

# **Mathematics of Bigdata Analysis: An Introduction**

**S. Lakshmivarahan  
School of Computer Science  
University of Oklahoma  
Norman, OK-73019  
varahan@ou.edu**

# Abundance of data

- Thanks to the advances in technology of
  - Sensors
  - Wireless Communication
  - Mass storage devices
  - Large super computers
- Shift from data sparse to data rich regime – amount of data doubles in every few years.

# Data organization

- **Time Series** : Number of daily new covid infection in a city.
- **Spatial**: Number of infected in every country on a given time.
- **Spatial-temporal**: Monthly rain fall in each of 50 states in the US.
- **Data Matrix**:  $X: [x_1, x_2, x_3 \dots x_n]$ ,  $x_i \in \mathbb{R}^d$ , - Represents n points in d- dimensional space.

# “Big” in Bigdata

- In the matrix form  $x \in \mathbb{R}^{d \times n}$  : Two variables.
- $n$  - is the number of data (columns).
- $d$  - is the dimension of the space (rows)
- In general: either  $n$  or  $d$  or both can be large.
- Similar measures apply for other data organization.

# Classical Statistics

- In classical mathematical statistics there are a number of asymptotic results obtained by fixing  $d$  and letting the number of samples to increase without bound such that the ratio

$$\frac{d}{n} \rightarrow 0$$

- This asymptotic theory provides the basis for estimation theory.

# Examples 1

- **Law of Large Numbers (LLN):** If  $x_i$ ,  $1 \leq i \leq n$  is i.i.d sequence of random variables from, say normal distribution  $N(m, \sigma^2)$  with unknown  $m$ .
- $\bar{x}(n) = \frac{1}{n} \sum_{i=1}^n x_i$  is an unbiased estimate.
- LLN:  $\text{prob} [ |\bar{x}(n) - m| > \varepsilon ] \rightarrow 0$  or  $n \rightarrow \infty$  -----> 1
- This is called asymptotic consistency.
- Also known as **measure concentration**.

# Examples 2

- **Central Limit Theorem (CLT):**
- In addition to (1), the following stronger result hold:

$$\frac{\sqrt{n} (\bar{x}(n) - m)}{\sigma} \xrightarrow{\text{in distribution}} N(0,1) \xrightarrow{\text{(2)}}$$

- That is, centered and scaled estimate converges in distribution to a standard normal Gaussian variable.

# High dimensional data

- Consider a set of  $n = 100$  color images of a human retina with  $256 \times 256 = 65,536$  pixels in each of the three frames representing Red, Blue and Green with a total of  $d = 65,536 \times 3 = 196,608$  pixels.

- Here  $x \in \mathbb{R}^{d \times n}$  where  $d \gg n$

- In here,  $\frac{d}{n} = \alpha > 0$

-----> (3)

# Implications of $\frac{d}{n} = \alpha > 0$

- Many of the known results from classical statistics when applied to this case,  $\frac{d}{n} = \alpha > 0$  give only “suboptimal” guarantees.
- To address this challenge a new specialty is emerging.
- M.J. Wainwright (2019) High-Dimensional Statistics: A non- asymptotic viewpoint, Cambridge university Press.
- R.Vershynin (2020) High-Dimensional Probability: An Introduction with Application in Data Science Cambridge University Press.

# Curse of dimensionality

- Coined by Richard Bellman (1920 – 1984) when developing.
- R.Bellman (1952) “Theory of Dynamic Programming”, Proc of NAS, pp 716-719.
- Finding optimal solution for multistage decision process often require  $2^d$  computation.
- The popular Reinforcement Learning (RL) is based on the theory of Markov Decision Process is an example of the application of DP.

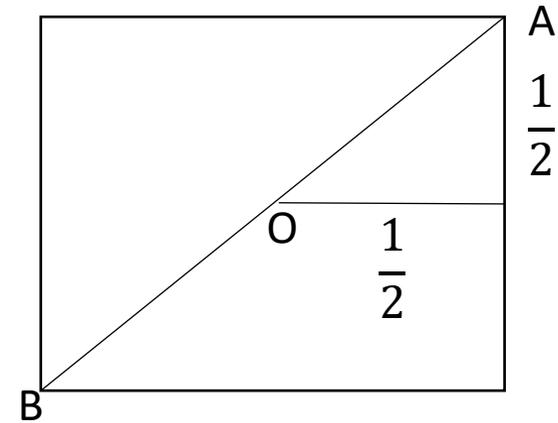
# Counter intuitive results in High dimension

- Empty space – High dimensional geometry.
- Concentration of distances, measures, functions.
- Statistical two class classification.
- Estimation of covariance matrices.

# Hyper cube $V_c(d,a)$ in $R^d$

- $V_c(d,a)$  – hypercube of side “a” in  $R^d$ .
- Diagonal AB in  $V_c(2,1)$ :

$$AB = 2 * OA = 2 \left[ \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]^{1/2} = \sqrt{2}$$



$V(2,1)$

- Diagonal AB in  $V_c(d,1)$ :

$$AB = 2 * OA = 2 \left[ \sum_{i=1}^d \left(\frac{1}{2}\right)^2 \right]^{1/2} = \sqrt{d}$$

-----> (4)

- Diagonal increases as  $\sqrt{d}$  while the side of the cube remains constraint as d increases.

# Empty space in $\mathbb{R}^d$

- Volume of  $V_c(d,a) = a^d$  .
- If we double the side :  $V_c(d, 2a) = 2^d V_c(d, a)$  -----> (5)
- Volume of the cube grows exponentially when you double its side.
- Creates a lot of empty space.

# Spheres in $\mathbb{R}^d$ : $V_s(d,r)$

- $V_s(d,r)$  – a sphere of radius  $r$  in  $\mathbb{R}^d$ .

- $\text{Vol}[V_s(d,r)] = \frac{\pi^{d/2}}{\Gamma(\frac{d}{2}+1)} r^d$  -----> (6)

- For integer  $k$ :  $\Gamma(k+1) = k \Gamma(k)$  and  $\Gamma(k+1) = k!$  -----> (7)

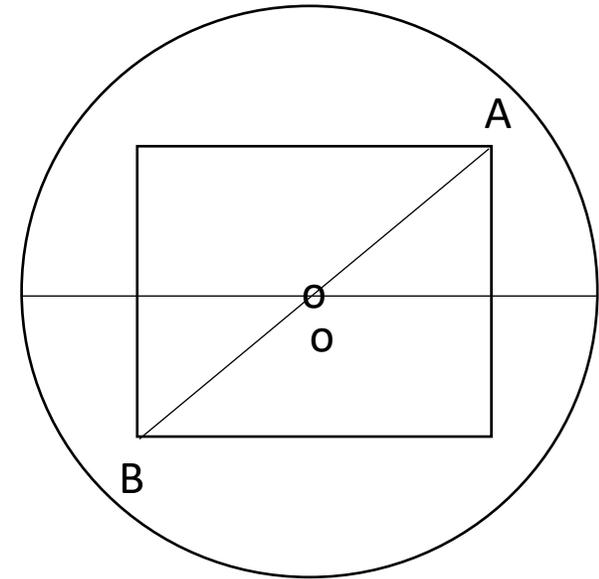
$$\Gamma(1/2) = \sqrt{\pi}$$

# Unit Sphere : $V_s(d,1)$

- $\text{Vol} [V_s(d,1)] = \frac{\pi^{d/2}}{\Gamma(\frac{d}{2}+1)}$   $\rightarrow 0$  as  $d \rightarrow \infty$
- $\text{Vol} [V_s(3,1)] = \frac{4}{3}\pi = 4.1867$   
 $\text{Vol} [V_s(10,1)] = \frac{\pi^{10}}{10!} = 0.0258$
- **Question** : For what values of  $r$ ,  $\text{Vol} [V_s(d,r)] = 1$
- Using Stirling's approximation to  $n!$  :  
$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$$
- Verify  $r = O(\sqrt{d})$  for  $\text{Vol} [V_s(d,r)] = 1$
- Empty space syndrome.

# Cube inside a cube

- Consider a unit cube inside a concentric unit sphere in  $\mathbb{R}^d$ .
- Have seen  $AB = \sqrt{d}$
- For  $d < 4$ ,  $AB < 2$  and inside the sphere.  
     $d = 4$ ,  $AB = 2$  and  $AB$  is a diameter.  
     $d > 4$ ,  $AB > 2$  and punches through the sphere.
- For large  $d$ ,  $2^d$  diagonals get out of the sphere.
- It looks like the picture of the COVID virus.



$$V_c(d,1) \subseteq V_s(d,1)$$

# Sphere in a Sphere

- Let  $r < R$ , concentric spheres of radii  $r$  and  $R$ .

- $$\frac{V_S(d,R) - V_S(d,r)}{V_S(d,R)} = 1 - \frac{V_S(d,r)}{V_S(d,R)} \quad \text{-----> (8)}$$

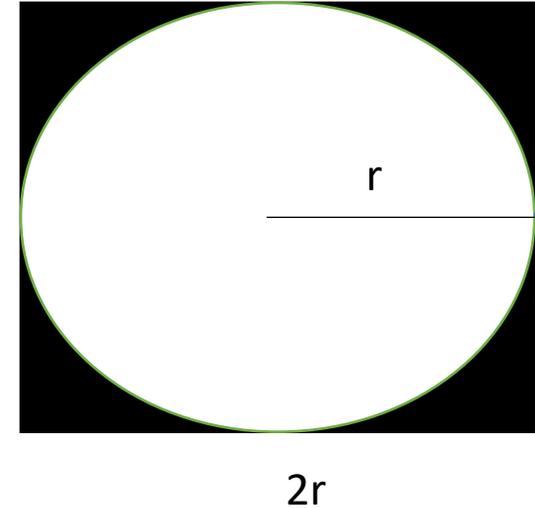
$$= 1 - \left(\frac{r}{R}\right)^d \quad \text{-> 1 as } d \text{ increases.}$$

(i.e.) Volume of the sphere reside near the empty space shell.

# Sphere in a cube

- Ratio  $\alpha = \frac{V_s(d,r)}{V_c(d,2r)}$   
$$= \frac{\pi^{d/2} \Gamma^d}{\Gamma(\frac{d}{2}+1)} \frac{1}{(2r)^d} = \left(\frac{\pi}{4}\right)^{d/2} \frac{1}{\Gamma(\frac{d}{2}+1)} \rightarrow 0$$
  
as  $d$  increases -----> (9)

- Fraction of the volume of the cube trapped inside the sphere goes to zero as  $d$  increases.
- Empty space at the center and volume of the cube is concentrated at its  $2^d$  corners.



# Pairwise distances in $R^2$

- Consider  $V_c(2,1)$ : Generate 1001 independent, identically distributed in  $V_c(2,1)$ .

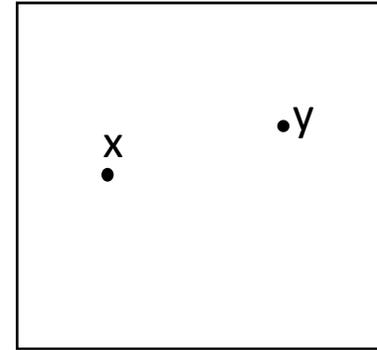
- Fix one of the point and call it  $x = (x_1, x_2)^T$ .

- Compute for each of the rest of 1000 points

$$D^2(x, y) = [(x_1 - y_1)^2 + (x_2 - y_2)^2] \quad (y \neq x).$$

- Clearly  $0 \leq D^2(x, y) \leq 2$  for all  $y \neq x$  since  $|x_1 - y_1| \leq 1$  and  $|x_2 - y_2| \leq 1$ .

- Histogram of  $D^2(x, y)$  is fully supported on  $[0,2]$ .



$V_c(2,1)$

# Pairwise distances in $R^2$ : $d = 100$

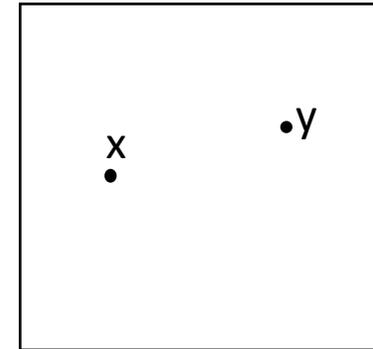
- Repeat the above experiment in  $V_c(d, 1)$ .

- Here  $\begin{cases} x = (x_1, x_2, \dots, x_d)^T \\ y = (y_1, y_2, \dots, y_d)^T \end{cases}$  with  $|x_i - y_i| \leq 1$

- $D^2(x, y) = \sum_{i=1}^d (x_i - y_i)^2 \text{ -----} \rightarrow (10)$

- Clearly  $0 \leq D^2(x, y) \leq 100$ .

- A lot more is true – thanks to the law of large numbers.



$V_c(d, 1)$

# Concentration of distances

- Clearly  $x_i$ 's and  $y_i$ 's ,  $(x_i - y_i)^2$  are i.i.d random variables with finite mean and variance.
- $D^2(x, y) = \sum_{i=1}^d (x_i - y_i)^2$  is the sum of i.i.d random variables.
- By the law of large numbers, the distribution of  $D^2(x, y)$  is concentrated in the interval  $[0, 100]$  around the mean.
- For small  $d$ , this distribution is spread out in  $[0, d]$  but for large  $d$ , it gets concentrated.

# Gaussian distribution in $\mathbb{R}^d$

- $x \in \mathbb{R}^d, m \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d}$ .

- $X \sim N(m, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x - m)^T \Sigma^{-1} (x - m) \right] \text{-----} \rightarrow (11)$

- $X \sim N(0, \sigma^2 I) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{x_i^2}{2\sigma^2} \right] \text{-----} \rightarrow (12)$

- $E[\|x\|^2] = d E(x_1^2) = d \sigma^2 \text{-----} \rightarrow (13)$

Since  $x_i$  are i.i.d  $N(0, \sigma^2)$ .

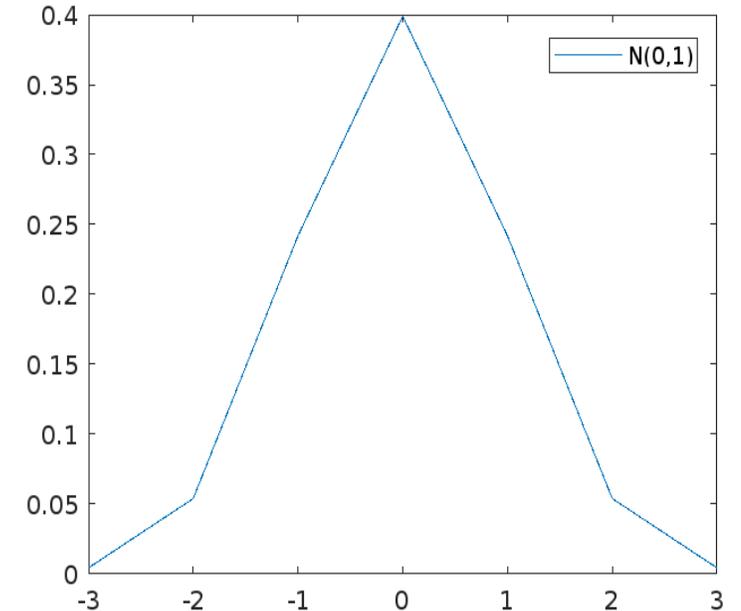
- For large  $d$ , the random variable  $\|x\|^2$  is concentrated about its mean  $d \sigma^2$ .

- $\sigma\sqrt{d}$  is called the radius of the Gaussian.

# Tail probability of N(0,1) in $R^1$

- Consider N(0,1)
- Let  $r(a) = \frac{1}{\sqrt{2\pi}} \int_{-a}^a \exp\left(-\frac{x^2}{2}\right) dx = \text{Area under N(0,1) between } -a \text{ and } a.$

a	r(a)	Tail : $1 - r(a)$
1	0.683	0.317
2	0.955	0.045
3	0.997	0.003



# Tail probability of $N(0, I)$ in $R^d$

- Probability that lies outside a sphere of radius 1.

d	1	2	5	10	20	100
P	0.317	0.1353	0.5494	0.9473	0.999	1.0

- $N(0, I)$  still attains its maximum at  $x = 0$ .
- For large  $d$ , tail has more information.
- Probability of  $N(0, I)$  contained in a thin annulus around  $\|x\|^2 = d$   
 $P[\sqrt{d} - \beta \leq \|x\| \leq \sqrt{d} + \beta] \geq 1 - 3e^{-\alpha\beta^2}$ , where  $\beta < \sqrt{d}$  and  $\alpha > 0$  is a constant.

# Chi-square distribution of $\|x\|^2$

- Let  $x \in \mathbb{R}^k$ ,  $x_i \sim \text{i.i.d. } N(0,1)$  for  $1 \leq i \leq k$ .
- $Y = \|x\|^2 = \sum_{i=1}^k x_i^2$  - chi-square distributed with  $k$  degrees of freedom given by
- $f_Y(y) = \frac{1}{2^{k/2} \Gamma(\frac{k}{2})} y^{\frac{k}{2}-1} e^{-\frac{y}{2}}$  -----> (15)
- Mean of  $Y = E[\|x\|^2] = k$  -----> (16)
- Var of  $Y = \text{VAR}(\|x\|^2) = 2k$  -----> (16)

# Chi- distribution of $\|x\|$

- Let  $Z = \|x\|$
- $Z$  said to chi-distributed

$$f_Z(z) = \frac{1}{2^{\frac{k}{2}-1} \Gamma(\frac{k}{2})} z^{k-1} e^{-\frac{z^2}{2}} \text{-----} \rightarrow (17)$$

- Mean of  $z = E[\|x\|] = \sqrt{2} \frac{\Gamma(\frac{k}{2}+1)}{\Gamma(\frac{k}{2})} \text{-----} \rightarrow (18)$

- Var of  $Z = k - \mu^2 \text{-----} \rightarrow (18)$

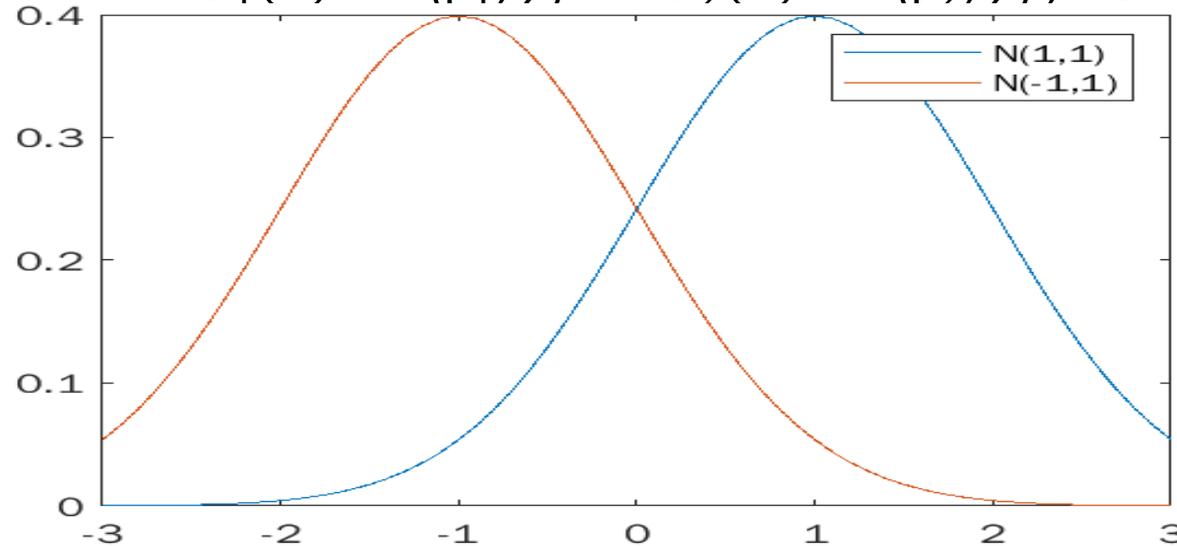
# Properties of $\|x\|$ : concentration of $\|x\|$

- Setting  $n = k + 1$ .
- $E(\|x\|) = \sqrt{n-1} \left[1 - \frac{1}{4n}\right]$
- $\text{Var}(\|x\|) = \frac{n-1}{2n} \approx \frac{1}{2} \text{-----} \rightarrow (19)$

k	n	$E(\ x\ )$	$\text{Var}(\ x\ )$
10	11	3.09	0.4545
50	51	7.106	0.4902
100	101	10.035	0.4905
500	501	22.35	0.4995

# Impact of high dimension in statistics: Linear discriminant analysis : Population based analysis

- Two Gaussian distribution  $P_1(x) = N(\mu_1, \sigma)$  and  $P_2(x) = N(\mu_2, \sigma)$ ,  $x \in \mathbb{R}^d$ .



- Mixture :  $P(x) = p_1 P_1(x) + p_2 P_2(x)$ ,  $p_1 > 0$  and  $p_1 + p_2 = 1$ .
- A sample is drawn from  $P(x)$  and need to identify which class it belongs to.

# Standard Algorithm

- Compute  $L = \log \left( \frac{P_2(x)}{P_1(x)} \right)$
- $L = \Psi(x) = \langle \mu_2 - \mu_1, \Sigma^{-1} (x - \frac{\mu_2 + \mu_1}{2}) \rangle \text{-----} \rightarrow (20)$
- Linear statistic.
- Optimum decision rule is based on thresholding  $\Psi(x)$ .
- When  $\mu_1 = 1$  and  $\mu_2 = -1$ :  $T = 0$  is a good threshold.

# Error probability

- Set  $p_1 = p_2 = \frac{1}{2}$
- Error ( $\Psi$ ) =  $\frac{1}{2} [P_1[\Psi(x') \leq 0] + P_2[\Psi(x'') > 0]]$
- $x'$  and  $x''$  are drawn from  $P_1(x)$  and  $P_2(x)$ .
- Error ( $\Psi$ ) =  $\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\frac{r}{2}} e^{-\frac{t^2}{2}} dt = \Phi\left(-\frac{r}{2}\right)$  -----> (21)
- $r^2 = (\mu_1 - \mu_2) \Sigma^{-1} (\mu_1 - \mu_2)$  : Mahalanobis Distance.

# Sample Counterpart

- We do not know the conditional distributions.
- Given a set of labelled samples:  $\{x_1, x_2, \dots, x_{n_1}\}$  from  $P_1(x)$ ,  $\{x_{n_1+1}, x_{n_1+2}, \dots, x_{n_1+n_2}\}$  from  $P_2(x)$
- Sample mean :  $\widehat{\mu}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} x_i$  and  $\widehat{\mu}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} x_{i+n_1}$
- Pooled sample covariance:
- $\widehat{\Sigma} = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (x_i - \widehat{\mu}_1)(x_i - \widehat{\mu}_1)^T + \frac{1}{n_2-1} \sum_{i=1}^{n_2} (x_{i+n_1} - \widehat{\mu}_2)(x_{i+n_1} - \widehat{\mu}_2)^T$

# Fisher's Linear discriminant function

- $\hat{\Psi}(x) = \langle \hat{\mu}_1 - \hat{\mu}_2, \hat{\Sigma}^{-1} (x - \frac{\hat{\mu}_1 + \hat{\mu}_2}{2}) \rangle \text{ -----} \rightarrow (22)$

- Assume  $n_i > d$  and  $\hat{\Sigma}$  is invertible.

- Error ( $\hat{\Psi}$ ) =  $\frac{1}{2} [P_1[\hat{\Psi}(x') \leq 0] + P_2[\hat{\Psi}(x'') > 0]] \text{ -----} \rightarrow (23)$

where  $x'$  and  $x''$  are samples from  $P_1(x)$  and  $P_2(x)$ .

# Kolmogorov's analysis (1960's)

- Assume  $\sum = I$  and  $\widehat{\Psi}_{id}(x) = \langle \widehat{\mu}_1 - \widehat{\mu}_2, x - \frac{\widehat{\mu}_1 + \widehat{\mu}_2}{2} \rangle$
- When  $(n_1 = n_2, d)$  and grow without bound with ratios  $\frac{d}{n} \rightarrow \alpha > 0$ .
- Let  $\|\widehat{\mu}_1 - \widehat{\mu}_2\| \rightarrow$  a constant  $\Upsilon > 0$ .

# Kolmogorov's Analysis Continued

- In this scaling:

$$\text{Error } (\hat{\Psi}_{id}) \rightarrow \phi \left( -\frac{r^2}{2\sqrt{r^2+\alpha}} \right) \text{ in probability} \text{ -----} \rightarrow (24)$$

- Since  $\frac{r^2}{2\sqrt{r^2+\alpha}} < \frac{r}{2}$ , Error  $(\hat{\Psi}_{id})$  is larger than when  $\alpha = 0$ .
- Clear demonstration of high- dimensional effect and resulting sub optimality.
- When  $\frac{d}{n} = \alpha = 0$ , we get the classical asymptotic result.

# Covariance estimation: Effect of high dimension

- Let  $\{x_1, x_2, \dots, x_n\}$  be an i.i.d samples from a distribution with zero mean where  $x_i \in \mathbb{R}^d$  .
- That is, we have n points chosen at random in  $\mathbb{R}^d$  .
- Let  $x = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^{d \times n}$  – Data matrix.

# Estimate Covariance matrix

- Sample Covariance :  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T = \frac{1}{n} x x^T \in \mathbb{R}^{d \times d}$
- $\hat{\Sigma}$  is unbiased :  $E(\hat{\Sigma}) = \Sigma$ .
- $\hat{\Sigma} \rightarrow \Sigma$ , the population covariance as  $n \rightarrow \infty$  when  $d$  is fixed – classical convergence.

# Measure of distance between $\hat{\Sigma}$ and $\Sigma$

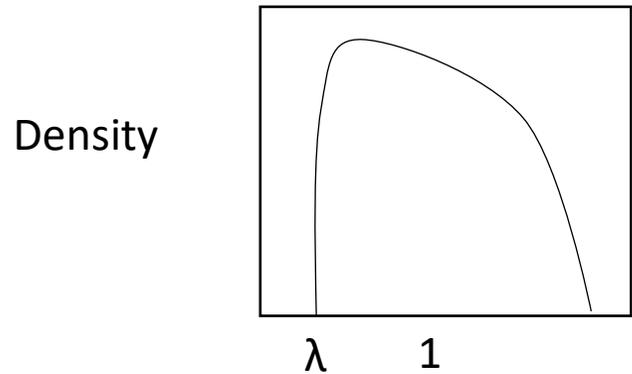
- Matrix norm – spectral norm, can be used  $\|\hat{\Sigma} - \Sigma\|_2 = \sup_{\|u\|_2=1} \|(\hat{\Sigma} - \Sigma)u\|_2 \text{ -----} \rightarrow (25)$
- It can be proved :  $\|\hat{\Sigma} - \Sigma\|_2 \rightarrow 0$  and  $n \rightarrow \infty$  .
- That is, sample covariance is strongly consistent estimate of  $\Sigma$  in classical setting.

# High dimensional effect

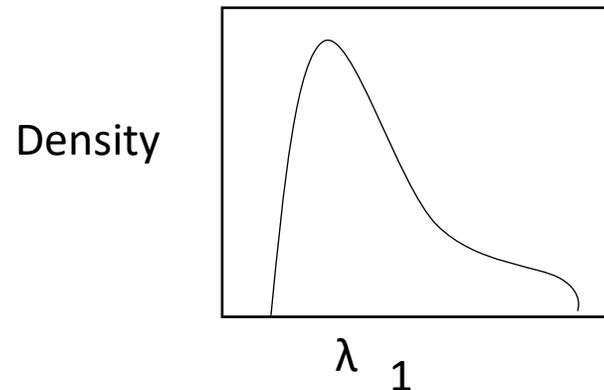
- Let  $n$  and  $d$  grow, but  $\frac{d}{n} = \alpha \in (0,1)$ .
- Estimate  $\hat{\Sigma}$  and compute its spectrum.
- Let  $\lambda_{max}(\hat{\Sigma}) = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d = \lambda_{min}(\hat{\Sigma}) \geq 0$ .

# Special case $\Sigma = I$

- In this special case when  $\frac{d}{n} = \alpha \in (0,1)$  eigen values  $\lambda_i$  are all dispersed around 1.



$n = 4000, d = 800, \alpha = 0.2$



$n = 4000, d = 2000, \alpha = 0.5$

- Empirical distribution of  $\lambda$  's for  $\alpha = 0.2$  and 0.5.

# Marcenko – Pastur law (1967) : Impact of High - dimension

- M-P law : They proved that the density of distribution of  $\lambda$  's is supported on the interval  $[t_{min}(\alpha), t_{max}(\alpha)]$  where  $t_{min}(\alpha) = (1 - \sqrt{\alpha})^2$  and  $t_{max}(\alpha) = (1 + \sqrt{\alpha})^2$ .
- This law allows  $(d, n)$  to increase but  $\frac{d}{n} = \alpha \in (0,1)$  - has a non – classical flavor.

# References

[1] M.J. Wainwright (2019) High- Dimensional Statistics: a non – asymptotic viewpoint, Cambridge University Press.

[2] R. Vershynin (2020) High- Dimensional Probability: An introduction with application in Data Science, Cambridge University Press.

[3] A. Blum, J. Hopcroft and R. Kannan (2020) Foundation of Data Science, Cambridge University Press.

[4] V.Shikhman and D. Mueller (2021) Mathematical Foundation of Big Data Analysis, Springer.

[5] J.Wang (2012) Geometric Structure of High- Dimensional Data and Dimensionality reduction, Springer.