

CURSE OF DIMENSIONALITY

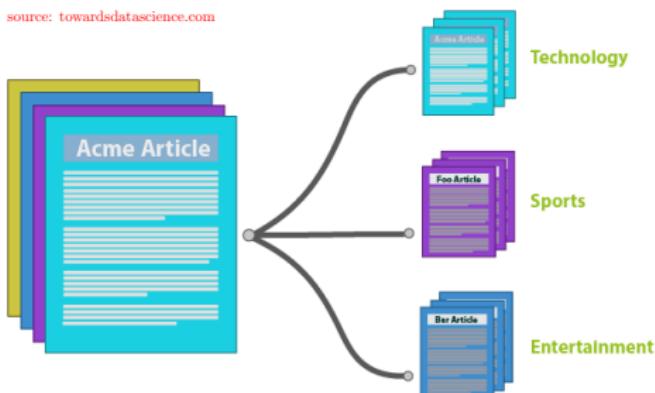
- Imdadullah Khan

High Dimensional Data

High Dimensional Data is common in many applications

- Text represented as set or bag or TF-IDF of words
- 1000's of unigram, millions of bigrams plus contextual attributes

source: [towardsdatascience.com](https://towardsdatascience.com/introduction-to-high-dimensional-data-1f3a2a2a2a)



Three text snippets are shown, each with a downward arrow pointing to its corresponding binary word vector:

Text Snippet	Binary Vector
The elephant sneezed at the sight of potatoes.	1 0 0 0 0 0 1 1 0 1 0 0 1 1 0 1 0 0 0
Bats can see via echolocation. See the bat sight sneeze!	0 1 0 1 0 0 0 0 1 0 1 1 0 1 0 1 0 0 0
Wondering, she opened the door to the studio.	1 0 1 0 0 0 1 1 0 1 0 1 0 1 0 0 0 0 0

Bengfort, Miller & Odrzy - Applied Text Analysis with Python

Three text snippets are shown, each with a downward arrow pointing to its corresponding TF-IDF word vector:

Text Snippet	TF-IDF Vector
The elephant sneezed at the sight of potatoes.	0 2 1 0 1 0 0 0 0 2 0 1 1 0 1 0 1 0
Bats can see via echolocation. See the bat sight sneeze!	0 1 0 1 0 0 0 0 1 0 1 1 0 1 0 1 0 0
Wondering, she opened the door to the studio.	0 0 1 0 0 0 1 1 0 0 1 0 1 0 1 0 0 0

Bengfort, Miller & Odrzy - Applied Text Analysis with Python

Three text snippets are shown, each with a downward arrow pointing to its corresponding TF-IDF word vector:

Text Snippet	TF-IDF Vector
The elephant sneezed at the sight of potatoes.	0 0 0 0.3 0 0 0 0.3 0 0 0.3 0 0 0.4 0 0 0 0.3
Bats can see via echolocation. See the bat sight sneeze!	0 1 0 0 0 0 0 1 0 0 0 0 0.4 0 0 0 0 0
Wondering, she opened the door to the studio.	1 0 1 0 0 0 1 1 0 1 0 1 0 1 0 0 0 0 0

Bengfort, Miller & Odrzy - Applied Text Analysis with Python

High Dimensional Data

High Dimensional Data is common in many applications

- Utility matrix for recommenders (Amazon product catalogue)
- The netflix prize training set: $\sim 1M$ ratings of the form $\langle \text{user}, \text{movie}, \text{date of grade}, \text{grade} \rangle$

- 480,189 users, 17,770 movies

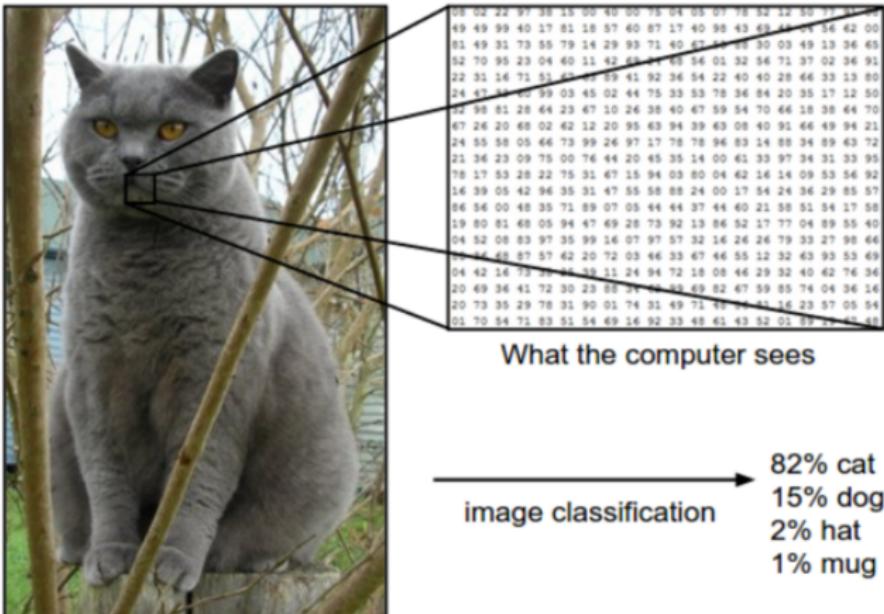


	p_1	p_2	p_3	p_j				p_m			
u_1	1	2	1	4		2	3	2	5		2
u_2		1			2	1		2		1	
u_3	1	1	2			1				1	2
		3		2		5		2		3	4
	1		2						5		
u_i		3	2	1	4	5	?	1	3	1	2
		4								4	
	5			1						5	
	1		4				1	3	5	1	2
u_n		3		1	1	2	1		4		5

High Dimensional Data

High Dimensional Data is common in many applications

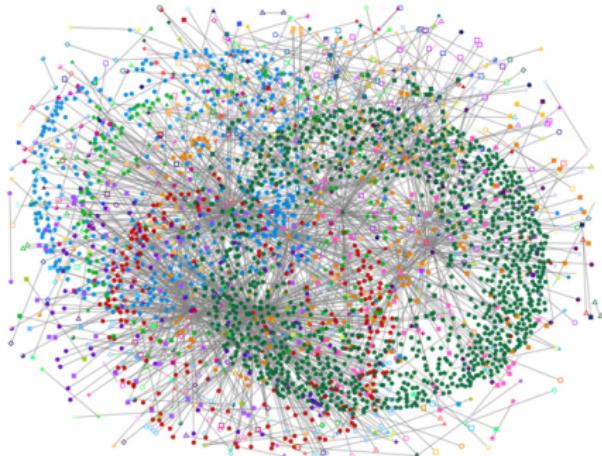
- Images and videos from multi-mega pixels digital cameras



High Dimensional Data

High Dimensional Data is common in many applications

- Social networks as adjacency matrix
- A row of facebook graph's adjacency matrix contains more than a billion dimensions



Curse of dimensionality

Richard Bellman coined the phrase, referring to difficulty of dynamic optimization with many variables

Broadly the following issues are faced when working with high dimensional data

- Computational challenging, processing, storing, communication
- In general as number of features increases redundancy also increases
 - More noise added to data than signal
- Hard to visualize and interpret

Proximity Problems

Given a set X of m -dim vectors, with $|X| = n$

Two generic proximity computation problems are building blocks of almost all data analytics

1 Distance Matrix Computation

- Find $n \times n$ matrix with all pairwise distances

2 k -nearest neighbors

- Given a query point q in the same space as X , return the k closest points in X to q

Proximity Problems: Applications

Given a set X of m -dim vectors, with $|X| = n$

Distance Matrix Computation

- Find $n \times n$ matrix with all pairwise distances

Near-duplicates detection

- Find all pairs of points with distance less than δ , or all pairs with distance less than 2σ from the mean distance
- Find mirror webpages, News Aggregation
- Plagiarism Detection

The distance matrix is input for

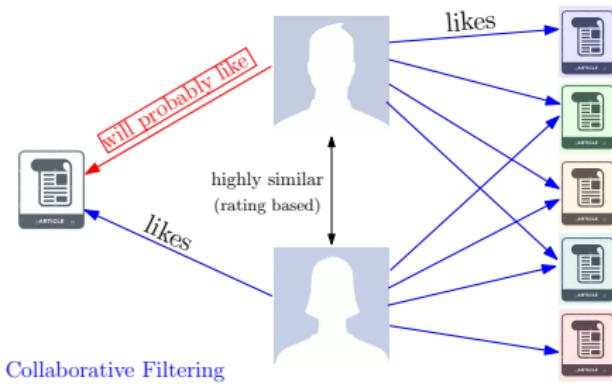
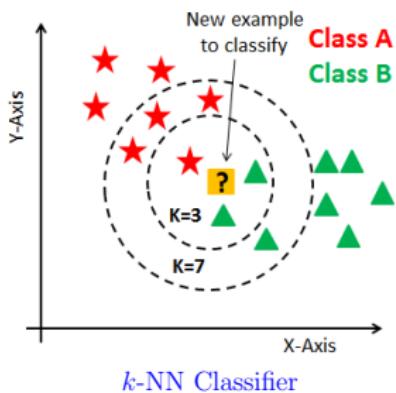
- Agglomerative clustering
- Principal Component Analysis
- Spectral Clustering
- Multi-dimensional Scaling

Proximity Problems: Applications

Given a set X of m -dim vectors, with $|X| = n$

k -nearest neighbors (k -NN) problem

- Given a query point q in the same space as X , return the k closest points in X to q



Proximity Problems: Fixed Radius Nearest Neighbors

Given a set X of m -dim vectors, with $|X| = n$

k -nearest neighbors (k -NN) problem

- Given a query point q in the same space as X , return the k closest points in X to q

A variant of the k -NN problem is

Fixed radius nearest neighbors problem:

- Given a query point q in the same space as X and a radius $r > 0$, find all points in X to within radius r from q

This variant is the same as the k -NN problem, in the sense that they are reducible to each other

Curse of Dimensionality: Computational Complexity

Given a set X of m -dim vectors, with $|X| = n$

Almost all $d(x, y)$ measures require traversal of all coordinates of x and y

Runtime of the brute force algorithms for D matrix computation

$$O(n^2 \times m)$$

Runtime of the brute force algorithms for $k\text{-NN}(q)$ is

$$O(n \times m)$$

Both runtimes grow linearly with dimensionality

Data Sparsity

As dimensionality increases the relative input space covered by a fixed-size training set diminishes

Many methods require a sizeable number of examples/samples in every region of the space to support a hypothesis or train a generalizable model

1000 students (discretized) scores in course $\in \{0, 25, 50, 75, 100\}\%$

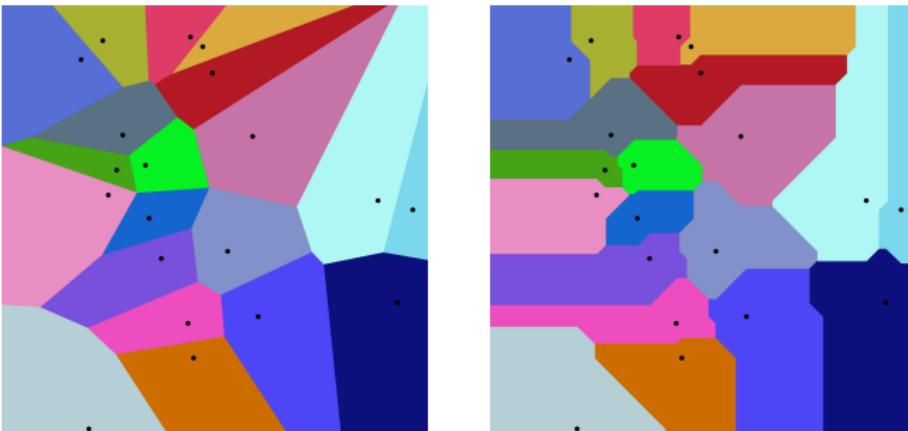
- Two courses c_1 and $c_2 \rightarrow 5 \times 5$ grades combination
 - Each combination has average $1000/25 = 40$ students
 - Good enough sample size, can infer rules like
 - if $grade(c_1) \leq 50 \wedge grade(c_2) \geq 75$, then student is Math major
- For four course, number of grade combinations is $5^4 = 625$
 - 1.6 students per combination
- For 10 course, average students per combination is 0.0001024
 - Almost all combinations are never observed

Approaches for nearest neighbor problem

- Store X in a list
- No preprocessing
- On query run a FINDMIN algorithm on distance to q
- Runtime is $O(n)$ distance computations
- For $m = 1$, store X in a sorted array
- Best data structure for 1-d $NN(q)$
- With Binary search for q runtime is $O(\log n)$ distance computations

Approaches for nearest neighbor problem

- Voronoi diagram ($m = 2$) Partition of plane into nearest neighbor regions
- Region R_i of a point $x_i \in X$ is the set of all points that are NN of x_i
- R_i : intersection of perp. bisectors of segments b/w x_i and other points
- For $m = 2$, Fortune's algorithm for voronoi diagram in $O(n \log n)$

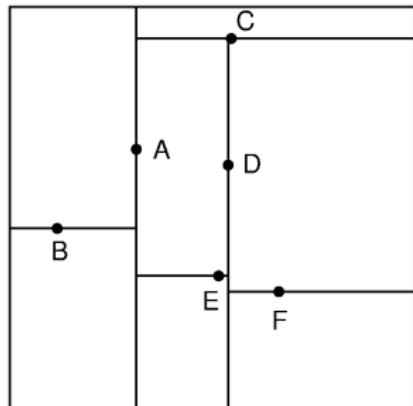


Voronoi diagrams of 20 points under (left) Euclidean and (right) Manhattan distance. source: Wikipedia

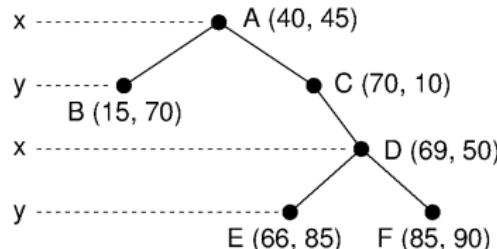
Hard to even describe in higher dimensions

Approaches for nearest neighbor problem

- *kd-tree* data structure: Partition the space into non-uniform cells
- A binary tree where each level compare 1 dimension (cutting dimension)
- Internal nodes correspond to hyperplanes splitting space in 2 half spaces
- Halve the points by a hyperplane perpendicular to one dimension
- Recursively construct *kd-tree* for the two halves, until one point remains
- Cycle through all dimensions



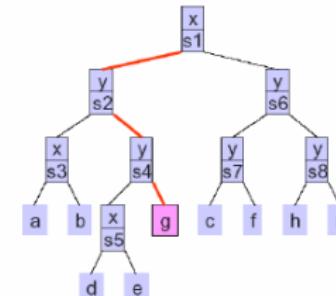
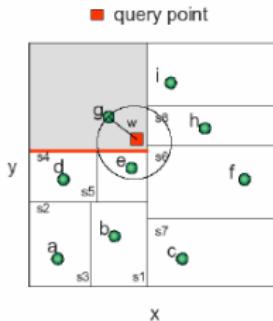
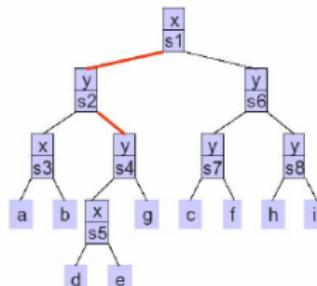
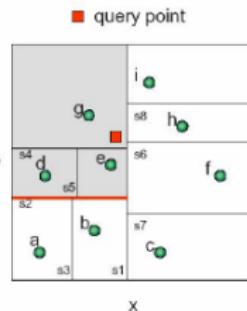
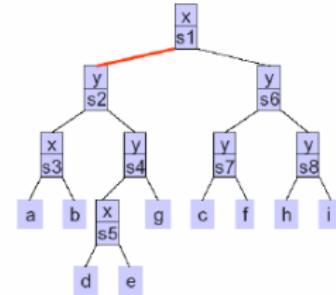
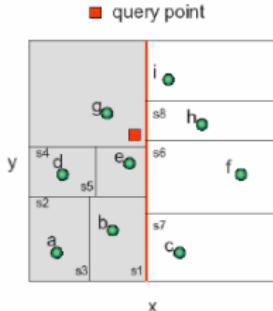
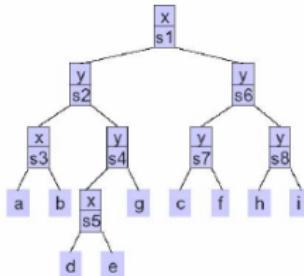
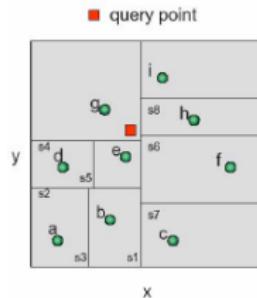
(a)



(b)

Approaches for nearest neighbor problem

Searching for nearest neighbor in kd -tree

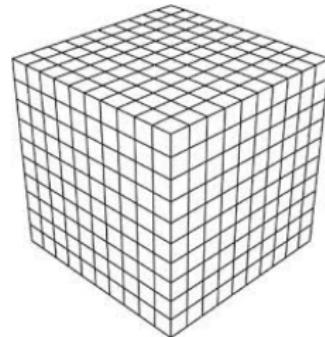


kd -trees are very effective for dimensions ≤ 10 or so

T. Nguyen @ Oregon State

Huge Search Space for Nearest neighbor

- For large dimensions partition the space into cells (grids or mesh)
- NN search performed in the cell containing q and ‘neighboring’ cells
- Number of ‘neighboring’ cells in 2-d is $3^2 = 9$, in 3-d 3^3 , in m -d, 3^m



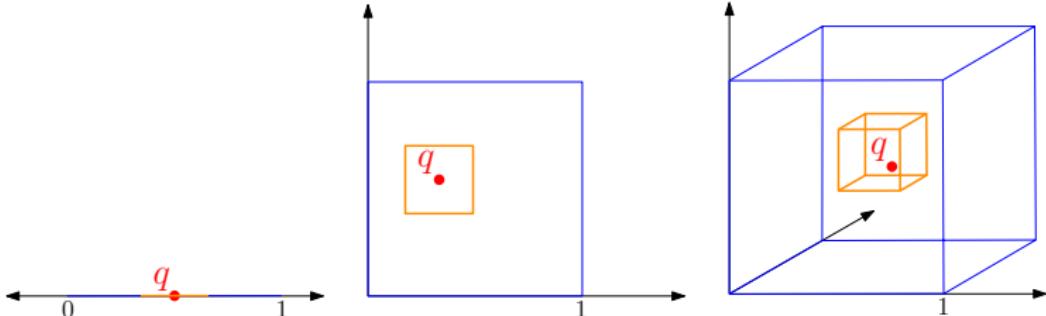
Grid can be non-uniform as in kd -tree

Huge Search Space for Nearest neighbor

Another way to look at this

Higher dimensional neighborhood is very large and not local. The notion of nearest neighbor breaks down

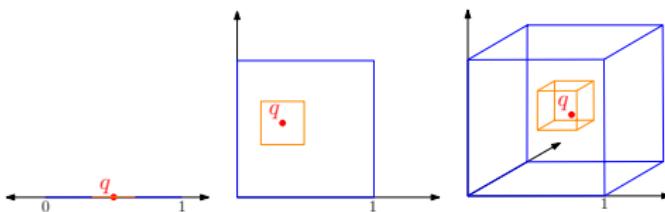
- Suppose n points are placed uniformly at random in $[0, 1]^m$
- Grow a hypercube around q to contain f fraction of points ($k = fn$)
- Expected edge length: $E_m(f) = (f^{1/m})$
- In $10d$ to get 10% points around q need cube with edge length 0.8
- To get only 1% point need to extend cube 0.63 along each dimension



Huge Search Space for Nearest neighbor

Another way to look at non-locality of higher dimensional neighborhoods

- Suppose 5000 points are randomly placed in $[0, 1]^m$. Let $q = \mathbf{0}$
- In 1d must go a distance $\sqrt[5]{5000} = 0.001$ on average to capture 5 NN
- In 2d must go $\sqrt[5]{5000} = 0.031$ units along both dimensions
- In 3d must go $\sqrt[5]{0.001} = 0.1 = 10\%$ of the total (unit) length
- In 4d must go $\sqrt[4]{0.001} = 0.177 = 17.7\%$ of unit length
- In 10d must go 50.1% of unit length along each dimension
- In md must go $(5/5000)^{1/m}$ along each dimension



In high dimensional space nobody can hear you scream

Diminishing Volume of m -ball

A manifestation of this phenomenon that points in higher dimensions are isolated is **the diminishing relative volume of the m -ball in m -cube**

The m -ball (m -d hypersphere) of radius r centered at origin

$$B_{m,r} := \{ \mathbf{x} \in \mathbb{R}^m : d(\mathbf{x}, \mathbf{0}) \leq r \} \implies \|\mathbf{x}\|_2 \leq r$$

Volume of $B_{m,r}$:

$$V_m(r) = \frac{\pi^{m/2}}{\Gamma(m/2 + 1)} r^m$$

$\Gamma(\cdot)$ essentially is factorial of fractional numbers

$$V_m(r) = \frac{\pi^{m/2}}{m/2!} r^m \quad \text{For simplicity assume } m \text{ is even}$$

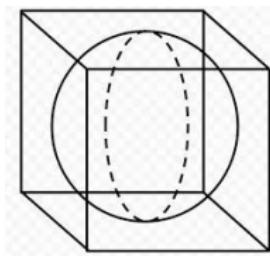
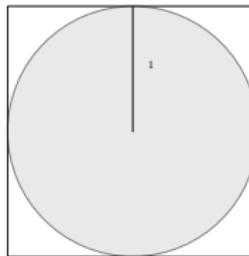
The m -cube (m -d hypercube) is the set $[-1, 1]^m$ (note edge length is 2)

Volume of m -cube: 2^m

Diminishing Volume of m -ball

In m -d ratio of volume of unit m -Ball to that of m -cube (edge length 2)

$\frac{\pi^{m/2}/m/2!}{2^m}$ approaches 0 very fast

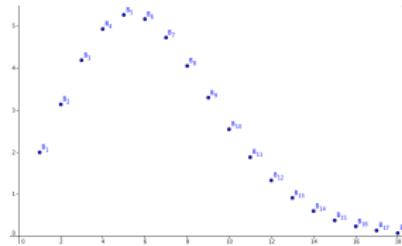


dim m	volume of m -ball	volume of m -cube	ratio
2	π	2^2	~ 0.785
3	$4/3\pi$	2^3	~ 0.523
4	$\pi^2/2$	2^4	~ 0.308
6	$\pi^3/6$	2^6	~ 0.080
m	$\frac{\pi^{m/2}}{m/2!}$	2^m	$\rightarrow 0$

Diminishing Volume of m -ball

Ratio of volumes of unit m -Ball and $[-1, 1]^m$

$$\frac{\pi^{m/2}/m/2!}{2^m}$$



- In higher dimensions all the volume is in ‘corners’
- Points in high dimensional spaces are isolated (empty surrounding)
- The probability that a randomly generated point is within r radius of q approaches 0 as dimensionality increases
- **The probability of a close nearest neighbor in a data set is very small**
- Caveat: Real datasets are not random
- Overcome this by getting larger training set (exponential in m)

Diminishing Volume of m -ball

ratio of volumes of unit m -Ball and $[-1, 1]^m$

$$\frac{\pi^{m/2}/m!/2^m}{2^m}$$



- In higher dimensions all the volume is in 'corners'
- Probability of a close nearest neighbor in random data set is very small
- Overcome this by getting larger training set (exponential in m)

To cover $[-1, 1]^m$ with $B_{m,1}$'s, the number of balls n must be

$$n \geq \frac{2^m}{V_m(1)} = \frac{2^m}{\pi^{m/2}/m!/2^m} = \frac{m!/2^m}{\pi^{m/2}} \quad m \rightarrow \infty \quad \sqrt{m\pi} \left(\frac{m2^{m/2}}{2\pi e} \right)^{m/2}$$

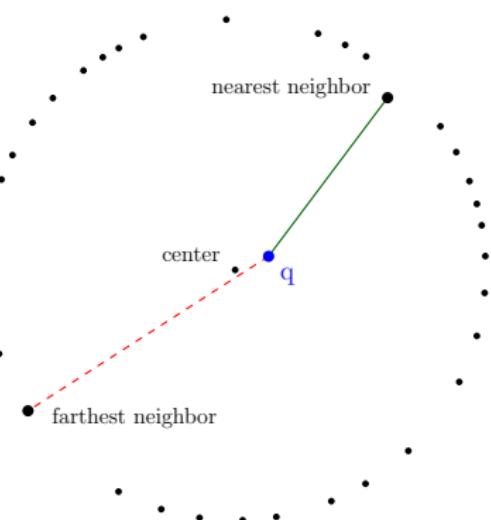
For $m = 16$ (a very small number) this n is substantially bigger than 2^{58}

Instability of Nearest neighbor

In higher dimension the notion of nearest neighbor breaks down

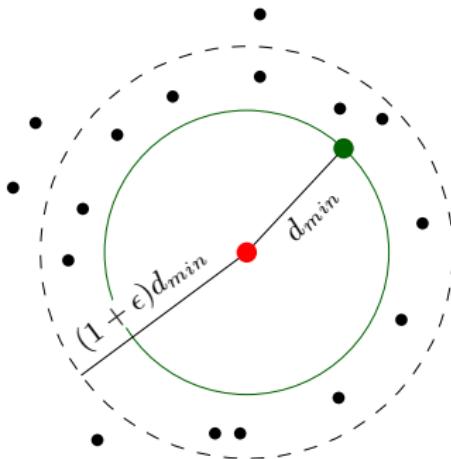
No difference (contrast) between nearest and farthest neighbors

A points nearest neighbor loses meaning



Instability of Nearest neighbor

A nearest neighbor query is ϵ -unstable ($\epsilon > 0$), if the distance from q and most other points are at most $(1 + \epsilon)$ times the distance from q to its 1NN



We show that as dimensionality increases the probability of all nearest neighbors queries becoming unstable increase (distance concentration)

Distance Concentration

Another facet of curse of dimensionality is the phenomenon of distance concentration

Assume points in \mathbb{R}^m and ℓ_2 distance measure

- As m increases, almost all pairs of points have their ℓ_2 distances
 - similar to distance of other pairs and
 - and very high
- normalized distance is close to 1 (both high and similar are encompassed)

We demonstrate it by observing distribution of pairwise distances for n points in \mathbb{R}^m (again real-life datasets are not random...)

Distance Concentration

Another facet of curse of dimensionality is the phenomenon of distance concentration

All pairwise distances are very high

Consequences:

- Distance measure loses its meaning
- We discussed it earlier that proximity measure is the building block of data analytics, when it becomes meaningless the building collapses
- Nearest neighbor is as good as farthest neighbor
 - e.g. in such cases very hard to build clusters,
 - no justification to group a pair of points and not another

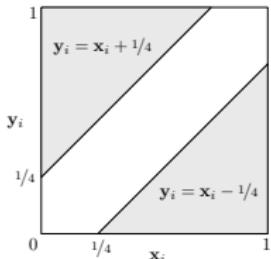
Distance Concentration Analytical Bounds

- Generate a set \mathcal{X} of n points at random in $[0, 1]^m$
- Maximum possible distance b/w a pair $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ is $d(\mathbf{x}, \mathbf{y}) \leq \sqrt{m}$
- Consider the squared- ℓ_2 distance (for convenience)
- $d^2(\mathbf{x}, \mathbf{y}) := \|\mathbf{x} - \mathbf{y}\|^2 \leq m$

Distance Concentration Analytical Bounds

Generate a set \mathcal{X} of n points at random in $[0, 1]^m$

For a fixed coordinate $i < m$, $\Pr[|\mathbf{x}_i - \mathbf{y}_i| \geq 1/4] > 1/2$



Let $V_i = \begin{cases} 1 & \text{if } |\mathbf{x}_i - \mathbf{y}_i| \geq 1/4 \\ 0 & \text{else} \end{cases}$ \triangleright Indicator if coordinate difference is big

Let $V = \sum_{i=1}^m V_i = |\{i : |\mathbf{x}_i - \mathbf{y}_i| \geq 1/4\}|$

$E(V) \geq m/2$ \triangleright linearity of expectation

On average at least half coordinates differences are $\geq 1/4$ ('big difference')

Distance Concentration Analytical Bounds

Theorem (Chernoff Bound (tail inequality))

Let $V = V_1 + V_2 + \dots + V_m$ be the sum of m independent Bernoulli random variables and let $E(V) = \mu$. The (loose) Chernoff bounds are:

- $\Pr(V \geq (1 + \delta)\mu) \leq e^{-\delta^2\mu/3}$ for $0 < \delta < 1$
- $\Pr(V \geq (1 + \delta)\mu) \leq e^{-\delta\mu/3}$ for $\delta > 1$
- $\Pr(V \leq (1 - \delta)\mu) \leq e^{-\delta^2\mu/2}$ for $0 < \delta < 1$

For fixed \mathbf{x}, \mathbf{y} $[V \geq \frac{m}{4} \implies \|\mathbf{x} - \mathbf{y}\|^2 \geq \frac{m}{64}]$ w.p $\geq 1 - e^{-\frac{m}{16}}$ $\triangleright \delta = \frac{1}{2}$

From this using union bound we get the following result

If $m = \Omega(\log n)$, then w.h.p for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ we have $d^2(\mathbf{x}, \mathbf{y}) \geq \frac{m}{64}$

This means all pairs are far ($\text{dist} \geq \sqrt{m}/8$)

Distance Concentration Simulation

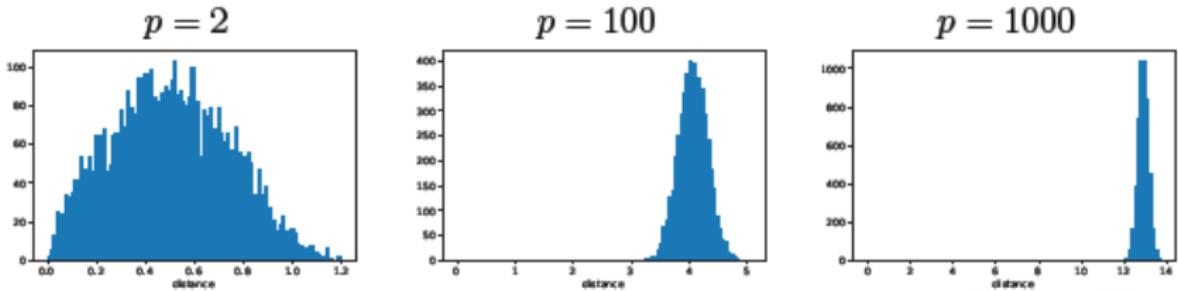


Figure: Histograms of pairwise-distances between $n = 100$ points sampled uniformly in the hypercube $[0, 1]^p$

Julie Delon @ Uni. Paris Descartes

Angle Concentration

- In large dimensions (at least for random points) the distance measure (at least ℓ_2 distance) is more or less meaningless
 - Can we use cosine distance?
 - The same concentration phenomenon is observed for pairwise angles
-
- Max num of pairwise orthogonal vectors ($\mathbf{x} \cdot \mathbf{y} = 0, \theta_{x,y} = 90^\circ$) in \mathbb{R}^2 is 2
 - Max num of pairwise orthogonal vectors in \mathbb{R}^3 is 3
 - Max number of pairwise almost orthogonal vectors
 $(\mathbf{x} \cdot \mathbf{y} \leq \epsilon, \theta_{x,y} = 90^\circ \pm \epsilon)$ is $e^{\Omega(m)}$

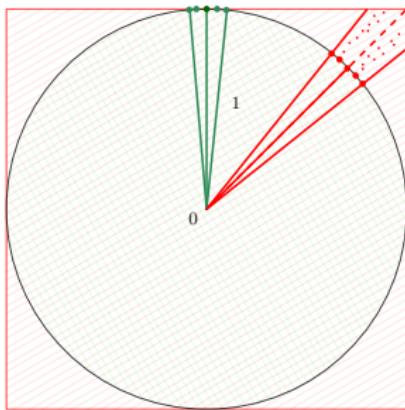
Angle Concentration: Random Direction

Generating a random direction in \mathbb{R}^m

- Equivalently a random unit vector in \mathbb{R}^m
- We will need it in subsequent sessions
- It is not a straight-forward task in higher dimensions

An immediate way to pick a random unit vector:

choose a random point in $\mathbf{v} \in [-1, 1]^m$ and normalize it as $\hat{\mathbf{v}} = \mathbf{v}/\|\mathbf{v}\|$



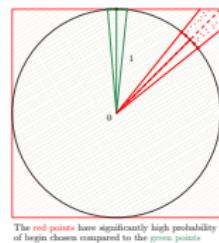
The red points have significantly high probability of being chosen compared to the green points

Clearly the distribution is skewed towards the diagonal directions

Angle Concentration: Random Direction

Generating a random direction in \mathbb{R}^m

- choose a random point in $\mathbf{v} \in [-1, 1]^m$
- normalize it as $\hat{\mathbf{v}} = \mathbf{v}/\|\mathbf{v}\|$
- distribution skewed towards diagonal directions

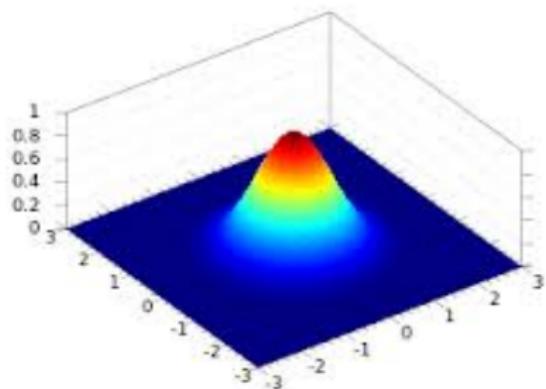
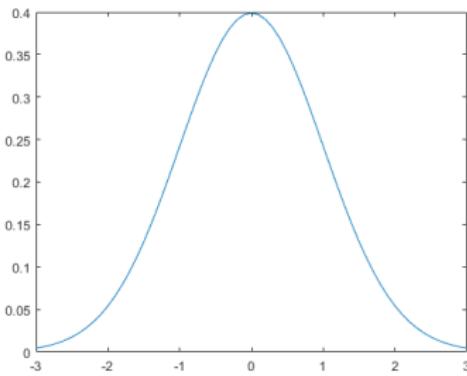


A quick fix, due to Marsaglia & Zaman

- Generate $\mathbf{v} \in [-1, 1]^m$
- If \mathbf{v} is outside the unit hypersphere ($v_1^2 + v_2^2 + \dots + v_m^2 > 1$) discard it
- Normalize any non-discarded \mathbf{v}
 - we get a point on the surface of the unit-ball equally likely
- Computationally expansive ▷ diminishing volume of unit ball
- Just in 2d choose a random number in $[0, \pi]$ and make a unit vector

Angle Concentration: Random Direction

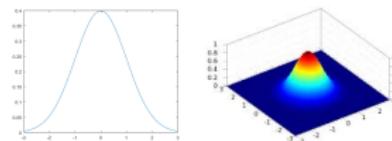
Generating a random direction in \mathbb{R}^m



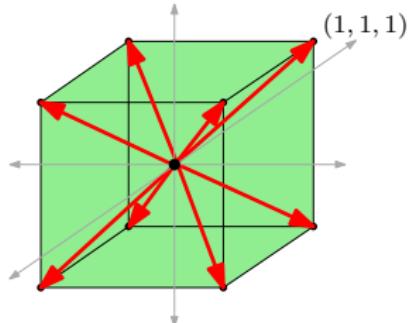
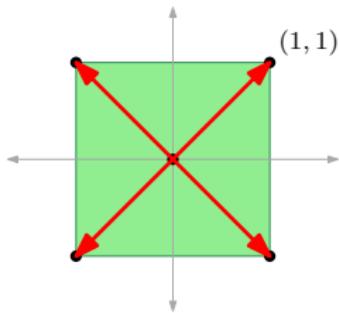
- Use spherical symmetry of the standard normal distribution
- Pick each coordinate v_i independently from $\mathcal{N}(0, 1)$ and normalize v
- Known to be uniformly distributed over the surface of the unit m -ball

Angle Concentration: Approximate Random Direction

Generating a random direction in \mathbb{R}^m



- Approximately generate unit directions
 - generate directions towards corners of the m -cubes $[-1, 1]^m$
- For $m \gg 1$, these 2^m directions approximately cover surface of m -ball
- Achlioptas (2003), Database-friendly random projections: ...



Angle Concentration: Analytical Bounds

Generate a set \mathcal{X} of n vectors in $[-1, 1]^m$ ▷ and normalize them

\mathbf{x} and \mathbf{y} are orthogonal if $\cos \theta_{\mathbf{x}, \mathbf{y}} = \langle \mathbf{x}, \mathbf{y} \rangle = \sum_i \mathbf{x}_i \mathbf{y}_i \sim 0$

For a fixed \mathbf{x} , let $V_i = \mathbf{x}_i \mathbf{y}_i$ and let $V = \sum_{i=1}^m V_i = \cos \theta_{\mathbf{x}, \mathbf{y}}$

$$\frac{-\mathbf{x}_i}{m} \leq V_i \leq \frac{\mathbf{x}_i}{m} \quad \text{and} \quad E(V_i) = 0$$

On average the vector \mathbf{x} is orthogonal to any vector \mathbf{y}

Distance Concentration Analytical Bounds

Theorem (Hoeffding's Inequality)

If X_i 's are random variables bounded by the interval $[a_i, b_i]$. Let $S = \sum_{i=1}^m X_i$. Then

$$\Pr(|S - E[S]| \geq t) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

Using this we get that

$$\Pr(V \geq \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^m (2x_i/m)^2}\right) = 2e^{-\epsilon^2 n/2}$$

From this using union bound we get the following result

If $m = \Omega(\log n)$, then w.h.p for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ we have $\cos \theta_{\mathbf{x}, \mathbf{y}} \leq \epsilon$

This means all pairs are almost orthogonal (angle $\leq \arccos(\epsilon)$)

Curse of dimensionality: Summary

Issues with Higher Dimensional Data

- Computational and Storage Challenges
 - Complexity of exact algorithms for proximity computation problems
- Data Sparsity (Sparse training set generalization)
- Issues for Nearest Neighbors
 - Huge Search Space
 - Diminishing volume of n -ball
 - Stability of nearest neighbors
- Distance Concentration
- Angle Concentration