} , and assume we have N balls of radius r . We would like to completely cover the object with the minimum number of balls. Let us denote this number as $N(\mathcal{T}, r)$. Now how to find $N(\mathcal{T}, r)$ for a given \mathcal{T} and r ?

Example of ACT: Covering Number of Balls

- Def: The covering number of a set $\mathcal{T} \subset \mathbb{R}^d$ and given radius r , is the smallest number of Euclidean balls with radius r that completely cover the set \mathcal{T} .
- The covering number $N(\mathcal{T}, r)$ is a descriptor of the complexity of \mathcal{T} . The less the number is, the easier is to describe \mathcal{T} using only $N(\mathcal{T}, r)$ balls of radius r
- Other meaning: The covering is simply the quantization of \mathcal{T} , whereas the covering number gives us the minimum number of quantized points that describe $N(\mathcal{T}, r)$ with a quantization error r .

Example of ACT: Covering Number of Balls

- More accurate quantization for small r , but more points.
- More accurate quantization for larger r , but less points.

Example of ACT: Covering Number of Balls

- The covering number $N(\mathcal{T}, r)$ suffers from the curse of dimensionality: The higher the dimension d of $\mathcal{T} \subset \mathbb{R}^d$ is, exponentially more balls are needed to cover \mathcal{T} , i.e., it is exponentially harder to describe it via balls of fixed radius.
- Fact: If $\mathcal{T} \subset \mathbb{R}^d$ is the unit (radius one) Euclidean ball in dimension d , and we need to cover it by balls of radius $r = 1/2$, then we would need $N(\mathcal{T}, 1/2) \geq 2^d$ number of balls.
- Note, the above number is exponential in d , hence, CD hits!
- Proof of the Fact: The volume of a r -radius ball in d dimension is

$$V(d, r) = \frac{\pi^{d/2}}{\Gamma(d/2 + 1)} r^d$$

Hence,

$$\frac{V \text{ of big ball}}{V \text{ of small ball}} = \frac{V(d, 1)}{V(d, 1/2)} = \frac{1}{1/2^d} = 2^d$$

On the other hand $V(d, 1) \leq N(\mathcal{T}, 1/2)V(d, 1/2)$, from where we obtain $N(\mathcal{T}, 1/2) \geq V(d, 1)/V(d, 1/2) = 2^d$.

Example of CD: Covering Number of Balls

- Note on the influence of the dimension on the volume of balls:
- Comparing the unit ball with any other ball of radius $r < 1$, we obtain

$$\frac{V \text{ of smaller ball}}{V \text{ of unit ball}} = \frac{V(d, r)}{V(d, 1)} = r^d \rightarrow 0 \text{ as } d \rightarrow \infty$$

- Comparing the unit ball with any other ball of radius $r > 1$, we obtain

$$\frac{V \text{ of larger ball}}{V \text{ of unit ball}} = \frac{V(d, r)}{V(d, 1)} = r^d \rightarrow \infty \text{ as } d \rightarrow \infty$$

- Hence, when the dimension increases, the balls bifurcate into two classes of balls compared to the unit ball
 - One class is balls with $r < 1$: These balls are almost infinitely smaller than the unit ball
 - Second class is balls with $r > 1$: These balls are almost infinitely larger than the unit ball

Example of CD: Covering Number of Balls

- The above CD happens not just for balls, but for many other objects in high-dimensions. In fact, for almost all high-dimensional objects.
 - This means that it is very likely we to work with high-dimensional data that suffers from the CD, since almost all high-dimensional objects suffer from the CD.
- There are only few high-dimensional objects that do not suffer from the CD. One of these is the polytope.
 - This means that it is very unlikely we to work with high-dimensional data that does not suffers from the CD, since there are few high-dimensional objects that do not suffer from the CD.
- Then how can we avoid CD? We avoid it by not seeking to obtain exactly correct results. Instead, we seek to obtain approximately correct results.

Example of Overcoming CD - Numerical Integration Dd

- Numerical integration in dimension d : Let us have a d -dimensional function $g(x_1, x_2, \dots, x_d)$. We want to numerically compute the following integral

$$I_d = \int_0^1 \int_0^1 \dots \int_0^1 g(x_1, x_2, \dots, x_d) dx_1 dx_2 \dots dx_d$$

Example of Overcoming CD - Numerical Integration D1

- When $d = 1$, we have

$$I_1 = \int_0^1 g(x) dx$$

- To approximately compute I_1 , we pick n equally-distant points $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n$, each in D1, with distance ϵ between neighboring points, and compute

$$\frac{1}{n} \sum_{i=1}^n g(\hat{x}_i) \approx I_1$$

- Note that we have used $n = 1/\epsilon$ number of points
- Figure:

Example of Overcoming CD - Numerical Integration D2

- When $d = 2$, we have

$$I_2 = \int_0^1 \int_0^1 g(x_1, x_2) dx_1 dx_2$$

- To approximately compute I_2 , we pick n equally-distant points $\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_n$ in D2, with distance ϵ between neighboring points in each dimension, and perform

$$\frac{1}{n} \sum_{i=1}^n g(\hat{\mathbf{x}}_i) = \frac{1}{n} \sum_{i=1}^n g(\hat{x}_{1i}, \hat{x}_{2i}) \approx I_2$$

- Note that each $\hat{\mathbf{x}}_i$ represents a vector of two points given by $\hat{\mathbf{x}}_i = [\hat{x}_{1i}, \hat{x}_{2i}]$, where \hat{x}_{1i} is the point from first dimension and \hat{x}_{2i} is the point from second dimension.
- Note that $\epsilon = x_{1i+1} - x_{1i} = x_{2i+1} - x_{2i}$
- Note that we have used $n = 1/\epsilon^2$ number of points, which can be seen in the next figure.

Example of Overcoming CD - Numerical Integration D2

- Figure:

Example of Overcoming CD - Numerical Integration Dd

- When $d = d$, we have

$$I_d = \int_0^1 \int_0^1 \dots \int_0^1 g(x_1, x_2, \dots, x_d) dx_1 dx_2 \dots dx_d$$

To approximate compute I_d , we pick n equally-distant points $\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_n$ in Dd , with distance ϵ between neighboring points in each dimension, and perform

$$\frac{1}{n} \sum_{i=1}^n g(\hat{\mathbf{x}}_i) = \frac{1}{n} \sum_{i=1}^n g(\hat{x}_{1i}, \hat{x}_{2i}, \dots, \hat{x}_{di}) \approx I_d$$

- Note that each $\hat{\mathbf{x}}_i$ represents a vector of d points given by $\hat{\mathbf{x}}_i = [\hat{x}_{1i}, \hat{x}_{2i}, \dots, \hat{x}_{di}]$, where \hat{x}_{ki} is the point in the k -th dimension of $\hat{\mathbf{x}}_i$, for $k = 1, 2, \dots, d$ and $i = 1, 2, \dots, n$.
- Note that we have used $n = 1/\epsilon^d$ number of points, and this cannot be even illustrated by a figure.
- If d is large we will be out of memory \rightarrow Uncomputable for large d !

Example of Overcoming CD - Numerical Integration Dd

- Now let's try to compute the d -dimensional integral using probabilistic tools.
- Instead of picking ϵ -distant $1/\epsilon^d$ points, let's pick n random points $\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_n$, where each $\hat{\mathbf{x}}_i = [\hat{x}_{1i}, \hat{x}_{2i}, \dots, \hat{x}_{di}]$ is chosen independently and identically according to the multivariate uniform distribution

$$f_{\mathbf{X}}(\mathbf{x}) = \prod_{k=1}^d f_X(x_k),$$

where

$$f_X(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1 \\ 0 & \text{if otherwise} \end{cases}$$

- Then let's compute

$$\frac{1}{n} \sum_{i=1}^n g(\hat{\mathbf{x}}_i) = \frac{1}{n} \sum_{i=1}^n g(\hat{x}_{1i}, \hat{x}_{2i}, \dots, \hat{x}_{di})$$

- How close is the above sum to I_d ?

Example of Overcoming CD - Numerical Integration

- If we repeat the above process many times, we can consider that $\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_n$ are n independent and identically distributed (i.i.d.) random vectors.
- Then we can compute the average of the sum as

$$E \left[\frac{1}{n} \sum_{i=1}^n g(\hat{\mathbf{x}}_i) \right] = \frac{1}{n} \sum_{i=1}^n E[g(\hat{\mathbf{x}}_i)] \stackrel{(a)}{=} E[g(\hat{\mathbf{x}})],$$

where (a) follows due to the i.i.d. construction of the $\hat{\mathbf{x}}_i$'s. Next,

$$\begin{aligned} E[g(\hat{\mathbf{x}})] &= \int_0^1 \int_0^1 \dots \int_0^1 g(x_1, x_2, \dots, x_d) \prod_{k=1}^d f_X(x_k) dx_1 dx_2 \dots dx_d \\ &\stackrel{(b)}{=} \int_0^1 \int_0^1 \dots \int_0^1 g(x_1, x_2, \dots, x_d) dx_1 dx_2 \dots dx_d = I_d \end{aligned}$$

where (b) follows due to the definition of $f_X(x)$.

- We have obtained that the constructed sum is an unbiased estimator of I_d .

Example of Overcoming CD - Numerical Integration Dd

- We have obtained that the average of the constructed sum is equal to I_d . But, what is the error we are making? We can compute the mean squared error (MSE) as

$$\begin{aligned}
 E \left[\left(\frac{1}{n} \sum_{i=1}^n g(\hat{\mathbf{x}}_i) - I_d \right)^2 \right] &= \frac{1}{n^2} E \left[\left(\sum_{i=1}^n g(\hat{\mathbf{x}}_i) - I_d \right)^2 \right] \\
 &= \frac{1}{n^2} VAR \left[\sum_{i=1}^n g(\hat{\mathbf{x}}_i) \right] \stackrel{(a)}{=} \frac{1}{n^2} \sum_{i=1}^n VAR[g(\hat{\mathbf{x}}_i)] = \frac{1}{n^2} n VAR[g(\hat{\mathbf{x}})] \\
 &= \frac{1}{n} VAR[g(\hat{\mathbf{x}})], \tag{1}
 \end{aligned}$$

where (a) is due to the i.i.d.

- Hence, the MSE decays with $1/n$. This means that the absolute error, defined as $\sqrt{\text{MSE}}$, decays with $1/\sqrt{n}$
- Note, the absolute error (and MSE) is independent from the dimension d ! We have avoided the curse of dimensionality!

Example of Overcoming CD - Numerical Integration Dd

- Figure of the absolute error:

Example of Overcoming CD - Computational Geometry

Def: A set $\mathcal{T} \subset \mathbb{R}^d$ is called a convex set if $\forall x \in \mathcal{T}$ and $\forall y \in \mathcal{T}$, the point $\lambda x + (1 - \lambda)y \in \mathcal{T}$, for any $0 \leq \lambda \leq 1$.

- The above simply means that any point on the line between x and y also belongs to the set \mathcal{T} .

Example of Overcoming CD - Computational Geometry

- A way of transforming a non-convex set into a convex set is called taking The Convex Hull
- Def: The convex hull of $\mathcal{T} \subset \mathbb{R}^d$, denoted by $\text{conv}(T)$ is the smallest convex set that contains T .
- Note $\mathcal{T} \subset \text{conv}(T)$.
- Examples:
 - $\mathcal{T} = \{a, b\}$, then $(\mathcal{T}) = [a, b]$, which denotes the line segment between a
 - $\mathcal{T} = \{a, b, c\}$, then

Example of Overcoming CD - Computational Geometry

- Fact: (Convex Combination) If $\mathbf{x} \in \text{conv}(\mathcal{T})$, then $\mathbf{x} = \sum_{i=1}^m \lambda_i \mathbf{x}_i$, where $\mathbf{x}_i \in \mathcal{T}$ and $\lambda_i \geq 0$, $\forall i$, and $\sum_{i=1}^m \lambda_i = 1$.
- Caratheodori Thm: $\forall \mathbf{x} \in \text{conv}(\mathcal{T})$, where $\mathcal{T} \subset \mathbb{R}^d$, can be expressed as a convex combination of at most $d + 1$ points (read vectors) in \mathcal{T} .
- The above means that any $\mathbf{x} \in \text{conv}(\mathcal{T})$ can be described using at most $d + 1$ other vectors in \mathcal{T} , even though the set $\text{conv}(\mathcal{T})$ itself contains infinitely many vectors.
- Hence, we can store infinitely many vectors just by storing at most $d + 1$ other vectors, which is a form of lossless compression.
- Other applications are in optimization: see paper “Fast and Accurate Least-Mean-Squares”

Example of Overcoming CD - Computational Geometry

Ex: $\mathcal{T} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\}$, where $\forall \mathbf{x}_i \in \mathbb{R}^d$. Then $\text{conv}(\mathcal{T})$ is rectangle with ∞ points.

Example of Overcoming CD - Computational Geometry

- What if we want to describe/store all points in $\text{conv}(\mathcal{T})$ with less than $d + 1$ points/vectors, in order to save storage space, i.e., we want to do lossy compression?
- What kind of an error we would make in that compression process?
- In fact, it turns out that we can describe \mathbf{x} with an error by choosing $m < d$, and m is independent of d . Avoided the CD!
- This is described in the following

Example of Overcoming CD - Computational Geometry

- Approximate Caratheodori Th (ACT): Let $\mathcal{T} \subset \mathbb{R}^d$, and let $\text{diam}(\mathcal{T}) \leq 1$. Then, $\forall \mathbf{x} \in \text{conv}(\mathcal{T})$ and $\forall k \in \mathbb{N}$, $\exists \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ such that

$$\left\| \mathbf{x} - \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i \right\|_2 \leq \frac{1}{\sqrt{2k}}$$

- $\text{diam}(\mathcal{T}) \leq 1$ means that the length between any two points in \mathcal{T} is at most one. This is just scaling, since we can always divide the points in \mathcal{T} by a scalar such that $\text{diam}(\mathcal{T}) \leq 1$ holds.
 - Note the equal weights above, i.e., $\lambda_i = 1/k$, $\forall i$.
 - Note $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d x_i^2}$, where x_i is the i -th element of \mathbf{x} .
- Compression:** Hence, instead of storing the infinite number of vectors in $\text{conv}(\mathcal{T})$ using at most $d + 1$ vectors in \mathcal{T} , I can store $\text{conv}(\mathcal{T})$ using $k \ll d$ vectors in \mathcal{T} , and the squared error that I would make is at most $1/\sqrt{2k}$.

Example of Overcoming CD - Computational Geometry

- How can we achieve the (ACT)? Using probability, and this is how:
- First, another definition of $VAR[X]$:

$$VAR[X] = \frac{1}{2}E[(X - \hat{X})^2],$$

where \hat{X} is i.i.d. as X , and X is one dimensional. Prove this at home!

- Similarly, for random vectors: If $\mathbf{X} \in \mathbb{R}^d$, then

$$VAR[\mathbf{X}] = E[\|\mathbf{X} - E[\mathbf{X}]\|_2^2] = \frac{1}{2}E[\|\mathbf{X} - \hat{\mathbf{X}}\|_2^2], \quad (2)$$

where $\hat{\mathbf{X}}$ is i.i.d. as \mathbf{X} . Again, prove at home!

Example of Overcoming CD - Computational Geometry

- Proof of ACT: Fix $\mathbf{x} \in \text{conv}(\mathcal{T})$. Then, we know that

$$\mathbf{x} = \sum_{i=1}^{d+1} \lambda_i \mathbf{z}_i,$$

where $\mathbf{z}_i \in \mathcal{T}$, $\lambda_i \geq 0$ and $\sum_{i=1}^{d+1} \lambda_i = 1$.

We now interpret λ_i 's as probabilities and the above sum as expectation.

Let \mathbf{Z} be a discrete random vector with PMF $\Pr\{\mathbf{Z} = \mathbf{z}_i\} = \lambda_i$. Then,

$$\mathbf{x} = E[\mathbf{Z}] = \sum_{i=1}^{d+1} \lambda_i \mathbf{z}_i \tag{3}$$

Example of Overcoming CD - Computational Geometry

Now let us take k RVs that are i.i.d. as \mathbf{Z} , given by $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_k$.

Then, we know that

$$\begin{aligned}
 E \left[\left\| \frac{1}{k} \sum_{i=1}^k \mathbf{Z}_i - E[\mathbf{Z}] \right\|_2^2 \right] &\stackrel{(a)}{=} E \left[\left\| \frac{1}{k} \sum_{i=1}^k \mathbf{Z}_i - \mathbf{x} \right\|_2^2 \right] \\
 &= E \left[\left\| \frac{1}{k} \sum_{i=1}^k (\mathbf{Z}_i - \mathbf{x}) \right\|_2^2 \right] \stackrel{(a)}{=} \frac{1}{k^2} E \left[\left\| \sum_{i=1}^k (\mathbf{Z}_i - E[\mathbf{Z}_i]) \right\|_2^2 \right] \\
 &= \frac{1}{k^2} VAR \left[\sum_{i=1}^k \mathbf{Z}_i \right] = \frac{1}{k^2} \sum_{i=1}^k VAR[\mathbf{Z}_i] = \frac{1}{k} VAR[\mathbf{Z}] \\
 &\stackrel{(b)}{=} \frac{1}{2k} E [\|\mathbf{Z} - \mathbf{Z}_1\|_2^2] \stackrel{(c)}{\leq} \frac{1}{2k} \text{diam}(\mathcal{T})^2 \stackrel{(d)}{\leq} \frac{1}{2k},
 \end{aligned}$$

where (a) is due to (3), (b) is due to (2), and c is due to the fact that the difference between any two points cannot be larger than $\text{diam}(\mathcal{T})$, and (d) is due to $\text{diam}(\mathcal{T}) \leq 1$

Example of Overcoming CD - Computational Geometry

Now , since

$$E \left[\left\| \frac{1}{k} \sum_{i=1}^k \mathbf{Z}_i - \mathbf{x} \right\|_2^2 \right] \leq \frac{1}{2k} \quad (4)$$

then $\exists \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k \in \mathcal{T}$ (read there exist k realizations of \mathbf{Z}) such that

$$\left\| \frac{1}{k} \sum_{i=1}^k \mathbf{z}_i - \mathbf{x} \right\|_2^2 \leq \frac{1}{2k} \quad (5)$$

from which

$$\left\| \mathbf{x} - \frac{1}{k} \sum_{i=1}^k \mathbf{z}_i \right\|_2 \leq \frac{1}{\sqrt{2k}} \quad (6)$$

Example of ACT

- We are given n images represented as vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, where $\mathbf{x}_i \in \mathbb{R}^d, \forall i$. Let x_{ji} be the j -th element of \mathbf{x}_i , where x_{ji} represents the j -th pixel of the image \mathbf{x}_i .
- We would like to create a new image \mathbf{x} (of some desired shape), by mixing some or all of the given images $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, i.e., by convex combination of the given images. Mathematically,

$$\mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{x}_i,$$

where $\lambda_i \geq 0$ and $\sum_{i=1}^n \lambda_i = 1$.

- We must assume that the desired image \mathbf{x} can indeed be formed as a convex combination of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. More precisely, \mathbf{x} belongs to the convex hull created from $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, i.e., $\mathbf{x} \in \text{conv}(\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\})$. Otherwise, the desired image cannot be created.

Example of ACT

- What the Caratheodori Theorem tells us, is that we need to mix at most $d + 1$ images to create the desire image \mathbf{x} . If the image is HD, thereby contains, $d = 1440 \times 1080 = 1,555,200$ pixels, we need to mix $1,555,200 + 1$ images to obtain the desire image.
- What the Approximate Caratheodori Theorem tells us, is that we need to mix at most k images to create the desire image \mathbf{x} . If we mix $k = 100$ images, then the absolute mean error between the original and the approximated image will be $1/(2 * 10) = 0.05 = 5\%$, which might be a very close approximation that the eye cannot notice.
- Note: The Approximate Caratheodori Theorem shows only existence of the points that make the approximations, but do not shows us how to obtain these points, other by brute-force search. Hence, the search for such an algorithm is ongoing.