

Norm of a Gaussian vector

Consider $X \sim \mathcal{N}(0, \mathbb{I}_d)$. We want to understand the distribution of $\|X\|_2$.

Observation 1

$$\|X\|_2^2 = \sum_{i=1}^d X_i^2, \text{ where}$$

$X_i \sim \mathcal{N}(0, 1)$ are iid.

However, X_i^2 are not bounded, and not even sub-Gaussian (why should one not expect them to be sub-Gaussian?) so the usual Chernoff-Hoeffding bounds do not apply. (and are not even true).

However, it is true that $\|X\|_2$ concentrates around its expectation.

This is in fact a consequence of a general result called the Gaussian concentration inequality: If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is L -Lipschitz and $X \sim \mathcal{N}(0, I_d)$ then

$\forall t \geq 0$

$$\mathbb{P}_X [f(X) - \mathbb{E}f(X) \geq t] \leq \exp\left(-\frac{t^2}{2L^2}\right).$$

Ex:- Show that this immediately implies

$$\mathbb{P}_X [\|X\| > \sqrt{d} + t] \leq \exp\left(-\frac{t^2}{2}\right)$$

[Note that $\|X\|$ indeed is sub-Gaussian, although $\|X\|^2$ is not]

\times

We will use a less general direct approach. Note that we have

$$\mathbb{P}_X [\|X\|^2 \geq d+t] \leq \frac{\mathbb{E}[\exp(aX_1^2)]^d}{\exp(a(d+t))} \quad \forall a \geq 0.$$

$$\text{Now } \mathbb{E}[\exp(aX_1^2)] = \frac{1}{\sqrt{1-2a}} \quad \forall a < \frac{1}{2}$$

[This is a direct calculation]

An aside := In 1-dimension, if $Z \sim N(0, 1)$,
then $\forall t > 0$

$$\Pr [Z > t] \leq \frac{1}{t\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right).$$

(Why?).

This gives $\forall a \geq 0, t \geq 0$

$$\Pr [\|X\|^2 \geq d+t] \leq \exp(-a(d+t) - \frac{d}{2} \log(1-2a))$$

The exponent is minimized at

$$a = \frac{t}{2(d+t)}, \text{ so that we}$$

have

$$\Pr [\|X\|^2 \geq d+t] \leq \exp\left(-\frac{t}{2} + \frac{d}{2} \log\left(1 + \frac{t}{d}\right)\right)$$

Using $\log(1+x) \leq x - \frac{x^2}{3}$ for $x \in [0, \frac{1}{2}]$,

we get, for $t \in [0, \frac{d}{2}]$,

$$\Pr [\|X\|^2 \geq d+t] \leq \exp\left(-\frac{t^2}{6d}\right)$$

Similar calculations show that for $t \in [0, d]$,

$$\Pr [\|X\|^2 \leq d-t] \leq \exp\left(-\frac{t^2}{4d}\right)$$

In particular for $t = 6\sqrt{\alpha d \log K}$; $s.t. 6\sqrt{\alpha \log K} \leq \frac{\sqrt{d}}{2}$

$$\Pr [\|X\|^2 \geq d + 6\sqrt{\alpha d \log K}] \leq K^{-6\alpha}.$$

and

$$\Pr [\|X\|^2 \leq d - 6\sqrt{\alpha d \log K}] \leq K^{-6\alpha}.$$

This also implies that (for $t \leq \frac{\sqrt{d}}{4}$)

$$\Pr [\|X\| \geq \sqrt{d} + t] \leq \exp\left(-\frac{t^2}{2}\right)$$

$$\Pr [\|X\| \leq \sqrt{d} - t] \leq \exp\left(-\frac{t^2}{2}\right).$$

Clustering Gaussians

We now consider the following problem: We have $2n$ points in \mathbb{R}^d , each sampled independently w.p. $\frac{1}{2}$ from $\mathcal{N}(\mu_1, I_d)$ and with probability $\frac{1}{2}$ from $\mathcal{N}(\mu_2, I_d)$ where $\mu_2, \mu_1 \in \mathbb{R}^d$ are "far". The goal is to label each pt. with the Gaussian it came from.

A Heuristic :- One might expect that there is a threshold d_x s.t. if $X, Y \sim \mathcal{N}(\mu_1, I_d)$, then w.h.p.

$$\|X - Y\|_2^2 \leq d_x^2 \text{ while if}$$

$$X \sim \mathcal{N}(\mu_1, I_d) \quad Y \sim \mathcal{N}(\mu_2, I_d)$$

$$\|X - Y\|_2^2 \gg d_x^2$$

We now show that when $\|\mu_1 - \mu_2\| \geq \Omega(d^{1/4})$, such a threshold exists.

Case 1 :- $X, Y \sim \mathcal{N}(\mu_i, I_d)$ $i=1$ or 2 .

$$\text{Then } Z = \frac{X - Y}{\sqrt{2}} \sim \mathcal{N}(0, I_d).$$

Hence, from above concentration results, when $\log(2n) \leq \frac{\sqrt{d}}{12}$,

$$|\|X - Y\|_2^2 - 2d| \leq 12 \sqrt{d \log(2n)} \text{ w.p. } \geq 1 - \frac{1}{4n^3}$$

Thus, taking a union bound, we see that for all pairs (X, Y) which were drawn from the same cluster,

$$\left| \|X - Y\|_2^2 - 2d \right| \leq 12 \sqrt{d \log(2n)} \quad \text{w.p.} \geq 1 - \frac{1}{2n}.$$

— (1)

Case 2 :- $X \sim \mathcal{N}(\mu_1, I_d)$
 $Y \sim \mathcal{N}(\mu_2, I_d)$ independently.

$\mu_1 - \mu_2 = \Delta.$ Then,

$X - Y \stackrel{d}{=} \sqrt{2}Z + \Delta$ where $Z \sim \mathcal{N}(0, I_d)$

(i.e. $X - Y$ has the same distribution as $\sqrt{2}Z + \Delta$)

Then,

$$\|X - Y\|_2^2 \stackrel{d}{=} 2\|Z\|^2 + \Delta^2 + 2\sqrt{2}\Delta \cdot Z$$

$$\geq 2\|Z\|^2 - 2\sqrt{2}\|\Delta\| \cdot \|\hat{\Delta} \cdot Z\| + \|\Delta\|^2, \quad \text{where}$$

$\hat{\Delta}$ is the unit vector in the direction of Δ .

Now $Z \sim N(0, I_d)$ so,

$Z \cdot \hat{\Delta} \sim N(0, 1)$, Therefore,

$$\Pr(|Z \cdot \hat{\Delta}| > 2d^{1/4} \sqrt{\log(2n)}) \leq \frac{2}{(2n)^{2\sqrt{d}}}$$

(Assuming $n \geq 2, d \geq 4$).

$$\leq \frac{1}{n^{2\sqrt{d}}}$$

Let $\Delta = c d^{1/4} \sqrt{\log 2n}$ where $c > 6$ is a constant. Then we get that

$$\|X - Y\|_2 \geq 2d - 12\sqrt{d \log 2n} + 2c\sqrt{d \log 2n}$$

with probability $\geq 1 - \frac{5}{4n^3}$.

Put $c = 15$. Then.

$$\|X - Y\|_2 \geq 2d + 18\sqrt{d \log 2n} \text{ w.p.}$$
$$\geq 1 - \frac{5}{4n^3}$$

provided $\|\mu_1 - \mu_2\| \geq 15d^{1/4} \sqrt{\log 2n}$

— (2)

From ① and ② via an union bound, we get: w.p. $\geq 1 - \frac{1}{n}$,

we have

$$\|X - Y\|_{-2d}^2 \begin{cases} \geq 18 \sqrt{d \log(2n)} & \text{if } X, Y \text{ sampled} \\ & \text{from } \mathcal{N}(\mu_1, I_d), \\ & \mathcal{N}(\mu_2, I_d) \\ & \text{respectively} \\ \leq 12 \sqrt{d \log(2n)} & \text{if } X, Y \text{ sampled} \\ & \text{from the} \\ & \text{same} \\ & \mathcal{N}(\mu_i, I_d); \end{cases}$$

provided.

$$\|\mu_1 - \mu_2\| \geq 15 d^{1/4} \sqrt{\log n}.$$

Thus we can choose a threshold d^* as above when

$$\|\mu_1 - \mu_2\| \geq d^{1/4} \sqrt{\log 2n}.$$