

2 High-Dimensional Space

2.1 Introduction

High dimensional data has become very important. However, high dimensional space is very different from the two and three dimensional spaces we are familiar with. Generate n points at random in d -dimensions where each coordinate is a zero mean, unit variance Gaussian. For sufficiently large d , with high probability the distances between all pairs of points will be essentially the same. Also the volume of the unit ball in d -dimensions, the set of all points \mathbf{x} such that $|\mathbf{x}| \leq 1$, goes to zero as the dimension goes to infinity. The volume of a high dimensional unit ball is concentrated near its surface and is also concentrated at its equator. These properties have important consequences which we will consider.

2.2 The Law of Large Numbers

If one generates random points in d -dimensional space using a Gaussian to generate coordinates, the distance between all pairs of points will be essentially the same when d is large. The reason is that the square of the distance between two points \mathbf{y} and \mathbf{z} ,

$$|\mathbf{y} - \mathbf{z}|^2 = \sum_{i=1}^d (y_i - z_i)^2,$$

can be viewed as the sum of d independent samples of a random variable x that is distributed as the squared difference of two Gaussians. In particular, we are summing independent samples $x_i = (y_i - z_i)^2$ of a random variable x of bounded variance. In such a case, a general bound known as the Law of Large Numbers states that with high probability, the average of the samples will be close to the expectation of the random variable. This in turn implies that with high probability, the sum is close to the sum's expectation.

Specifically, the Law of Large Numbers states that

$$\text{Prob} \left(\left| \frac{x_1 + x_2 + \cdots + x_n}{n} - E(x) \right| \geq \epsilon \right) \leq \frac{\text{Var}(x)}{n\epsilon^2}. \quad (2.1)$$

The larger the variance of the random variable, the greater the probability that the error will exceed ϵ . Thus the variance of x is in the numerator. The number of samples n is in the denominator since the more values that are averaged, the smaller the probability that the difference will exceed ϵ . Similarly the larger ϵ is, the smaller the probability that the difference will exceed ϵ and hence ϵ is in the denominator. Notice that squaring ϵ makes the fraction a dimensionless quantity.

We use two inequalities to prove the Law of Large Numbers. The first is Markov's inequality that states that the probability that a nonnegative random variable exceeds a is bounded by the expected value of the variable divided by a .

Theorem 2.1 (Markov's inequality) *Let x be a nonnegative random variable. Then for $a > 0$,*

$$\text{Prob}(x \geq a) \leq \frac{E(x)}{a}.$$

Proof: For a continuous nonnegative random variable x with probability density p ,

$$\begin{aligned} E(x) &= \int_0^{\infty} xp(x)dx = \int_0^a xp(x)dx + \int_a^{\infty} xp(x)dx \\ &\geq \int_a^{\infty} xp(x)dx \geq a \int_a^{\infty} p(x)dx = a\text{Prob}(x \geq a). \end{aligned}$$

Thus, $\text{Prob}(x \geq a) \leq \frac{E(x)}{a}$. ■

The same proof works for discrete random variables with sums instead of integrals.

Corollary 2.2 $\text{Prob}(x \geq bE(x)) \leq \frac{1}{b}$

Markov's inequality bounds the tail of a distribution using only information about the mean. A tighter bound can be obtained by also using the variance of the random variable.

Theorem 2.3 (Chebyshev's inequality) *Let x be a random variable. Then for $c > 0$,*

$$\text{Prob}(|x - E(x)| \geq c) \leq \frac{\text{Var}(x)}{c^2}.$$

Proof: $\text{Prob}(|x - E(x)| \geq c) = \text{Prob}(|x - E(x)|^2 \geq c^2)$. Let $y = |x - E(x)|^2$. Note that y is a nonnegative random variable and $E(y) = \text{Var}(x)$, so Markov's inequality can be applied giving:

$$\text{Prob}(|x - E(x)| \geq c) = \text{Prob}(|x - E(x)|^2 \geq c^2) \leq \frac{E(|x - E(x)|^2)}{c^2} = \frac{\text{Var}(x)}{c^2}.$$
■

The Law of Large Numbers follows from Chebyshev's inequality together with facts about independent random variables. Recall that:

$$\begin{aligned} E(x + y) &= E(x) + E(y), \\ \text{Var}(x - c) &= \text{Var}(x), \\ \text{Var}(cx) &= c^2\text{Var}(x). \end{aligned}$$

Also, if x and y are independent, then $E(xy) = E(x)E(y)$. These facts imply that if x and y are independent then $Var(x + y) = Var(x) + Var(y)$, which is seen as follows:

$$\begin{aligned} Var(x + y) &= E(x + y)^2 - E^2(x + y) \\ &= E(x^2 + 2xy + y^2) - (E^2(x) + 2E(x)E(y) + E^2(y)) \\ &= E(x^2) - E^2(x) + E(y^2) - E^2(y) = Var(x) + Var(y), \end{aligned}$$

where we used independence to replace $E(2xy)$ with $2E(x)E(y)$.

Theorem 2.4 (Law of Large Numbers) *Let x_1, x_2, \dots, x_n be n independent samples of a random variable x . Then*

$$Prob\left(\left|\frac{x_1 + x_2 + \dots + x_n}{n} - E(x)\right| \geq \epsilon\right) \leq \frac{Var(x)}{n\epsilon^2}$$

Proof: By Chebychev's inequality

$$\begin{aligned} Prob\left(\left|\frac{x_1 + x_2 + \dots + x_n}{n} - E(x)\right| \geq \epsilon\right) &\leq \frac{Var\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right)}{\epsilon^2} \\ &= \frac{1}{n^2\epsilon^2} Var(x_1 + x_2 + \dots + x_n) \\ &= \frac{1}{n^2\epsilon^2} (Var(x_1) + Var(x_2) + \dots + Var(x_n)) \\ &= \frac{Var(x)}{n\epsilon^2}. \end{aligned}$$

■

The Law of Large Numbers is quite general, applying to any random variable x of finite variance. Later we will look at tighter concentration bounds for spherical Gaussians and sums of 0-1 valued random variables.

One observation worth making about the Law of Large Numbers is that the size of the universe does not enter into the bound. For instance, if you want to know what fraction of the population of a country prefers tea to coffee, then the number n of people you need to sample in order to have at most a δ chance that your estimate is off by more than ϵ depends only on ϵ and δ and not on the population of the country.

As an application of the Law of Large Numbers, let \mathbf{z} be a d -dimensional random point whose coordinates are each selected from a zero mean, $\frac{1}{2\pi}$ variance Gaussian. We set the variance to $\frac{1}{2\pi}$ so the Gaussian probability density equals one at the origin and is bounded below throughout the unit ball by a constant.¹ By the Law of Large Numbers, the square of the distance of \mathbf{z} to the origin will be $\Theta(d)$ with high probability. In particular, there is

¹If we instead used variance 1, then the density at the origin would be a decreasing function of d , namely $(\frac{1}{2\pi})^{d/2}$, making this argument more complicated.

vanishingly small probability that such a random point \mathbf{z} would lie in the unit ball. This implies that the integral of the probability density over the unit ball must be vanishingly small. On the other hand, the probability density in the unit ball is bounded below by a constant. We thus conclude that the unit ball must have vanishingly small volume.

Similarly if we draw two points \mathbf{y} and \mathbf{z} from a d -dimensional Gaussian with unit variance in each direction, then $|\mathbf{y}|^2 \approx d$ and $|\mathbf{z}|^2 \approx d$. Since for all i ,

$$E(y_i - z_i)^2 = E(y_i^2) + E(z_i^2) - 2E(y_i z_i) = \text{Var}(y_i) + \text{Var}(z_i) + 2E(y_i)E(z_i) = 2,$$

$|\mathbf{y} - \mathbf{z}|^2 = \sum_{i=1}^d (y_i - z_i)^2 \approx 2d$. Thus by the Pythagorean theorem, the random d -dimensional \mathbf{y} and \mathbf{z} must be approximately orthogonal. This implies that if we scale these random points to be unit length and call \mathbf{y} the North Pole, much of the surface area of the unit ball must lie near the equator. We will formalize these and related arguments in subsequent sections.

We now state a general theorem on probability tail bounds for a sum of independent random variables. Tail bounds for sums of Bernoulli, squared Gaussian and Power Law distributed random variables can all be derived from this. The table in Figure 2.1 summarizes some of the results.

Theorem 2.5 (Master Tail Bounds Theorem) *Let $x = x_1 + x_2 + \dots + x_n$, where x_1, x_2, \dots, x_n are mutually independent random variables with zero mean and variance at most σ^2 . Let $0 \leq a \leq \sqrt{2n\sigma^2}$. Assume that $|E(x_i^s)| \leq \sigma^2 s!$ for $s = 3, 4, \dots, \lfloor (a^2/4n\sigma^2) \rfloor$. Then,*

$$\text{Prob}(|x| \geq a) \leq 3e^{-a^2/(12n\sigma^2)}.$$

The proof of Theorem 2.5 is elementary. A slightly more general version, Theorem 12.5, is given in the appendix. For a brief intuition of the proof, consider applying Markov's inequality to the random variable x^r where r is a large even number. Since r is even, x^r is nonnegative, and thus $\text{Prob}(|x| \geq a) = \text{Prob}(x^r \geq a^r) \leq E(x^r)/a^r$. If $E(x^r)$ is not too large, we will get a good bound. To compute $E(x^r)$, write $E(x)$ as $E(x_1 + \dots + x_n)^r$ and expand the polynomial into a sum of terms. Use the fact that by independence $E(x_i^{r_i} x_j^{r_j}) = E(x_i^{r_i}) E(x_j^{r_j})$ to get a collection of simpler expectations that can be bounded using our assumption that $|E(x_i^s)| \leq \sigma^2 s!$. For the full proof, see the appendix.

2.3 The Geometry of High Dimensions

An important property of high-dimensional objects is that most of their volume is near the surface. Consider any object A in R^d . Now shrink A by a small amount ϵ to produce a new object $(1 - \epsilon)A = \{(1 - \epsilon)x | x \in A\}$. Then the following equality holds:

$$\text{volume}((1 - \epsilon)A) = (1 - \epsilon)^d \text{volume}(A).$$

	Condition	Tail bound
Markov	$x \geq 0$	$\text{Prob}(x \geq a) \leq \frac{E(x)}{a}$
Chebychev	Any x	$\text{Prob}(x - E(x) \geq a) \leq \frac{\text{Var}(x)}{a^2}$
Chernoff	$x = x_1 + x_2 + \dots + x_n$ $x_i \in [0, 1]$ i.i.d. Bernoulli;	$\text{Prob}(x - E(x) \geq \varepsilon E(x))$ $\leq 3e^{-c\varepsilon^2 E(x)}$
Higher Moments	r positive even integer	$\text{Prob}(x \geq a) \leq E(x^r)/a^r$
Gaussian Annulus	$x = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$ $x_i \sim N(0, 1)$; $\beta \leq \sqrt{n}$ indep.	$\text{Prob}(x - \sqrt{n} \geq \beta) \leq 3e^{-c\beta^2}$
Power Law for x_i ; order $k \geq 4$	$x = x_1 + x_2 + \dots + x_n$ x_i i.i.d ; $\varepsilon \leq 1/k^2$	$\text{Prob}(x - E(x) \geq \varepsilon E(x))$ $\leq (4/\varepsilon^2 kn)^{(k-3)/2}$

Figure 2.1: Table of Tail Bounds. The Higher Moments bound is obtained by applying Markov to x^r . The Chernoff, Gaussian Annulus, and Power Law bounds follow from Theorem 2.5 which is proved in the appendix.

To see that this is true, partition A into infinitesimal cubes. Then, $(1 - \varepsilon)A$ is the union of a set of cubes obtained by shrinking the cubes in A by a factor of $1 - \varepsilon$. When we shrink each of the $2d$ sides of a d -dimensional cube by a factor f , its volume shrinks by a factor of f^d . Using the fact that $1 - x \leq e^{-x}$, for any object A in R^d we have:

$$\frac{\text{volume}((1 - \varepsilon)A)}{\text{volume}(A)} = (1 - \varepsilon)^d \leq e^{-\varepsilon d}.$$

Fixing ε and letting $d \rightarrow \infty$, the above quantity rapidly approaches zero. This means that nearly all of the volume of A must be in the portion of A that does not belong to the region $(1 - \varepsilon)A$.

Let S denote the unit ball in d dimensions, that is, the set of points within distance one of the origin. An immediate implication of the above observation is that at least a $1 - e^{-\varepsilon d}$ fraction of the volume of the unit ball is concentrated in $S \setminus (1 - \varepsilon)S$, namely in a small annulus of width ε at the boundary. In particular, most of the volume of the d -dimensional unit ball is contained in an annulus of width $O(1/d)$ near the boundary. If the ball is of radius r , then the annulus width is $O(\frac{r}{d})$.

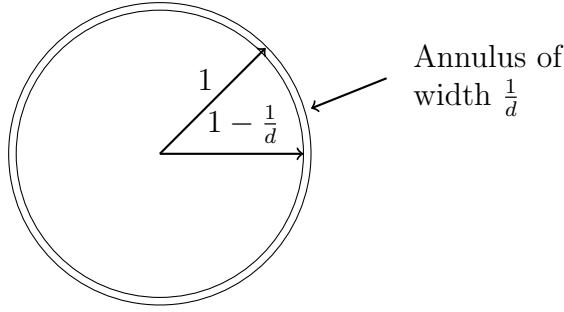


Figure 2.2: Most of the volume of the d -dimensional ball of radius r is contained in an annulus of width $O(r/d)$ near the boundary.

2.4 Properties of the Unit Ball

We now focus more specifically on properties of the unit ball in d -dimensional space. We just saw that most of its volume is concentrated in a small annulus of width $O(1/d)$ near the boundary. Next we will show that in the limit as d goes to infinity, the volume of the ball goes to zero. This result can be proven in several ways. Here we use integration.

2.4.1 Volume of the Unit Ball

To calculate the volume $V(d)$ of the unit ball in R^d , one can integrate in either Cartesian or polar coordinates. In Cartesian coordinates the volume is given by

$$V(d) = \int_{x_1=-1}^{x_1=1} \int_{x_2=-\sqrt{1-x_1^2}}^{x_2=\sqrt{1-x_1^2}} \cdots \int_{x_d=-\sqrt{1-x_1^2-\cdots-x_{d-1}^2}}^{x_d=\sqrt{1-x_1^2-\cdots-x_{d-1}^2}} dx_d \cdots dx_2 dx_1.$$

Since the limits of the integrals are complicated, it is easier to integrate using polar coordinates. In polar coordinates, $V(d)$ is given by

$$V(d) = \int_{S^d} \int_{r=0}^1 r^{d-1} dr d\Omega.$$

Since the variables Ω and r do not interact,

$$V(d) = \int_{S^d} d\Omega \int_{r=0}^1 r^{d-1} dr = \frac{1}{d} \int_{S^d} d\Omega = \frac{A(d)}{d}$$

where $A(d)$ is the surface area of the d -dimensional unit ball. For instance, for $d = 3$ the surface area is 4π and the volume is $\frac{4}{3}\pi$. The question remains, how to determine the

surface area $A(d) = \int_{S^d} d\Omega$ for general d .

Consider a different integral

$$I(d) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{-(x_1^2 + x_2^2 + \cdots + x_d^2)} dx_d \cdots dx_2 dx_1.$$

Including the exponential allows integration to infinity rather than stopping at the surface of the sphere. Thus, $I(d)$ can be computed by integrating in both Cartesian and polar coordinates. Integrating in polar coordinates will relate $I(d)$ to the surface area $A(d)$. Equating the two results for $I(d)$ allows one to solve for $A(d)$.

First, calculate $I(d)$ by integration in Cartesian coordinates.

$$I(d) = \left[\int_{-\infty}^{\infty} e^{-x^2} dx \right]^d = (\sqrt{\pi})^d = \pi^{\frac{d}{2}}.$$

Here, we have used the fact that $\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$. For a proof of this, see Section 12.3 of the appendix. Next, calculate $I(d)$ by integrating in polar coordinates. The volume of the differential element is $r^{d-1} d\Omega dr$. Thus,

$$I(d) = \int_{S^d} d\Omega \int_0^{\infty} e^{-r^2} r^{d-1} dr.$$

The integral $\int_{S^d} d\Omega$ is the integral over the entire solid angle and gives the surface area,

$A(d)$, of a unit sphere. Thus, $I(d) = A(d) \int_0^{\infty} e^{-r^2} r^{d-1} dr$. Evaluating the remaining integral gives

$$\int_0^{\infty} e^{-r^2} r^{d-1} dr = \int_0^{\infty} e^{-t} t^{\frac{d-1}{2}} \left(\frac{1}{2} t^{-\frac{1}{2}} dt \right) = \frac{1}{2} \int_0^{\infty} e^{-t} t^{\frac{d}{2} - 1} dt = \frac{1}{2} \Gamma\left(\frac{d}{2}\right)$$

and hence, $I(d) = A(d) \frac{1}{2} \Gamma\left(\frac{d}{2}\right)$ where the Gamma function $\Gamma(x)$ is a generalization of the factorial function for noninteger values of x . $\Gamma(x) = (x-1) \Gamma(x-1)$, $\Gamma(1) = \Gamma(2) = 1$, and $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$. For integer x , $\Gamma(x) = (x-1)!$.

Combining $I(d) = \pi^{\frac{d}{2}}$ with $I(d) = A(d) \frac{1}{2} \Gamma\left(\frac{d}{2}\right)$ yields

$$A(d) = \frac{\pi^{\frac{d}{2}}}{\frac{1}{2} \Gamma\left(\frac{d}{2}\right)}$$

establishing the following lemma.

Lemma 2.6 *The surface area $A(d)$ and the volume $V(d)$ of a unit-radius ball in d dimensions are given by*

$$A(d) = \frac{2\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})} \quad \text{and} \quad V(d) = \frac{2\pi^{\frac{d}{2}}}{d \Gamma(\frac{d}{2})}.$$

To check the formula for the volume of a unit ball, note that $V(2) = \pi$ and $V(3) = \frac{2}{3} \frac{\pi^{\frac{3}{2}}}{\Gamma(\frac{3}{2})} = \frac{4}{3}\pi$, which are the correct volumes for the unit balls in two and three dimensions. To check the formula for the surface area of a unit ball, note that $A(2) = 2\pi$ and $A(3) = \frac{2\pi^{\frac{3}{2}}}{\frac{1}{2}\sqrt{\pi}} = 4\pi$, which are the correct surface areas for the unit ball in two and three dimensions. Note that $\pi^{\frac{d}{2}}$ is an exponential in $\frac{d}{2}$ and $\Gamma(\frac{d}{2})$ grows as the factorial of $\frac{d}{2}$. This implies that $\lim_{d \rightarrow \infty} V(d) = 0$, as claimed.

2.4.2 Volume Near the Equator

An interesting fact about the unit ball in high dimensions is that most of its volume is concentrated near its “equator”. In particular, for any unit-length vector \mathbf{v} defining “north”, most of the volume of the unit ball lies in the thin slab of points whose dot-product with \mathbf{v} has magnitude $O(1/\sqrt{d})$. To show this fact, it suffices by symmetry to fix \mathbf{v} to be the first coordinate vector. That is, we will show that most of the volume of the unit ball has $|x_1| = O(1/\sqrt{d})$. Using this fact, we will show that two random points in the unit ball are with high probability nearly orthogonal, and also give an alternative proof from the one in Section 2.4.1 that the volume of the unit ball goes to zero as $d \rightarrow \infty$.

Theorem 2.7 *For $c \geq 1$ and $d \geq 3$, at least a $1 - \frac{2}{c}e^{-c^2/2}$ fraction of the volume of the d -dimensional unit ball has $|x_1| \leq \frac{c}{\sqrt{d-1}}$.*

Proof: By symmetry we just need to prove that at most a $\frac{2}{c}e^{-c^2/2}$ fraction of the half of the ball with $x_1 \geq 0$ has $x_1 \geq \frac{c}{\sqrt{d-1}}$. Let A denote the portion of the ball with $x_1 \geq \frac{c}{\sqrt{d-1}}$ and let H denote the upper hemisphere. We will then show that the ratio of the volume of A to the volume of H goes to zero by calculating an upper bound on $\text{volume}(A)$ and a lower bound on $\text{volume}(H)$ and proving that

$$\frac{\text{volume}(A)}{\text{volume}(H)} \leq \frac{\text{upper bound volume}(A)}{\text{lower bound volume}(H)} = \frac{2}{c}e^{-\frac{c^2}{2}}.$$

To calculate the volume of A , integrate an incremental volume that is a disk of width dx_1 and whose face is a ball of dimension $d-1$ and radius $\sqrt{1-x_1^2}$. The surface area of the disk is $(1-x_1^2)^{\frac{d-1}{2}}V(d-1)$ and the volume above the slice is

$$\text{volume}(A) = \int_{\frac{c}{\sqrt{d-1}}}^1 (1-x_1^2)^{\frac{d-1}{2}} V(d-1) dx_1$$

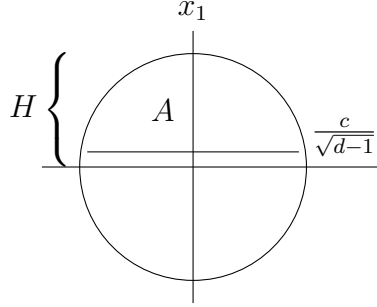


Figure 2.3: Most of the volume of the upper hemisphere of the d -dimensional ball is below the plane $x_1 = \frac{c}{\sqrt{d-1}}$.

To get an upper bound on the above integral, use $1 - x \leq e^{-x}$ and integrate to infinity. To integrate, insert $\frac{x_1\sqrt{d-1}}{c}$, which is greater than one in the range of integration, into the integral. Then

$$\text{volume}(A) \leq \int_{\frac{c}{\sqrt{d-1}}}^{\infty} \frac{x_1\sqrt{d-1}}{c} e^{-\frac{d-1}{2}x_1^2} V(d-1) dx_1 = V(d-1) \frac{\sqrt{d-1}}{c} \int_{\frac{c}{\sqrt{d-1}}}^{\infty} x_1 e^{-\frac{d-1}{2}x_1^2} dx_1$$

Now

$$\int_{\frac{c}{\sqrt{d-1}}}^{\infty} x_1 e^{-\frac{d-1}{2}x_1^2} dx_1 = -\frac{1}{d-1} e^{-\frac{d-1}{2}x_1^2} \Big|_{\frac{c}{\sqrt{d-1}}}^{\infty} = \frac{1}{d-1} e^{-\frac{c^2}{2}}$$

Thus, an upper bound on $\text{volume}(A)$ is $\frac{V(d-1)}{c\sqrt{d-1}} e^{-\frac{c^2}{2}}$.

The volume of the hemisphere below the plane $x_1 = \frac{c}{\sqrt{d-1}}$ is a lower bound on the entire volume of the upper hemisphere and this volume is at least that of a cylinder of height $\frac{1}{\sqrt{d-1}}$ and radius $\sqrt{1 - \frac{1}{d-1}}$. The volume of the cylinder is $V(d-1)(1 - \frac{1}{d-1})^{\frac{d-1}{2}} \frac{1}{\sqrt{d-1}}$. Using the fact that $(1-x)^a \geq 1-ax$ for $a \geq 1$, the volume of the cylinder is at least $\frac{V(d-1)}{2\sqrt{d-1}}$ for $d \geq 3$.

Thus,

$$\text{ratio} \leq \frac{\text{upper bound above plane}}{\text{lower bound total hemisphere}} = \frac{\frac{V(d-1)}{c\sqrt{d-1}} e^{-\frac{c^2}{2}}}{\frac{V(d-1)}{2\sqrt{d-1}}} = \frac{2}{c} e^{-\frac{c^2}{2}}$$

■

One might ask why we computed a lower bound on the total hemisphere since it is one half of the volume of the unit ball which we already know. The reason is that the volume of the upper hemisphere is $\frac{1}{2}V(d)$ and we need a formula with $V(d-1)$ in it to cancel the $V(d-1)$ in the numerator.

Near orthogonality. One immediate implication of the above analysis is that if we draw two points at random from the unit ball, with high probability their vectors will be nearly orthogonal to each other. Specifically, from our previous analysis in Section 2.3, with high probability both will be close to the surface and will have length $1 - O(1/d)$. From our analysis above, if we define the vector in the direction of the first point as “north”, with high probability the second will have a projection of only $\pm O(1/\sqrt{d})$ in this direction, and thus their dot-product will be $\pm O(1/\sqrt{d})$. This implies that with high probability, the angle between the two vectors will be $\pi/2 \pm O(1/\sqrt{d})$. In particular, we have the following theorem that states that if we draw n points at random in the unit ball, with high probability all points will be close to unit length and each pair of points will be almost orthogonal.

Theorem 2.8 *Consider drawing n points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ at random from the unit ball. With probability $1 - O(1/n)$*

1. $|\mathbf{x}_i| \geq 1 - \frac{2 \ln n}{d}$ for all i , and
2. $|\mathbf{x}_i \cdot \mathbf{x}_j| \leq \frac{\sqrt{6 \ln n}}{\sqrt{d-1}}$ for all $i \neq j$.

Proof: For the first part, for any fixed i by the analysis of Section 2.3, the probability that $|\mathbf{x}_i| < 1 - \epsilon$ is less than $e^{-\epsilon d}$. Thus

$$\text{Prob}(|\mathbf{x}_i| < 1 - \frac{2 \ln n}{d}) \leq e^{-(\frac{2 \ln n}{d})d} = 1/n^2.$$

By the union bound, the probability there exists an i such that $|\mathbf{x}_i| < 1 - \frac{2 \ln n}{d}$ is at most $1/n$.

For the second part, Theorem 2.7 states that the probability $|\mathbf{x}_i| > \frac{c}{\sqrt{d-1}}$ is at most $\frac{2}{c} e^{-\frac{c^2}{2}}$. There are $\binom{n}{2}$ pairs i and j and for each such pair if we define \mathbf{x}_i as “north”, the probability that the projection of \mathbf{x}_j onto the “north” direction is more than $\frac{\sqrt{6 \ln n}}{\sqrt{d-1}}$ is at most $O(e^{-\frac{6 \ln n}{2}}) = O(n^{-3})$. Thus, the dot-product condition is violated with probability at most $O(\binom{n}{2} n^{-3}) = O(1/n)$ as well. ■

Alternative proof that volume goes to zero. Another immediate implication of Theorem 2.7 is that as $d \rightarrow \infty$, the volume of the ball approaches zero. Specifically, consider a small box centered at the origin of side length $\frac{2c}{\sqrt{d-1}}$. Using Theorem 2.7, we show that for $c = 2\sqrt{\ln d}$, this box contains over half of the volume of the ball. On the other hand, the volume of this box clearly goes to zero as d goes to infinity, since its volume is $O((\frac{\ln d}{d-1})^{d/2})$. Thus the volume of the ball goes to zero as well.

By Theorem 2.7 with $c = 2\sqrt{\ln d}$, the fraction of the volume of the ball with $|x_1| \geq \frac{c}{\sqrt{d-1}}$ is at most:

$$\frac{2}{c} e^{-\frac{c^2}{2}} = \frac{1}{\sqrt{\ln d}} e^{-2 \ln d} = \frac{1}{d^2 \sqrt{\ln d}} < \frac{1}{d^2}.$$

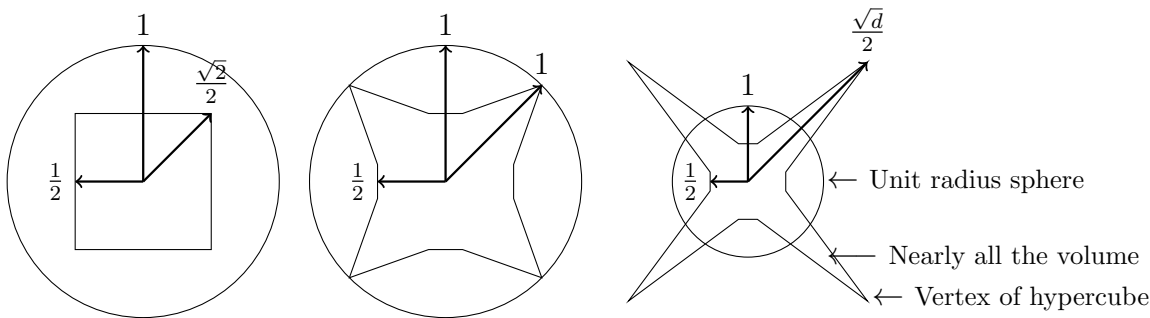


Figure 2.4: Illustration of the relationship between the sphere and the cube in 2, 4, and d -dimensions.

Since this is true for each of the d dimensions, by a union bound at most a $O(\frac{1}{d}) \leq \frac{1}{2}$ fraction of the volume of the ball lies outside the cube, completing the proof.

Discussion. One might wonder how it can be that nearly all the points in the unit ball are very close to the surface and yet at the same time nearly all points are in a box of side-length $O(\frac{\ln d}{d-1})$. The answer is to remember that points on the surface of the ball satisfy $x_1^2 + x_2^2 + \dots + x_d^2 = 1$, so for each coordinate i , a typical value will be $\pm O(\frac{1}{\sqrt{d}})$. In fact, it is often helpful to think of picking a random point on the sphere as very similar to picking a random point of the form $(\pm \frac{1}{\sqrt{d}}, \pm \frac{1}{\sqrt{d}}, \pm \frac{1}{\sqrt{d}}, \dots, \pm \frac{1}{\sqrt{d}})$.

2.5 Generating Points Uniformly at Random from a Ball

Consider generating points uniformly at random on the surface of the unit ball. For the 2-dimensional version of generating points on the circumference of a unit-radius circle, independently generate each coordinate uniformly at random from the interval $[-1, 1]$. This produces points distributed over a square that is large enough to completely contain the unit circle. Project each point onto the unit circle. The distribution is not uniform since more points fall on a line from the origin to a vertex of the square than fall on a line from the origin to the midpoint of an edge of the square due to the difference in length. To solve this problem, discard all points outside the unit circle and project the remaining points onto the circle.

In higher dimensions, this method does not work since the fraction of points that fall inside the ball drops to zero and all of the points would be thrown away. The solution is to generate a point each of whose coordinates is an independent Gaussian variable. Generate x_1, x_2, \dots, x_d , using a zero mean, unit variance Gaussian, namely, $\frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$ on the

real line.² Thus, the probability density of \mathbf{x} is

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}}} e^{-\frac{x_1^2 + x_2^2 + \dots + x_d^2}{2}}$$

and is spherically symmetric. Normalizing the vector $\mathbf{x} = (x_1, x_2, \dots, x_d)$ to a unit vector, namely $\frac{\mathbf{x}}{|\mathbf{x}|}$, gives a distribution that is uniform over the surface of the sphere. Note that once the vector is normalized, its coordinates are no longer statistically independent.

To generate a point \mathbf{y} uniformly over the ball (surface and interior), scale the point $\frac{\mathbf{x}}{|\mathbf{x}|}$ generated on the surface by a scalar $\rho \in [0, 1]$. What should the distribution of ρ be as a function of r ? It is certainly not uniform, even in 2 dimensions. Indeed, the density of ρ at r is proportional to r for $d = 2$. For $d = 3$, it is proportional to r^2 . By similar reasoning, the density of ρ at distance r is proportional to r^{d-1} in d dimensions. Solving $\int_{r=0}^{r=1} cr^{d-1} dr = 1$ (the integral of density must equal 1) one should set $c = d$. Another way to see this formally is that the volume of the radius r ball in d dimensions is $r^d V(d)$. The density at radius r is exactly $\frac{d}{dr}(r^d V_d) = dr^{d-1} V_d$. So, pick $\rho(r)$ with density equal to dr^{d-1} for r over $[0, 1]$.

We have succeeded in generating a point

$$\mathbf{y} = \rho \frac{\mathbf{x}}{|\mathbf{x}|}$$

uniformly at random from the unit ball by using the convenient spherical Gaussian distribution. In the next sections, we will analyze the spherical Gaussian in more detail.

2.6 Gaussians in High Dimension

A 1-dimensional Gaussian has its mass close to the origin. However, as the dimension is increased something different happens. The d -dimensional spherical Gaussian with zero mean and variance σ^2 in each coordinate has density function

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right).$$

The value of the density is maximum at the origin, but there is very little volume there. When $\sigma^2 = 1$, integrating the probability density over a unit ball centered at the origin yields almost zero mass since the volume of such a ball is negligible. In fact, one needs

²One might naturally ask: “how do you generate a random number from a 1-dimensional Gaussian?” To generate a number from any distribution given its cumulative distribution function P , first select a uniform random number $u \in [0, 1]$ and then choose $x = P^{-1}(u)$. For any $a < b$, the probability that x is between a and b is equal to the probability that u is between $P(a)$ and $P(b)$ which equals $P(b) - P(a)$ as desired. For the 2-dimensional Gaussian, one can generate a point in polar coordinates by choosing angle θ uniform in $[0, 2\pi]$ and radius $r = \sqrt{-2 \ln(u)}$ where u is uniform random in $[0, 1]$. This is called the Box-Muller transform.

to increase the radius of the ball to nearly \sqrt{d} before there is a significant volume and hence significant probability mass. If one increases the radius much beyond \sqrt{d} , the integral barely increases even though the volume increases since the probability density is dropping off at a much higher rate. The following theorem formally states that nearly all the probability is concentrated in a thin annulus of width $O(1)$ at radius \sqrt{d} .

Theorem 2.9 (Gaussian Annulus Theorem) *For a d -dimensional spherical Gaussian with unit variance in each direction, for any $\beta \leq \sqrt{d}$, all but at most $3e^{-c\beta^2}$ of the probability mass lies within the annulus $\sqrt{d} - \beta \leq |\mathbf{x}| \leq \sqrt{d} + \beta$, where c is a fixed positive constant.*

For a high-level intuition, note that $E(|\mathbf{x}|^2) = \sum_{i=1}^d E(x_i^2) = dE(x_1^2) = d$, so the mean squared distance of a point from the center is d . The Gaussian Annulus Theorem says that the points are tightly concentrated. We call the square root of the mean squared distance, namely \sqrt{d} , the radius of the Gaussian.

To prove the Gaussian Annulus Theorem we make use of a tail inequality for sums of independent random variables of bounded moments (Theorem 12.5).

Proof (Gaussian Annulus Theorem): Let $\mathbf{x} = (x_1, x_2, \dots, x_d)$ be a point selected from a unit variance Gaussian centered at the origin, and let $r = |\mathbf{x}|$. $\sqrt{d} - \beta \leq |\mathbf{y}| \leq \sqrt{d} + \beta$ is equivalent to $|r - \sqrt{d}| \geq \beta$. If $|r - \sqrt{d}| \geq \beta$, then multiplying both sides by $r + \sqrt{d}$ gives $|r^2 - d| \geq \beta(r + \sqrt{d}) \geq \beta\sqrt{d}$. So, it suffices to bound the probability that $|r^2 - d| \geq \beta\sqrt{d}$.

Rewrite $r^2 - d = (x_1^2 + \dots + x_d^2) - d = (x_1^2 - 1) + \dots + (x_d^2 - 1)$ and perform a change of variables: $y_i = x_i^2 - 1$. We want to bound the probability that $|y_1 + \dots + y_d| \geq \beta\sqrt{d}$. Notice that $E(y_i) = E(x_i^2) - 1 = 0$. To apply Theorem 12.5, we need to bound the s^{th} moments of y_i .

For $|x_i| \leq 1$, $|y_i|^s \leq 1$ and for $|x_i| \geq 1$, $|y_i|^s \leq |x_i|^{2s}$. Thus

$$\begin{aligned} |E(y_i^s)| &= E(|y_i|^s) \leq E(1 + x_i^{2s}) = 1 + E(x_i^{2s}) \\ &= 1 + \sqrt{\frac{2}{\pi}} \int_0^\infty x^{2s} e^{-x^2/2} dx \end{aligned}$$

Using the substitution $2z = x^2$,

$$\begin{aligned} |E(y_i^s)| &= 1 + \frac{1}{\sqrt{\pi}} \int_0^\infty 2^s z^{s-(1/2)} e^{-z} dz \\ &\leq 2^s s!. \end{aligned}$$

The last inequality is from the Gamma integral.

Since $E(y_i) = 0$, $Var(y_i) = E(y_i^2) \leq 2^2 \cdot 2 = 8$. Unfortunately, we do not have $|E(y_i^s)| \leq 8s!$ as required in Theorem 12.5. To fix this problem, perform one more change of variables, using $w_i = y_i/2$. Then, $Var(w_i) \leq 2$ and $|E(w_i^s)| \leq 2s!$, and our goal is now to bound the probability that $|w_1 + \dots + w_d| \geq \frac{\beta\sqrt{d}}{2}$. Applying Theorem 12.5 where $\sigma^2 = 2$ and $n = d$, this occurs with probability less than or equal to $3e^{-\frac{\beta^2}{96}}$. ■

In the next sections we will see several uses of the Gaussian Annulus Theorem.

2.7 Random Projection and Johnson-Lindenstrauss Lemma

One of the most frequently used subroutines in tasks involving high dimensional data is nearest neighbor search. In nearest neighbor search we are given a database of n points in \mathbf{R}^d where n and d are usually large. The database can be preprocessed and stored in an efficient data structure. Thereafter, we are presented “query” points in \mathbf{R}^d and are asked to find the nearest or approximately nearest database point to the query point. Since the number of queries is often large, the time to answer each query should be very small, ideally a small function of $\log n$ and $\log d$, whereas preprocessing time could be larger, namely a polynomial function of n and d . For this and other problems, dimension reduction, where one projects the database points to a k -dimensional space with $k \ll d$ (usually dependent on $\log d$) can be very useful so long as the relative distances between points are approximately preserved. We will see using the Gaussian Annulus Theorem that such a projection indeed exists and is simple.

The projection $f : \mathbf{R}^d \rightarrow \mathbf{R}^k$ that we will examine (many related projections are known to work as well) is the following. Pick k Gaussian vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ in \mathbf{R}^d with unit-variance coordinates. For any vector \mathbf{v} , define the projection $f(\mathbf{v})$ by:

$$f(\mathbf{v}) = (\mathbf{u}_1 \cdot \mathbf{v}, \mathbf{u}_2 \cdot \mathbf{v}, \dots, \mathbf{u}_k \cdot \mathbf{v}).$$

The projection $f(\mathbf{v})$ is the vector of dot products of \mathbf{v} with the \mathbf{u}_i . We will show that with high probability, $|f(\mathbf{v})| \approx \sqrt{k}|\mathbf{v}|$. For any two vectors \mathbf{v}_1 and \mathbf{v}_2 , $f(\mathbf{v}_1 - \mathbf{v}_2) = f(\mathbf{v}_1) - f(\mathbf{v}_2)$. Thus, to estimate the distance $|\mathbf{v}_1 - \mathbf{v}_2|$ between two vectors \mathbf{v}_1 and \mathbf{v}_2 in \mathbf{R}^d , it suffices to compute $|f(\mathbf{v}_1) - f(\mathbf{v}_2)| = |f(\mathbf{v}_1 - \mathbf{v}_2)|$ in the k -dimensional space since the factor of \sqrt{k} is known and one can divide by it. The reason distances increase when we project to a lower dimensional space is that the vectors \mathbf{u}_i are not unit length. Also notice that the vectors \mathbf{u}_i are not orthogonal. If we had required them to be orthogonal, we would have lost statistical independence.

Theorem 2.10 (The Random Projection Theorem) *Let \mathbf{v} be a fixed vector in \mathbf{R}^d and let f be defined as above. There exists constant $c > 0$ such that for $\varepsilon \in (0, 1)$,*

$$\text{Prob} \left(\left| |f(\mathbf{v})| - \sqrt{k}|\mathbf{v}| \right| \geq \varepsilon \sqrt{k}|\mathbf{v}| \right) \leq 3e^{-ck\varepsilon^2},$$

where the probability is taken over the random draws of vectors \mathbf{u}_i used to construct f .

Proof: By scaling both sides of the inner inequality by $|\mathbf{v}|$, we may assume that $|\mathbf{v}| = 1$. The sum of independent normally distributed real variables is also normally distributed where the mean and variance are the sums of the individual means and variances. Since $\mathbf{u}_i \cdot \mathbf{v} = \sum_{j=1}^d u_{ij}v_j$, the random variable $\mathbf{u}_i \cdot \mathbf{v}$ has Gaussian density with zero mean and unit variance, in particular,

$$\text{Var}(\mathbf{u}_i \cdot \mathbf{v}) = \text{Var}\left(\sum_{j=1}^d v_{ij}v_j\right) = \sum_{j=1}^d v_j^2 \text{Var}(u_{ij}) = \sum_{j=1}^d v_j^2 = 1$$

Since $\mathbf{u}_1 \cdot \mathbf{v}, \mathbf{u}_2 \cdot \mathbf{v}, \dots, \mathbf{u}_k \cdot \mathbf{v}$ are independent Gaussian random variables, $f(\mathbf{v})$ is a random vector from a k -dimensional spherical Gaussian with unit variance in each coordinate, and so the theorem follows from the Gaussian Annulus Theorem (Theorem 2.9) with d replaced by k . \blacksquare

The random projection theorem establishes that the probability of the length of the projection of a single vector differing significantly from its expected value is exponentially small in k , the dimension of the target subspace. By a union bound, the probability that any of $O(n^2)$ pairwise differences $|\mathbf{v}_i - \mathbf{v}_j|$ among n vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ differs significantly from their expected values is small, provided $k \geq \frac{3}{c\varepsilon^2} \ln n$. Thus, this random projection preserves all relative pairwise distances between points in a set of n points with high probability. This is the content of the Johnson-Lindenstrauss Lemma.

Theorem 2.11 (Johnson-Lindenstrauss Lemma) *For any $0 < \varepsilon < 1$ and any integer n , let $k \geq \frac{3}{c\varepsilon^2} \ln n$ with c as in Theorem 2.9. For any set of n points in R^d , the random projection $f : R^d \rightarrow R^k$ defined above has the property that for all pairs of points \mathbf{v}_i and \mathbf{v}_j , with probability at least $1 - 3/2n$,*

$$(1 - \varepsilon)\sqrt{k}|\mathbf{v}_i - \mathbf{v}_j| \leq |f(\mathbf{v}_i) - f(\mathbf{v}_j)| \leq (1 + \varepsilon)\sqrt{k}|\mathbf{v}_i - \mathbf{v}_j|.$$

Proof: Applying the Random Projection Theorem (Theorem 2.10), for any fixed \mathbf{v}_i and \mathbf{v}_j , the probability that $|f(\mathbf{v}_i - \mathbf{v}_j)|$ is outside the range

$$\left[(1 - \varepsilon)\sqrt{k}|\mathbf{v}_i - \mathbf{v}_j|, (1 + \varepsilon)\sqrt{k}|\mathbf{v}_i - \mathbf{v}_j| \right]$$

is at most $3e^{-ck\varepsilon^2} \leq 3/n^3$ for $k \geq \frac{3 \ln n}{c\varepsilon^2}$. Since there are $\binom{n}{2} < n^2/2$ pairs of points, by the union bound, the probability that any pair has a large distortion is less than $\frac{3}{2n}$. \blacksquare

Remark: It is important to note that the conclusion of Theorem 2.11 asserts for all \mathbf{v}_i and \mathbf{v}_j , not just for most of them. The weaker assertion for most \mathbf{v}_i and \mathbf{v}_j is typically less useful, since our algorithm for a problem such as nearest-neighbor search might return one of the bad pairs of points. A remarkable aspect of the theorem is that the number of dimensions in the projection is only dependent logarithmically on n . Since k is often much less than d , this is called a dimension reduction technique. In applications, the dominant term is typically the $1/\varepsilon^2$ term.

For the nearest neighbor problem, if the database has n_1 points and n_2 queries are expected during the lifetime of the algorithm, take $n = n_1 + n_2$ and project the database to a random k -dimensional space, for k as in Theorem 2.11. On receiving a query, project the query to the same subspace and compute nearby database points. The Johnson Lindenstrauss Lemma says that with high probability this will yield the right answer whatever the query. Note that the exponentially small in k probability was useful here in making k only dependent on $\ln n$, rather than n .

2.8 Separating Gaussians

Mixtures of Gaussians are often used to model heterogeneous data coming from multiple sources. For example, suppose we are recording the heights of individuals age 20-30 in a city. We know that on average, men tend to be taller than women, so a natural model would be a Gaussian mixture model $p(x) = w_1 p_1(x) + w_2 p_2(x)$, where $p_1(x)$ is a Gaussian density representing the typical heights of women, $p_2(x)$ is a Gaussian density representing the typical heights of men, and w_1 and w_2 are the *mixture weights* representing the proportion of women and men in the city. The *parameter estimation problem* for a mixture model is the problem: given access to samples from the overall density p (e.g., heights of people in the city, but without being told whether the person with that height is male or female), reconstruct the parameters for the distribution (e.g., good approximations to the means and variances of p_1 and p_2 , as well as the mixture weights).

There are taller women and shorter men, so even if one solved the parameter estimation problem for heights perfectly, given a data point, one couldn't necessarily tell which population it came from. That is, given a height, one couldn't necessarily tell if it came from a man or a woman. In this section, we will look at a problem that is in some ways easier and some ways harder than this problem of heights. It will be harder in that we will be interested in a mixture of two Gaussians in high-dimensions as opposed to the $d = 1$ case of heights. But it will be easier in that we will assume the means are quite well-separated compared to the variances. Specifically, our focus will be on a mixture of two spherical unit-variance Gaussians whose means are separated by a distance $\Omega(d^{1/4})$. We will show that at this level of separation, we can with high probability uniquely determine which Gaussian each data point came from. The algorithm to do so will actually be quite simple. Calculate the distance between all pairs of points. Points whose distance apart is smaller are from the same Gaussian, whereas points whose distance is larger are from different Gaussians. Later, we will see that with more sophisticated algorithms, even a separation of $\Omega(1)$ suffices.

First, consider just one spherical unit-variance Gaussian centered at the origin. From Theorem 2.9, most of its probability mass lies on an annulus of width $O(1)$ at radius \sqrt{d} . Also $e^{-|\mathbf{x}|^2/2} = \prod_i e^{-x_i^2/2}$ and almost all of the mass is within the slab $\{\mathbf{x} \mid -c \leq x_1 \leq c\}$, for $c \in O(1)$. Pick a point \mathbf{x} from this Gaussian. After picking \mathbf{x} , rotate the coordinate system to make the first axis align with \mathbf{x} . Independently pick a second point \mathbf{y} from

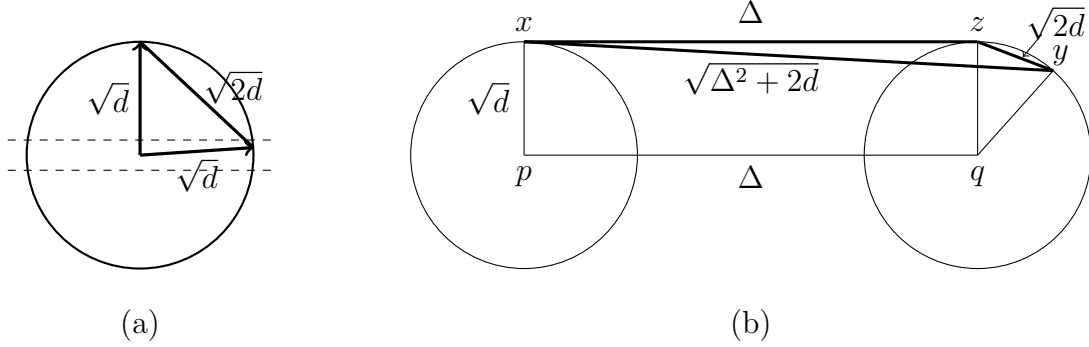


Figure 2.5: (a) indicates that two randomly chosen points in high dimension are surely almost nearly orthogonal. (b) indicates the distance between a pair of random points from two different unit balls approximating the annuli of two Gaussians.

this Gaussian. The fact that almost all of the probability mass of the Gaussian is within the slab $\{\mathbf{x} \mid -c \leq x_1 \leq c, c \in O(1)\}$ at the equator implies that \mathbf{y} 's component along \mathbf{x} 's direction is $O(1)$ with high probability. Thus, \mathbf{y} is nearly perpendicular to \mathbf{x} . So, $|\mathbf{x} - \mathbf{y}| \approx \sqrt{|\mathbf{x}|^2 + |\mathbf{y}|^2}$. See Figure 2.5(a). More precisely, since the coordinate system has been rotated so that \mathbf{x} is at the North Pole, $\mathbf{x} = (\sqrt{d} \pm O(1), 0, \dots, 0)$. Since \mathbf{y} is almost on the equator, further rotate the coordinate system so that the component of \mathbf{y} that is perpendicular to the axis of the North Pole is in the second coordinate. Then $\mathbf{y} = (O(1), \sqrt{d} \pm O(1), 0, \dots, 0)$. Thus,

$$(\mathbf{x} - \mathbf{y})^2 = d \pm O(\sqrt{d}) + d \pm O(\sqrt{d}) = 2d \pm O(\sqrt{d})$$

and $|\mathbf{x} - \mathbf{y}| = \sqrt{2d} \pm O(1)$ with high probability.

Consider two spherical unit variance Gaussians with centers \mathbf{p} and \mathbf{q} separated by a distance Δ . The distance between a randomly chosen point \mathbf{x} from the first Gaussian and a randomly chosen point \mathbf{y} from the second is close to $\sqrt{\Delta^2 + 2d}$, since $\mathbf{x} - \mathbf{p}$, $\mathbf{p} - \mathbf{q}$, and $\mathbf{q} - \mathbf{y}$ are nearly mutually perpendicular. Pick \mathbf{x} and rotate the coordinate system so that \mathbf{x} is at the North Pole. Let \mathbf{z} be the North Pole of the ball approximating the second Gaussian. Now pick \mathbf{y} . Most of the mass of the second Gaussian is within $O(1)$ of the equator perpendicular to $\mathbf{z} - \mathbf{q}$. Also, most of the mass of each Gaussian is within distance $O(1)$ of the respective equators perpendicular to the line $\mathbf{q} - \mathbf{p}$. See Figure 2.5 (b). Thus,

$$\begin{aligned} |\mathbf{x} - \mathbf{y}|^2 &\approx \Delta^2 + |\mathbf{z} - \mathbf{q}|^2 + |\mathbf{q} - \mathbf{y}|^2 \\ &= \Delta^2 + 2d \pm O(\sqrt{d}). \end{aligned}$$

To ensure that the distance between two points picked from the same Gaussian are closer to each other than two points picked from different Gaussians requires that the upper limit of the distance between a pair of points from the same Gaussian is at most

the lower limit of distance between points from different Gaussians. This requires that $\sqrt{2d} + O(1) \leq \sqrt{2d + \Delta^2} - O(1)$ or $2d + O(\sqrt{d}) \leq 2d + \Delta^2$, which holds when $\Delta \in \omega(d^{1/4})$. Thus, mixtures of spherical Gaussians can be separated in this way, provided their centers are separated by $\omega(d^{1/4})$. If we have n points and want to correctly separate all of them with high probability, we need our individual high-probability statements to hold with probability $1 - 1/\text{poly}(n)$,³ which means our $O(1)$ terms from Theorem 2.9 become $O(\sqrt{\log n})$. So we need to include an extra $O(\sqrt{\log n})$ term in the separation distance.

Algorithm for separating points from two Gaussians: Calculate all pairwise distances between points. The cluster of smallest pairwise distances must come from a single Gaussian. Remove these points. The remaining points come from the second Gaussian.

One can actually separate Gaussians where the centers are much closer. In the next chapter we will use singular value decomposition to separate points from a mixture of two Gaussians when their centers are separated by a distance $O(1)$.

2.9 Fitting a Spherical Gaussian to Data

Given a set of sample points, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, in a d -dimensional space, we wish to find the spherical Gaussian that best fits the points. Let f be the unknown Gaussian with mean $\boldsymbol{\mu}$ and variance σ^2 in each direction. The probability density for picking these points when sampling according to f is given by

$$c \exp \left(- \frac{(\mathbf{x}_1 - \boldsymbol{\mu})^2 + (\mathbf{x}_2 - \boldsymbol{\mu})^2 + \dots + (\mathbf{x}_n - \boldsymbol{\mu})^2}{2\sigma^2} \right)$$

where the normalizing constant c is the reciprocal of $\left[\int e^{-\frac{|\mathbf{x} - \boldsymbol{\mu}|^2}{2\sigma^2}} d\mathbf{x} \right]^n$. In integrating from $-\infty$ to ∞ , one can shift the origin to $\boldsymbol{\mu}$ and thus c is $\left[\int e^{-\frac{|\mathbf{x}|^2}{2\sigma^2}} d\mathbf{x} \right]^{-n} = \frac{1}{(2\pi)^{\frac{n}{2}}}$ and is independent of $\boldsymbol{\mu}$.

The *Maximum Likelihood Estimator* (MLE) of f , given the samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, is the f that maximizes the above probability density.

Lemma 2.12 *Let $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ be a set of n d -dimensional points. Then $(\mathbf{x}_1 - \boldsymbol{\mu})^2 + (\mathbf{x}_2 - \boldsymbol{\mu})^2 + \dots + (\mathbf{x}_n - \boldsymbol{\mu})^2$ is minimized when $\boldsymbol{\mu}$ is the centroid of the points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, namely $\boldsymbol{\mu} = \frac{1}{n}(\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_n)$.*

Proof: Setting the gradient of $(\mathbf{x}_1 - \boldsymbol{\mu})^2 + (\mathbf{x}_2 - \boldsymbol{\mu})^2 + \dots + (\mathbf{x}_n - \boldsymbol{\mu})^2$ with respect to $\boldsymbol{\mu}$ to zero yields

$$-2(\mathbf{x}_1 - \boldsymbol{\mu}) - 2(\mathbf{x}_2 - \boldsymbol{\mu}) - \dots - 2(\mathbf{x}_n - \boldsymbol{\mu}) = 0.$$

Solving for $\boldsymbol{\mu}$ gives $\boldsymbol{\mu} = \frac{1}{n}(\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_n)$. ■

³poly(n) means bounded by a polynomial in n .

To determine the maximum likelihood estimate of σ^2 for f , set $\boldsymbol{\mu}$ to the true centroid. Next, show that σ is set to the standard deviation of the sample. Substitute $\nu = \frac{1}{2\sigma^2}$ and $a = (\mathbf{x}_1 - \boldsymbol{\mu})^2 + (\mathbf{x}_2 - \boldsymbol{\mu})^2 + \cdots + (\mathbf{x}_n - \boldsymbol{\mu})^2$ into the formula for the probability of picking the points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. This gives

$$\frac{e^{-a\nu}}{\left[\int_x e^{-x^2\nu} dx \right]^n}.$$

Now, a is fixed and ν is to be determined. Taking logs, the expression to maximize is

$$-a\nu - n \ln \left[\int_x e^{-\nu x^2} dx \right].$$

To find the maximum, differentiate with respect to ν , set the derivative to zero, and solve for σ . The derivative is

$$-a + n \frac{\int_x |x|^2 e^{-\nu x^2} dx}{\int_x e^{-\nu x^2} dx}.$$

Setting $y = |\sqrt{\nu}\mathbf{x}|$ in the derivative, yields

$$-a + \frac{n}{\nu} \frac{\int_y y^2 e^{-y^2} dy}{\int_y e^{-y^2} dy}.$$

Since the ratio of the two integrals is the expected distance squared of a d -dimensional spherical Gaussian of standard deviation $\frac{1}{\sqrt{2}}$ to its center, and this is known to be $\frac{d}{2}$, we get $-a + \frac{nd}{2\nu}$. Substituting σ^2 for $\frac{1}{2\nu}$ gives $-a + nd\sigma^2$. Setting $-a + nd\sigma^2 = 0$ shows that the maximum occurs when $\sigma = \frac{\sqrt{a}}{\sqrt{nd}}$. Note that this quantity is the square root of the average coordinate distance squared of the samples to their mean, which is the standard deviation of the sample. Thus, we get the following lemma.

Lemma 2.13 *The maximum likelihood spherical Gaussian for a set of samples is the Gaussian with center equal to the sample mean and standard deviation equal to the standard deviation of the sample from the true mean.*

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be a sample of points generated by a Gaussian probability distribution. Then $\boldsymbol{\mu} = \frac{1}{n}(\mathbf{x}_1 + \mathbf{x}_2 + \cdots + \mathbf{x}_n)$ is an unbiased estimator of the expected value of the distribution. However, if in estimating the variance from the sample set, we use the estimate of the expected value rather than the true expected value, we will not get an unbiased estimate of the variance, since the sample mean is not independent of the sample set. One should use $\tilde{\boldsymbol{\mu}} = \frac{1}{n-1}(\mathbf{x}_1 + \mathbf{x}_2 + \cdots + \mathbf{x}_n)$ when estimating the variance. See Section 12.5.10 of the appendix.