

## Personal Notes



# Contents

<b>1</b>	<b>Linear Algebra</b>	<b>5</b>
1.1	Singular Value Decomposition . . . . .	5
1.2	Pseudo Inverse . . . . .	6
1.3	Low-Rank Matrix Approximation . . . . .	6
1.4	Important Matrix Decompositions . . . . .	7
1.5	Exercises . . . . .	7
<b>2</b>	<b>Calculus and Analysis For Machine Learning</b>	<b>9</b>
<b>3</b>	<b>Least Squares and Linear Regression</b>	<b>11</b>
3.1	Least Squares: Linear Algebra Perspective . . . . .	11
3.1.1	Standard Least Squares: $m > n$ . . . . .	11
3.1.2	Other case: $m < n$ . . . . .	12
3.1.3	Least Squares and the Pseudo Inverse . . . . .	12
3.2	Gauss-Markov Theorem: Least Squares and MLE . . . . .	13
3.3	Ridge Regression . . . . .	13
3.4	LASSO and Elastic Nets . . . . .	13
3.5	Bayesian Linear Regression . . . . .	13
3.6	LASSO . . . . .	15
<b>4</b>	<b>PCA</b>	<b>17</b>
<b>5</b>	<b>Other fun topics/one off topics</b>	<b>19</b>
<b>6</b>	<b>Other notes/Daily notes</b>	<b>21</b>
6.1	Gaussian Distribution: Normalizing Constant 5/24/24 . . . . .	21
6.2	Polar Change of Coordinates . . . . .	21
6.3	Change of Variables . . . . .	22
6.4	High dimensional spaces . . . . .	29
6.4.1	Volume of $d$ -ball . . . . .	29
6.4.2	Balls and Cubes . . . . .	30
6.4.3	Most of the mass is near the equator . . . . .	30
6.4.4	Most of the volume is near the surface . . . . .	30
6.4.5	High dimensional cubes . . . . .	31
6.4.6	Distances between randomly sampled points . . . . .	31
6.4.7	Digression: $L_p$ balls . . . . .	31
6.4.8	Discussion . . . . .	31
6.5	High Dimensional Gaussian . . . . .	32
6.5.1	Radial density of Gaussian . . . . .	32



# Chapter 1

## Linear Algebra

### 1.1 Singular Value Decomposition

Being able to diagonalize a matrix makes a lot of matrix manipulations a lot easier. Life becomes particularly easy if the transform used to diagonalize a matrix is orthogonal. The SVD provides an orthonormal diagonal decomposition for arbitrary (non-square) matrices. The catch: the basis vectors used for the domain/range can be different. So the SVD isn't necessarily suitable for applications with powers of matrices since you don't get  $A^2 = SAS^T SAS^T = SA^2 S^T$ . But in a lot of other applications, it's very useful.

Let  $A \in \mathbb{R}^{m \times n}$  and let  $r = \text{rank } A$ . We want orthonormal bases  $v_1, \dots, v_r$  for the row space and  $u_1, \dots, u_r$  for the column space such that

$$Av_i = \sigma_i u_i$$

for some  $\sigma_i \in \mathbb{R}$ ,  $i = 1, \dots, r$ . What is a good candidate for  $(v_i)_{i=1}^r$  and  $(u_i)_{i=1}^r$ ? How about the eigenvectors of  $A^T A \in \mathbb{R}^{n \times n}$  and  $AA^T \in \mathbb{R}^{m \times m}$ ? These matrices are symmetric positive semidefinite. To verify PSD, note that  $x^T A^T A x = \|A^T x\|^2 \geq 0 \forall x$ . Because these matrices are symmetric, they have an orthogonal set of eigenvectors. This is an immediate consequence of the spectral theorem—see Exercise 1.5. Let  $(v_i)_{i=1}^r$  be orthonormal eigenvectors of  $A^T A$  with associated eigenvalues given by  $\sigma_i^2$ . Note that

$$A^T A v_i = \sigma_i^2 v_i \implies v_i^T A^T A v_i = \sigma_i^2 v_i^T v_i \implies \|A v_i\|^2 = \sigma_i^2.$$

Now, let  $u_i := \frac{1}{\sigma_i} A v_i$ . By the previous display,  $u_i$  is a unit vector. Ideally, we would like the set of  $Av_i$  to map to an orthonormal basis in the column space. Let's see what happens. We have

$$A^T A v_i = \sigma_i^2 v_i \iff AA^T A v_i = \sigma_i^2 A v_i \iff AA^T u_i = \sigma_i^2 u_i.$$

So,  $u_i = \frac{1}{\sigma_i} A v_i$  is a unit eigenvector of  $AA^T$ . Thus, we have

- $(u_i)_{i=1}^r$  is an orthonormal basis for the column space
- $(v_i)_{i=1}^r$  is an orthonormal basis for the row space
- $Av_i = \sigma_i u_i$
- $\sigma_i \geq 0$ ,  $\sigma_i = \sqrt{\lambda_i}$  where  $\lambda_i \in \sigma(A^T A)$

Now, complete each basis with orthonormal basis for the the nullspace and left nullspace of  $A$  to get  $(v_1, \dots, v_r, \dots, v_n)$  and  $(u_1, \dots, u_r, \dots, u_m)$ . Let  $V \in \mathbb{R}^{n \times n}$  be formed columnwise from  $(v_i)_i$  and likewise form  $U \in \mathbb{R}^{m \times m}$  columnwise from  $(u_i)_i$ . Let  $\Sigma \in \mathbb{R}^{m \times n}$  be formed by placing  $\sigma_i$  on the diagonal. Of course, do this all with a consistent ordering so eigenvectors and eigenvalues in each matrix match. We have

$$Av_i = \sigma_i u_i \iff AV = U\Sigma \iff \underbrace{A = U\Sigma V^T}_{\text{SVD}}$$

This last expression is the SVD. The columns of  $U$  and  $V$  are known, respectively, as the left and right singular vectors of  $A$  and  $(\sigma_i)$  are the singular values.

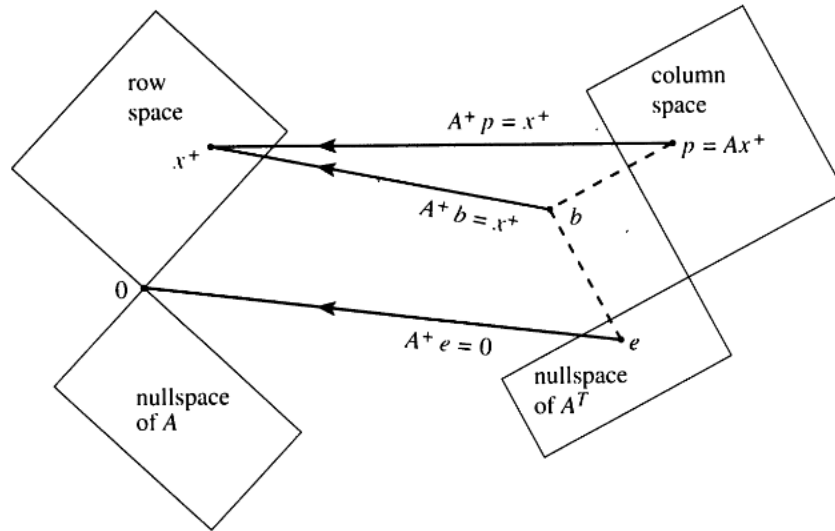


Figure 1.1: Action of pseudo-inverse (taken from Strang paper).

## 1.2 Pseudo Inverse

For a linear mapping  $A$ , row  $A$  and col  $A$  have the same dimension (see Exercise 1.5). This seems like a remarkable fact. The linear mapping  $A$  is a bijection between row  $A$  and col  $A$  (see Exercise 1.5 for the full rank case).<sup>1</sup> The pseudo-inverse, denoted by  $A^+$ , maps from col  $A$  back to row  $A$ . That is, for an element  $\mathbf{y}$  in col  $A$ ,  $A^+\mathbf{y}$  recovers the unique  $\mathbf{x}$  in row  $A$  such that  $A\mathbf{x} = \mathbf{y}$ . If  $\mathbf{y} \in (\text{col } A)^\perp = N(A^\top)$ , then  $A^+\mathbf{y} = 0$ . The pseudo-inverse coincides with  $A^{-1}$  when it exists.

Figure ??, taken from a Gilbert Strang paper, clearly illustrates the action of  $A^+$ .

The matrix  $A^+$  satisfying these properties can be clearly derived from the SVD. Let rank  $A = r$ , and let  $(u_i)_{i=1}^r$  and  $(v_i)_{i=1}^r$  be as defined in Section 1.1, so that these are orthonormal bases of the column and row spaces respectively. We have  $Av_i = \sigma u_i$ . We would like an inverse satisfying

$$A^+ \sigma u_i = v_i \iff A^+ u_i = \frac{1}{\sigma_i} v_i \iff A^+ U = V \Sigma^+$$

where  $\Sigma^+$  is the same dimensions as  $\Sigma^\top$ , but with  $1/\sigma_i$  on the diagonal where  $\sigma_i \neq 0$  holds by definition since we only define  $\sigma_i$  up to  $r = \text{rank } A$ . Hence, the pseudo-inverse is given by

$$A^+ = V \Sigma^+ U^\top$$

**TODO:** There are a couple of important subtleties I'm missing here. What to do with zero singular values? What's the dimension of  $\Sigma$ . Just lock that down and be consistent. Also, order  $\sigma_1 \geq \dots \geq \sigma_r > 0$

## 1.3 Low-Rank Matrix Approximation

**Problem Statement** Let  $A \in \mathbb{R}^{m \times n}$  be a real matrix with  $m < n$  and let  $d$  be a positive integer that will be used to represent a rank constraint. We would like to find a matrix  $A_d$  satisfying

$$\begin{aligned} A_d &\in \arg \min_{B \in \mathbb{R}^{m \times n}} \|A - B\| \\ \text{s.t. } \text{rank } B &\leq d. \end{aligned}$$

Note that we have not been specific about the precise norm used. The solution can vary based on the norm used (I think). But we'll see that the optimal solution follows readily from the SVD when we use the Frobenius or spectral norm.

<sup>1</sup>**TODO:** Full case?

**Theorem 1.** Let  $A = U\Sigma V^\top$  be the SVD of  $A$ . Let  $U_d \in \mathbb{R}^{m \times d}$  and  $V_d \in \mathbb{R}^{n \times d}$  be the matrices constructed from the first  $d$  columns of  $U$  and  $V$ , where we assume the convention that  $\Sigma$  is ordered as in ???. Then

$$A_d := U_r \Sigma_d V_d^\top$$

is an optimal rank  $d$  approximation in the sense that

$$\|A - A_d\|_F \leq \|A - B\|_F \quad \text{and} \quad \|A - A_d\|_2 \leq \|A - B\|_2 \quad \forall B \text{ s.t. } \text{rank } B \leq d.$$

**TODO** Add in proof for Frobenius norm. Can use Weyl's inequality. See [here](#) and link to Tao notes therein. Also, might be nice to spend some time looking at variational expression for eigenvalues and applications. Or can use a different proof from [here](#).

## 1.4 Important Matrix Decompositions

**TODO:** Cholesky, QR. Applications to numerical linear algebra.

## 1.5 Exercises

### Exercise 1:

Let  $A \in M_{m,n}$ . Show  $\text{rank } A = \text{rank } A^\top$ .

### Exercise 2:

Show  $\text{row } A \perp N(A)$ .

### Exercise 3:

**Rank-nullity theorem.** For  $A \in \mathbb{R}^{m \times n}$ , show that  $\text{rank } A + \dim N(A) = n$ .

### Exercise 4:

**Spectral theorem.** Show that if  $A$  is Hermitian (or just real symmetric for an easier version) that  $A$  may be decomposed as  $A = Q\Lambda Q^\top$ , where  $Q$  is orthonormal and  $\Lambda$  is diagonal.

### Exercise 5:

Show that any symmetric matrix is diagonalizable.

### Exercise 6:

A matrix  $A$  is said to be diagonalizable if it is similar to a diagonal matrix. Show that a matrix  $A$  is diagonalizable if and only if there is a set of  $n$  linearly independent vectors, each of which is an eigenvector of  $A$ .

### Exercise 7:

Singular values vs eigenvalues for symmetric matrices. Show the following: Suppose  $A$  is a symmetric matrix. If  $\lambda \in \sigma(A)$ , then  $\lambda^2 \in \sigma(A^\top A)$ , and, in particular,  $|\lambda|$  is a singular value of  $A$ .

Corollary: If  $A$  is PSD, then the set of eigenvectors and singular values coincide.

### Exercise 8:

Prove  $\sigma(A) = \sigma(A^\top)$ .

### Exercise 9:

Suppose that  $\hat{y} \in \text{col } X$ . Show there exists a unique  $\hat{w}$  solving  $X\hat{w} = \hat{y}$  (\*\*). Prove above claim and compute  $\hat{w}$  in terms of  $X$  and  $y$ .<sup>2</sup>

<sup>2</sup>Note that this implies that, restricted to  $\text{row } X$ ,  $X$  is injective to  $\text{col } X$ . It's easy to show surjectivity to  $\text{col } X$ , since  $\text{col } X$  is the range of  $X$ . For any  $y \in \text{col } X$ , there exists a  $w$  s.t.  $Xw = y$ . Use the fact that  $\text{row } X \perp N(X)$  to throw away the part of  $w \in N(X)$ . Now you have a point  $w$  in  $\text{row } X$  that maps to  $y$ .





## Chapter 2

# Calculus and Analysis For Machine Learning

Possible sections:

- Some basic calculus results (MVT, extreme value theorem, etc.)
- Fundamental theorem of calculus (then, implications? Like, how is this useful for...)
- Inverse and implicit function theorems (and implications)
- Divergence theorem (and relationship to fundamental theorem of calculus)
- Stokes theorem? (Does this apply to higher than  $\mathbb{R}^3$ ?)
- integration by parts (and some implications. Like for ML?)
- Multivariable chain rule
- Change of variables theorem



## Chapter 3

# Least Squares and Linear Regression

### 3.1 Least Squares: Linear Algebra Perspective

Let

$$\mathbf{x}_i \in \mathbb{R}^m, \quad y_i \in \mathbb{R}, \quad i = 1, \dots, m.$$

Here, we have  $m$  examples where  $x_i$  is a feature vector and  $y_i$  is a label. Let

$$X = \begin{pmatrix} -x_1 - \\ \vdots \\ -x_m - \end{pmatrix} \in \mathbb{R}^{m \times n} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^n$$

The matrix  $X$  is commonly known as the design matrix. We wish to solve

$$X\hat{\mathbf{w}} = \mathbf{y}. \tag{3.1}$$

In typical least squares applications, we have  $m > n$ . This means we have an overdetermined system—there are more equations than variables. We'll deal with that case first, then we'll consider connections to the pseudo inverse. Then we'll consider the complementary case where  $m < n$  to complete the picture.

#### 3.1.1 Standard Least Squares: $m > n$

Note that  $X$  is tall and  $\dim \text{col } X \leq n < m$ . The nullspace  $N(X)$  may or may not be empty. if we assume  $\text{rank } X = n$ , then

- $X$  has linearly independent columns
- $\dim \text{col } X = n$
- $N(X) = \{0\}$  (by rank-nullity theorem)

We'll assume  $\text{rank } X = n$ . Suppose for now that  $\hat{y} \in \text{col } X$ . There exists a unique  $\hat{w}$  solving  $X\hat{w} = \hat{y}$  (see Exercise 1.5). Now, suppose  $y$  may not lie in  $\text{col } X$ . We would like to find  $\hat{w}$  such that

$$\hat{w} \in \arg \min_w \|Xw - y\|_2^2. \tag{3.2}$$

(Hence the name “least squares solution.”) We're in a Hilbert space, so this is just the projection of  $y$  onto  $\text{col } X$ . In particular, the projection theorem (recall it?) gives us existence and uniqueness of a solution. Hence, the optimal  $\hat{w}$  satisfies  $(X\hat{w} - y) \perp \text{col } X$ . Equivalently,  $(X\hat{w} - y)^\top X = 0 \iff \hat{w}^\top X^\top X = \mathbf{y}^\top X \iff X^\top X\hat{w} = X^\top \mathbf{y}$ . We know that  $(X^\top X)^{-1}$  exists, hence our least squares solution is

$$\hat{w} = (X^\top X)^{-1} X^\top \mathbf{y}. \tag{3.3}$$

Critically, this is the least squares solution of (3.1) when  $X$  is tall and full rank. We will treat the case when  $X$  is fat in the next section. In Section 3.1.3 we will see that both of these cases can be seamlessly handled by applying the pseudo-inverse.

### 3.1.2 Other case: $m < n$

In Section 3.1.1 we could not solve (3.1) because  $X$  was tall and the system was overdetermined. No solution existed.<sup>1</sup> In this section, we treat the case where  $X$  is fat. The system is underdetermined. We have few equations and many unknowns. And  $N(X)$  is necessarily nonempty.

Editorial comment: I get why you care about the  $X$  tall case. But why do we care about solving the  $X$  fat case in practice?

To make life easier, assume for now that  $\text{rank } X = m$ . Then  $\dim \text{col } X = m$ . Viewed as a linear transformation,  $X$  is surjective and  $y \in \text{col } X$ . Since  $N(X) \neq \{0\}$ , there exists many solutions to (3.1). How can we narrow down the set of solutions and pick one? Some options:

- Pick  $\hat{w}$  so  $\|\hat{w}\|_2$  is small
- Pick sparsest  $w$ —i.e., so  $\|w\|_0$  is small

The first idea jives with the notion of regularization and picking a "low complexity" solution. The second is related to compressed sensing and LASSO.

For now, we'll go with option 1. Suppose  $\hat{w}$  and  $w$  both solve (3.1). Then  $\hat{w} - w \in N(X)$ . If  $\hat{w}$  lies in the row space of  $X$  then  $\hat{w} \perp (\hat{w} - w)$ , since  $\text{row } X \perp N(X)$ . Suppose  $\hat{w} \neq w$ . Then

$$\begin{aligned}\|\hat{w}\|^2 &= \langle w + (\hat{w} - w), w + (\hat{w} - w) \rangle \\ \|w\|^2 - \|\hat{w} - w\|^2 &> \|w\|^2.\end{aligned}$$

Hence, if there exists a solution in the row space, it has minimum norm and is unique.

Recall that the transformation  $X$  may be viewed as a bijection between  $\text{col } X$  and  $\text{row } X$  (e.g., see Exercise 1.5). Since  $y \in \text{col } X$ , and the solution we're looking for is  $\hat{w} \in \text{row } X$ , we may recover  $\hat{w}$  if we can compute the inverse of our bijective map restricted to these sets. This is precisely what the pseudo inverse does, as discussed in Section 3.1.3.

For now, we take a slightly more hands on approach to demonstrate an alternative and maybe more illustrative way to arrive at the solution when  $X$  is full rank. Note that  $\text{rank } X = \text{rank } X^T X = \text{rank } X X^T = \text{rank } X X^T$  (see Exercise ??). Since we assumed  $\text{rank } X = m$  ( $X$  fat), we get  $\text{rank } X X^T = m \iff (X X^T)^{-1}$  exists. We want to solve (3.1). Anything we can set  $\hat{w}$  to to make this work? Let

$$\hat{\mathbf{w}} = X^T (X X^T)^{-1} \mathbf{y}. \quad (3.4)$$

Then  $X \hat{\mathbf{w}} = X X^T (X X^T)^{-1} \mathbf{y} = \mathbf{y}$ . Note that if we let  $\mathbf{z} = (X X^T)^{-1} \mathbf{y}$ , then  $\hat{\mathbf{w}} = X^T \mathbf{z}$ , hence  $\hat{\mathbf{w}} \in \text{row } X$ . Hence,  $\hat{\mathbf{w}}$  is our unique row space solution to (3.1).

In Section 3.1.3 we'll use the pseudo inverse to generalize this to handle the rank-deficient case and tie it together with the tall  $X$  case of Section (3.1.1).

### 3.1.3 Least Squares and the Pseudo Inverse

A linear mapping  $A$  is a bijection from  $\text{row } A$  to  $\text{col } A$  (see Exercise 1.5). The pseudo inverse, studied in Section 1.2, gives the inverse map from  $\text{col } A$  to  $\text{row } A$ . It maps  $\text{col } A^\perp = N(A^T)$  to zero. This is the exact operation we required when deriving the solutions to (3.1) when  $X$  was tall and fat in (3.3) and (3.4) respectively. The pseudo-inverse generalizes these in the sense that

$$\hat{\mathbf{w}} = X^+ \mathbf{y}$$

is identical to (3.3) and (3.4) in the cases previously studied. Moreover, the pseudo inverse applies more generally and gives the solution to (3.1) when  $X$  is tall or fat and when  $X$  is rank deficient.

<sup>1</sup>If  $\text{rank } X \in \mathbb{R}^{m \times n}$  were less than  $n$ , then  $N(X)$  would be nonempty and a solution would exist.

## 3.2 Gauss-Markov Theorem: Least Squares and MLE

**TODO** Connect that MLE is optimal estimator in some sense in presence of gaussian noise

## 3.3 Ridge Regression

**TODO** Set up ridge regression. Explain it in terms of regularization. (Maybe see Tibshirani notes about why shrinking coefficients is helpful.) Also, connect to improving the condition number/stability of the pseudo inverse. Add a numerical example. Also,...

Ridge regression uses the following estimator instead of the pseudo inverse **TODO** Add eqref to psudo inverse section.

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T y. \quad (3.5)$$

This can be helpful for improving numerical stability when  $X$  has many rows that are almost colinear (so  $X^T X$  is ill conditioned). That was supposedly the original motivation. It can also be helpful for regularization in terms of restricting the hypothesis class. I need to do some work to understand better why this is. But it seems to be generally helpful up to a point. Also, it can be interpreted as a formulation of optimizing (3.2) with  $\mathbf{w}$  restricted to a ball of a given radius. (You interpret this as optimizing the Lagrangian of that problem, or something like that.)

## 3.4 LASSO and Elastic Nets

**TODO**: Add sections about these? To paint a more complete picture? Maybe just keep them really short for now so I remember they're relevant and what people claim about them?

## 3.5 Bayesian Linear Regression

Let  $\mathbf{w} \in \mathbb{R}^n$  and consider the standard linear regression model with Gaussian noise

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon$$

where  $\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$ . Here  $\mathbf{x} \in \mathbb{R}^n$  denotes a vector of covariates,  $y$  denotes a scalar output, and  $\mathcal{D} = \{(\mathbf{x}_i, y_i), i = 1, \dots, m\}$  denotes our data set of  $m$  observations. All outputs are assumed to be independent, and we make no assumptions about the distribution of inputs. Let the column vectors  $\mathbf{x}_i$  for all inputs be aggregated in the  $n \times m$  matrix  $X$ .<sup>2</sup> In the non-Bayesian settings discussed previously, our goal was to come up with a point estimate of  $\mathbf{w}$ . In the Bayesian setting, we will assume a prior

$$\mathbf{w} \sim \mathcal{N}(0, \Sigma_p).$$

**Goals:** Having a prior, we may compute a couple of important things that are impossible in the previous settings. We are going to try and compute the posterior distribution on  $\mathbf{w}$

$$p(\mathbf{w} | X, \mathbf{y}) \quad (3.6)$$

and the posterior predictive distribution

$$p(f_* | \mathbf{x}_*, X, \mathbf{y})$$

The first term is the distribution of  $\mathbf{w}$  conditioned on our observed data  $\mathcal{D}$ . This, of course, accounts for the prior as well. The second requires some clarification of the notation. For a given input of covariates  $\mathbf{x}_*$ , let  $f_*$  denote the random variable  $f_* = f(\mathbf{x}_*) + \varepsilon$ . Then, the second item in (3.6) may be thought of as the distribution of the output  $y$  if we were to try a new test point  $\mathbf{x}_*$  outside of our data set. (Note the relation to Gaussian processes.) Handy.

**TODO**: Add an example motivating why you would be interested in this.

<sup>2</sup>In stats literature the matrix  $X^T$  is typically referred to as the design matrix. Rasmussen uses the transpose of this.

We will apply Bayes rule to derive both of these. The first is relatively easy to derive. The second takes more work. We begin with the first.<sup>3</sup> Our goal is to express this quantity in terms of things that we know. We can compute the likelihood  $p(\mathbf{y}|X, \mathbf{w})$  (see below) and we know the prior  $p(\mathbf{w})$ . We will attempt to express it in these terms with an arbitrary normalizing constant. Because the result will be a Gaussian, we won't need to be able to compute the normalizing constant explicitly. To that end, note that

$$\begin{aligned} p(\mathbf{y}|X, \mathbf{w}) &= \prod_{i=1}^m p(y_i|\mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left(-\frac{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}{2\sigma_n^2}\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2} \|\mathbf{y} - X^T \mathbf{w}\|^2\right) \\ &= \mathcal{N}(X^T \mathbf{w}, \sigma_n^2 I). \end{aligned}$$

In the third line, we bring the product into the exp as a sum and note the equivalence to the squared Euclidean norm. Applying Bayes' rule we see that

$$p(\mathbf{w}|\mathbf{y}, X) = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|X)}. \quad (3.7)$$

Maybe there is an easier way to see this, but it confused me at first. The first steps are to use Bayes' rule to see that

$$\begin{aligned} p(\mathbf{w}|\mathbf{y}, X) &= \frac{p(\mathbf{w})p(\mathbf{y}, X|\mathbf{w})}{p(\mathbf{y}, X)} \\ &= \frac{\cancel{p(\mathbf{w})}p(\mathbf{y}|X, \mathbf{w})p(X, \mathbf{w})}{p(\mathbf{y}, X)\cancel{p(\mathbf{w})}} \end{aligned}$$

where in the second line we use  $p(\mathbf{y}, X|\mathbf{w}) = \frac{p(X, \mathbf{y}, \mathbf{w})}{p(\mathbf{w})} = \frac{p(\mathbf{y}|X, \mathbf{w})p(X, \mathbf{w})}{p(\mathbf{w})}$ . Follow your nose from there.

In (3.7) we have expressed the desired quantity in terms of known quantities and a normalizing constant. Writing only the terms from the likelihood and prior and “completing the square” we obtain

$$\begin{aligned} p(\mathbf{w}|X, \mathbf{y}) &\propto \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - X^T \mathbf{w})^\top (\mathbf{y} - X^T \mathbf{w})\right) \exp\left(-\frac{1}{2}\mathbf{w}^\top \Sigma_p^{-1} \mathbf{w}\right) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \bar{\mathbf{w}})^\top \left(\frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1}\right) (\mathbf{w} - \bar{\mathbf{w}})\right), \end{aligned}$$

where

$$\bar{\mathbf{w}} = \sigma_n^{-2} (\sigma_n^{-2} X X^T + \Sigma_p^{-1})^{-1} X \mathbf{y}. \quad (3.8)$$

This is a Gaussian with mean  $\bar{\mathbf{w}}$  and covariance  $A^{-1}$ , where  $A = \sigma_n^{-2} X X^T + \Sigma_p^{-1}$ . Note that the mean and mode coincide so that (3.8) is the MAP estimate of  $\mathbf{w}$ . The MAP estimate coincides with ridge regression (3.5), which, in the non-Bayesian setting is known as the penalized MLE. The name comes from the fact that the least squares estimator is an MLE in common noise settings (e.g., iid Gaussian—see Section 3.2) and ridge regression (3.5) is the least squares estimator (3.3) with a penalty term.<sup>4</sup>

The predictive distribution is obtained by averaging over all parameter values, weighted by their posterior distribution.

$$\begin{aligned} p(f_*|\mathbf{x}_*, X, \mathbf{y}) &= \int p(f_*|x_*, \mathbf{w})p(\mathbf{w}|X, \mathbf{y}) d\mathbf{w} \\ &= \mathcal{N}\left(\frac{1}{\sigma_n^2} \mathbf{x}_*^\top A^{-1} X \mathbf{y}, \mathbf{x}_*^\top A^{-1} \mathbf{x}_*\right) \end{aligned}$$

<sup>3</sup>This is basically copied from Rasmussen, p.9-10

<sup>4</sup>In that section, we defined  $X$  by filling rows with  $\mathbf{x}_i$ 's. Hence, it's the transpose of what we've used here.

Annoyingly, this isn't derived in Rasmussen. But the derivation follows directly from Bishop Section 2.2.3, and in particular (2.115) and (2.116). (TODO: Add as exercise? Go over this. What's the gist of how it's done?)

**TODO:** Add reference to notebook exploring this. Maybe drop in a figure and some kind of link to the notebook. Where to store the notebooks so that I can easily link to them?

## 3.6 LASSO

**TODO** This is a digression and I'm leaning away from doing it for now. Basically, in spite of all the data science tools out there, I feel like regression is what people use and rely on. Hence why this chapter is so valuable. And LASSO tells a valuable part of that story. Maybe put in a short section based on some Tibshirani notes? Mostly, just so I have a starting point/context if/when I do run into this. Then plan to expand it then? If I do it, I could make that note at the beginning of the section.





## Chapter 4

# PCA

**TODO** Add various perspectives on PCA.



## Chapter 5

# Other fun topics/one off topics

- Thompson sampling (just because)
- Gibbs sampling and other sampling techniques (let's be honest, this is it's own chapter)



## Chapter 6

# Other notes/Daily notes

### 6.1 Gaussian Distribution: Normalizing Constant 5/24/24

The standard Gaussian distribution is given by  $1/\sqrt{2\pi} \int e^{-\frac{1}{2}x^2/2}$ . How do we know this is a valid density function? Specifically, how do we know the normalizing constant is 1? Here's a quick proof.

Before starting, recall the change of variables formula

$$\int_{G^{-1}(\Omega)} f(x) dx = \int_{\Omega} f(G(x)) |\det G(x)| dx.$$

Here,  $G$  is an invertible diffeomorphism. Now, what we're essentially looking for is to show that

$$\int_{-\infty}^{\infty} \exp(-\frac{1}{2}x^2) dx = \sqrt{2\pi}.$$

The following trick allows us to compute it. Let  $I = \int_{-\infty}^{\infty} \exp(-\frac{1}{2}x^2) dx$ . We will compute  $I^2$  instead of  $I$ . I'm not sure if there is some good intuition for why this is easier. But it works. We have

$$\begin{aligned} I^2 &= \left( \int_{-\infty}^{\infty} \exp(-\frac{1}{2}x^2) dx \right) \left( \int_{-\infty}^{\infty} \exp(-\frac{1}{2}y^2) dy \right) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-\frac{1}{2}(x^2 + y^2)) dx dy. \end{aligned}$$

This is amenable to a conversion to polar coordinates. Now, I'm sure we could apply the change of variables formula directly at this point. But with a little bit of cheating, I'm going to note that what I'd ultimately like to do is integrate over  $r \geq 0$ . I'm going to observe that if we consider polar coordinates, then integrating a ring of fixed radius from  $\theta = 0$  to  $2\pi$ , we have a density of  $\int_0^{2\pi} e^{-r^2} d\theta = 2\pi r e^{-r^2}$ . So,

$$I^2 = \int_0^{\infty} r e^{-\frac{1}{2}r^2} dr.$$

Now, apply change of variables. With  $G(x) = \sqrt{2}x$  we see that

$$\begin{aligned} I^2 &= 2\pi \int_0^{\infty} \sqrt{x} e^{-x} x^{-1/2} dx \\ &= 2\pi \int_0^{\infty} e^{-x} dx. \end{aligned}$$

The last integral we know how to evaluate, and it evaluates to 1, so we see that  $I^2 = 2\pi$  or  $I = \sqrt{2\pi}$ .

### 6.2 Polar Change of Coordinates

When you do a polar change of coordinates, you change from  $dx dy$  to  $r dr d\theta$ . I get confused where the  $r$  actually comes from. And, honestly, the way change of variables is applied is a little counterintuitive

to me. So, I'm jotting it down. For completeness, recall the change of variables formula.

$$\int_{G^{-1}(\Omega)} f(x) dx = \int_{\Omega} f(G(x)) |\det DG(x)| dx.$$

As a running example, we'll consider how to integrate  $\int e^{-(x^2+y^2)}$ . To apply change of variables, let  $f(G(r, \theta)) = e^{-r^2}$ . The conversion from polar to cartesian is given by

$$G(r, \theta) = (r \cos \theta, r \sin \theta).$$

Let  $H := G^{-1}$  so that  $f(x, y) = f(G(H(x, y)))$ . Here, we have  $H(x, y) = (x^2 + y^2, \arctan 2(y/x))$ . So, we see that

$$f(x, y) = e^{x^2+y^2}.$$

In order to apply change of variables, note that

$$D_{(r, \theta)} G(r, \theta) = \begin{pmatrix} \cos \theta & r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix},$$

so that  $|\det DG| = r(\sin^2 \theta \cos^2 \theta) = r$ . Now, applying change of variables, we see that

$$\int_{\mathbb{R}^2} e^{-(x^2+y^2)} dL(\mathbb{R}^2) = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} e^{-r^2} r dL(\mathbb{R}^2),$$

where I'm trying to use that notation to denote the standard Lebesgue measure over  $\mathbb{R}^2$ . Note that if we let  $\Omega$  be the domain on the right hand side, then  $G^{-1}(\Omega)$  gives us the domain on the left hand side. I'm not sure if I'm actually doing that right, but this should be equivalent to

$$\int_{\mathbb{R}^2} e^{-(x^2+y^2)} dx dy = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} e^{-r^2} r dr d\theta.$$

And, of course, in general, you have

$$\int_{\mathbb{R}^2} f(x, y) dx dy = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} f(G(r, \theta)) r dr d\theta,$$

where we've carefully defined what the functions on each side of that mean.

## 6.3 Change of Variables

This section is about, in my opinion, one of the most useful and beautiful results in calculus.

**Theorem 2** (Change of Variables). *Suppose that  $\Omega$  is an open subset of  $\mathbb{R}^n$  and  $G : \Omega \rightarrow \mathbb{R}^n$  is a  $C^1$  diffeomorphism. If  $f$  is  $L^1$  integrable over  $G(\Omega)$ , then*

$$\int_{G(\Omega)} f(x) dx = \int_{\Omega} f(G(x)) |\det DG(x)| dx.$$

Above,  $DG(x)$  is the Jacobian determinant of  $G$ .

Sometimes, when integrating, it is helpful to do a change of variables in order to massage an integral into a form that is easier to handle. This change of variables is represented by  $G(x)$  inside the integral. However, simply modifying the argument to the integral changes the value of the integral. The change of variables formula tells us the necessary correction term inside the integral. This is  $|\det DG(x)|$ . We'll build intuition for what this correction term means later on. But first, we'll go over a few examples to illustrate how eminently useful this result is.

**Example:** Change of variables is the multivariable generalization of  $u$  substitution, familiar from introductory calculus. Anytime you use  $u$ -substitution to reformulate some tricky integral, you're using the change of variables formula.

**Example (Function of a RV):** Suppose that  $X \in \mathbb{R}^n$  is a random variable with density function  $f_X$ . Let  $Y \in \mathbb{R}^n$ , with  $Y = g(X)$  where  $g$  is a  $C^1$  diffeomorphism. What is the density function of  $Y$ ? This is an important question that arises often in ML applications and elsewhere. Ignoring questions of measurability, note that for a set  $A \subset \mathbb{R}^n$

$$\mathbb{P}(Y \in A) = \mathbb{P}(g(X) \in A) = \mathbb{P}(X \in g^{-1}(A)) = \int_{g^{-1}(A)} f_X(x) dx.$$

Let  $h = g^{-1}$  to make notation clearer. Applying change of variables we see that

$$\int_{h(A)} f_X(x) dx = \int_A f_X(h(y)) |\det Dh(y)| dy.$$

where in the second line we somewhat arbitrarily change the variable of integration from  $x$  to  $y$  to make the connection to the density  $f_Y$  clearer. Hence,

$$\mathbb{P}(Y \in A) = \int_A f_X(h(y)) |\det Dh(y)| dy$$

Since this holds for all  $A$  (ignoring measurability considerations), we see that the integrand on the RHS must be the density of  $Y$ , hence

$$f_Y(y) = f_X(h(y)) |\det Dh(y)|. \quad (6.1)$$

Any time you need an analytic expression for a function of a random variable, the change of variables formula is indispensable.

For a concrete example of this, suppose that  $X$  is uniformly distributed on  $[0, 1]$  and  $Y = X^2$ . What is the density of  $Y$ ? It's  $f_Y(y) = \frac{1}{2\sqrt{y}}$ . This is shown in Figure ?? . Also, see demo in Jupyter notebook.

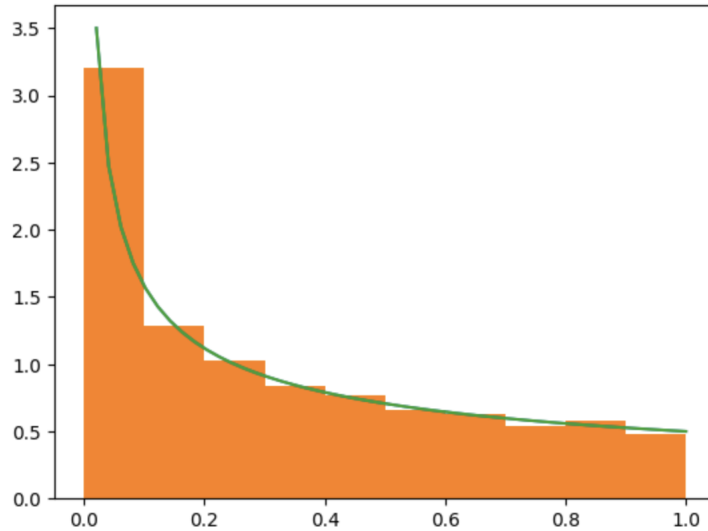


Figure 6.1: Density of  $Y = X^2$ ,  $X \sim U([0, 1])$

**Example (Normalizing Flow):** Suppose we have a set of points  $(z_i)_i$  where each  $z_i \sim p$ , i.i.d, where  $p$  is some density function. Suppose that we would like to model the density  $p$ . Consider a change of variables-based approach. Let  $f$  be some distribution we know how to sample from easily (e.g., normal or uniform density). Suppose we sample  $W \sim f$  and then pass it through some parametric function  $g_\theta(x)$ , which we assume to be a  $C^1$  diffeomorphism. We wish the random variable  $Z = g_\theta(W)$  to have a density as similar as  $p$  as possible. Letting  $h_\theta = g_\theta^{-1}$ , change of variables tells us that the density function of  $Z = g_\theta(W)$  is

$$f(z; \theta) = f(h_\theta(z)) |\det D_z h_\theta(z)| \quad (6.2)$$

We may fit  $f(z; \theta)$  to our dataset  $(z_i)$  using a maximum likelihood approach. If we assume  $(z_i)$  is iid, then the likelihood of the joint sample is

$$L(\theta) = \prod_i f(z_i; \theta) = \prod_i f(h_\theta(z_i)) |\det D_z h_\theta(z_i)|,$$

and the negative log likelihood is

$$\ell(\theta) = - \sum_i [\log f(h_\theta(z_i)) + \log |\det D_z h_\theta(z_i)|]. \quad (6.3)$$

Minimizing (6.3) simply maximizes the likelihood of the data under our model. This provides us with a function  $h_\theta$  that allows us to draw samples from an approximation to  $p$ . Once  $\ell$  is minimized, if we let  $Z = h_\theta^{-1}(X)$  for  $X \sim f$ , then the density of  $Z$  is approximately  $p$ . Methods exist for obtaining  $h_\theta^{-1}$ . Also, if we simply wish to evaluate the approximate density of  $p$  at test point, then we simply use (6.2). This method of approximating a density is called a normalizing flow, and is a current research area in ML. There is *a lot* we could say about this, but we'll ignore it to stay focused on CoV for now. Maybe save this for another time.

**Example (Law of the unconscious statistician):** Let  $X$  be a random variable with density function  $f_X$ . Let  $g$  be some measurable function. What is the expected value of  $g(X)$ ? It is commonly taken to be, by definition,

$$\mathbb{E}[g(X)] = \int g(X) f_X(x) dx. \quad (6.4)$$

*However*, this is not obvious. This is so commonly taken to be the definition of  $\mathbb{E}[g(X)]$  that it is known as the law of the unconscious statistician.

The definition of the expectation of a random variable is

$$\mathbb{E}[Z] := \int x f_Z(z) dz.$$

Consider the simple case where  $X$  is scalar-valued and  $g$  is a  $C^1$  diffeomorphism. Applying (6.1), we know that the density function for  $Y = g(X)$  is

$$f_Y(y) = f_X(g^{-1}(y)) |(g^{-1}(y))'|$$

and

$$\mathbb{E}[g(X)] = \int y f_X(g^{-1}(y)) |(g^{-1}(y))'| dy.$$

Now, apply change of variables again with  $y = g(x)$ . To do this, we just substitute  $g(x)$  in for  $y$  above and multiply by  $g'(x)$  inside the integral. We arrive at

$$\begin{aligned} \mathbb{E}[g(X)] &= \int g(x) f_X(x) \underbrace{\frac{d}{dy} g^{-1}(g(x)) g'(x)}_{=1} dx \\ &= \int g(x) f_X(x) dx, \end{aligned}$$

where in the last line, we recognize from the chain rule that  $\frac{d}{dy} g^{-1}(g(x)) g'(x) = \frac{d}{dx} g^{-1}(g(x)) = 1$ . Hence, (6.4) is correct. But the result isn't obvious. Here, we only showed it in a simple scalar-valued case with invertible and differentiable  $g$ . It does hold more broadly. A measure-theoretic treatment shows that it holds when  $g$  is a measurable function and the random variable  $X$  has a finite mean (I think... double check).

**Example (Reparameterization Trick):** The CoV formula is critical in using the reparameterization trick, which is the foundation of an important class of variational-inference based methods in ML. Save the details for another time, for the sake of brevity.

**Example (Polar change of coordinates/Gaussian normalization constant):** Consider evaluating the integral

$$I = \int e^{-\frac{1}{2}(x^2+y^2)} dx dy. \quad (6.5)$$



TODO: Just added in the  $1/2$  above. Make consistent through the rest! It should already be consistent after the easter egg. (We'll see at the end why this is actually an extremely useful integral to know.) This integral is radially symmetric and will be easier to evaluate in polar coordinates, where  $x = r \cos \theta$  and  $y = r \sin \theta$ . The conversion from polar to cartesian is given by the function

$$G(r, \theta) = (r \cos \theta, r \sin \theta).$$

Letting  $f(x, y)$  be the integrand above, we have

$$f(G(r, \theta)) = \exp(-(r^2 \cos^2 \theta + r^2 \sin^2 \theta)) = \exp(-r^2).$$

In order to apply change of variables, note that

$$D_{(r, \theta)} G(r, \theta) = \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix},$$

so that  $|\det DG| = r(\sin^2 \theta \cos^2 \theta) = r$ . Now, applying change of variables, we see that

$$\int_{\mathbb{R}^2} e^{-(x^2+y^2)} dL(\mathbb{R}^2) = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} e^{-r^2} r dL(\mathbb{R}^2),$$

where I'm trying to use that notation to denote the standard Lebesgue measure over  $\mathbb{R}^2$ . Note that if we let  $\Omega$  be the domain on the right hand side, then  $G^{-1}(\Omega)$  gives us the domain on the left hand side. I'm not sure if I'm actually doing that right, but this should be equivalent to

$$\int_{\mathbb{R}^2} e^{-(x^2+y^2)} dx dy = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} e^{-r^2} r dr d\theta.$$

And, of course, in general, you have

$$\int_{\mathbb{R}^2} f(x, y) dx dy = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} f(G(r, \theta)) r dr d\theta,$$

where we've carefully defined what the functions on each side of that mean.

As a final easter egg in this example, note that

$$I = \int_{\substack{r \geq 0 \\ \theta \in [0, 2\pi)}} e^{-\frac{1}{2}r^2} r dr d\theta = 2\pi \int \exp(-r^2) r dr.$$

Applying another change of variables where we let  $r = G(u) = u^{1/2}$  we get (I might have missed a minus sign below...)

$$I = 2\pi \int_0^\infty \exp(-\frac{1}{2}u) du = 2\pi.$$

So, we have an explicit evaluation of this integral. But, note that

$$\left( \int \exp(-\frac{1}{2}x^2) dx \right)^2 = I$$

You can see this by just expanding the square explicitly, and then combining everything under a double integral to get (6.5). So,

$$\int \exp(-\frac{1}{2}x^2) dx = \sqrt{2\pi}.$$

or, equivalently,

$$\frac{1}{\sqrt{2\pi}} \int \exp(-\frac{1}{2}x^2) dx = 1$$

This is how we know what the normalizing constant is for a normal random variable. End example.

**Example** A high-dimensional gaussian is concentrated on the surface of a sphere. Change of variables is useful in deriving this. (Will do later.)

**Bonus Example 1:** How can you sample from the unit 2sphere (embedded in  $\mathbb{R}^3$ )? Consider using polar coordinates, so  $\theta \in [0, 2\pi]$  is an azimuth and  $\phi \in [0, \pi]$  is an elevation (This is an angle measured from the north pole.) Can you just sample  $\theta, \phi$  uniformly from their sets? (Pause and consider.) Answer: No. To see why, consider a change of variables from uniform to polar coordinates. Suppose that  $A$  is some subset of the unit sphere. Then under a uniform distribution on the unit sphere, we have

$$\mathbb{P}(A) = \frac{1}{4\pi} \int_A 1 \, dx \, dy \, dz = \frac{1}{4\pi} \int_{g^{-1}(A)} \sin(\phi) \, d\phi \, d\theta,$$

where we have used the fact that the surface area of the unit sphere is  $4\pi$  and the fact that the spherical change of coordinates (see Figure 6.2) so that

$$\begin{aligned} x &= r \sin(\phi) \cos(\theta) \\ y &= r \sin(\phi) \sin(\theta) \\ z &= r \cos(\phi). \end{aligned}$$

The  $\sin(\phi)$  comes from computing  $\det DJ(\theta, \phi, r) = r^2 \sin(\theta)$ .<sup>1</sup> The integrand on the right hand side

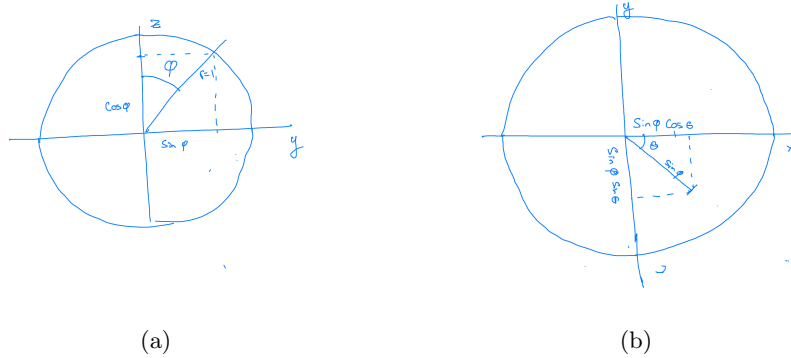


Figure 6.2: Spherical coordinates. I accidentally left out the  $r$  in this figure. So just multiply everything by a radius  $r$ .

gives us the density for  $\phi, \theta$  that corresponds to a uniform distribution on the sphere. In particular, the density to sample uniformly from the unit sphere in spherical coordinates is

$$f(\theta, \phi) = \frac{1}{4\pi} \sin(\phi) \mathbf{1}_{[0, 2\pi] \times [0, \pi]}.$$

How do you sample from this distribution? Here's a trick that works, that I haven't thought through how you'd come up with from scratch. Let  $u, v \sim U([0, 1]^2)$ . Let  $\theta = 2\pi u$  and  $v = \cos^{-1}(2v - 1)$ . Letting

$$h(\theta, \phi) = g^{-1}(\theta, \phi) = \begin{pmatrix} 1/2\pi\theta \\ (\cos(\phi) - 1)/2 \end{pmatrix}$$

and applying (6.1) we immediately get that  $f_{\theta, \phi}(\theta, \phi) = \frac{1}{4\pi} \sin(\phi)$ . The domain of the density comes from looking at the range  $g(u, v)$ .

It should be noted that this is the hard way of sampling from the unit sphere. The easy way is to sample from a gaussian and then normalizing the samples. But this examples builds intuition that I think can generally be helpful.

**Bonus Example 2:** Suppose you want to uniformly sample from the set of all rotation matrices in three dimensions. How can you do that? A rotation consists of an axis of rotation and an angle of rotation. You can uniformly pick the axis of rotation as above, and then sample an angle of rotation

<sup>1</sup>This actually raises an important subtlety. If we don't fix  $r = 1$ , then we are going from  $\mathbb{R}^3 \rightarrow \mathbb{R}^3$ . The Jacobian is invertible. If we do fix  $r = 1$ , then we have 3 cartesian coordinates and 2 polar coordinates. And change of variables can't be applied. We get around this by computing the change of variables for 3 to 3 coordinates, and then fixing  $r = 1$ , and then integrating only with respect to the two dimensional Lebesgue measure w.r.t. the angles. But this feels like it's playing a bit fast and loose.

uniformly from  $[0, 2\pi]$ . Question: How do you sample a random rotation matrix in higher dimensions? Not sure...

### Building Intuition

The previous examples have demonstrated the clear utility of being able to change variables inside of an integral. The key term in CoV is the correction term  $|\det DG(x)|$ . This gives us the exchange rate required to compensate for the introduction of  $g(x)$  inside the integrand. To build intuition about this, consider the example  $f(x) = \frac{1}{2}x^2$ . Suppose, for some made up reason, that we wish to evaluate the integral  $\int_{[0,1]} f(x) dx$ , but we only have access to  $f(g(x))$ , where  $g(x) = cx$ ,  $c > 1$ .<sup>2</sup> The

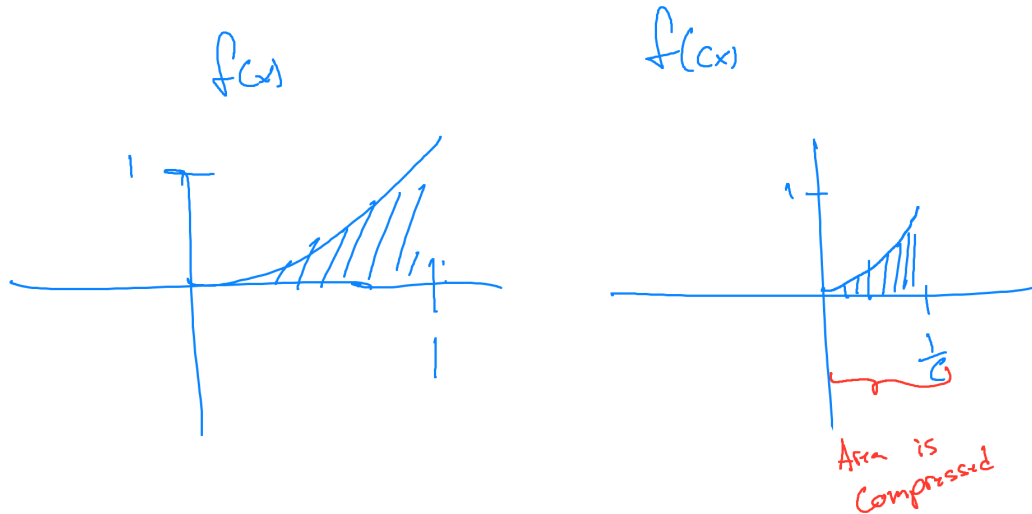


Figure 6.3: Illustration of "compression" effect.

first thing to notice is that if we're going to evaluate the same integral, we've got to make sure the underlying domains match up. This is illustrated in Figure ???. The function  $g(x) = cx$  has the effect of compressing the domain of integration. The end point for integration inside the function is reached when  $x = 1/c$ . We make sure the domains match up by taking the preimage of the desired domain under  $g$ . Equivalently, we may state this as they do in Theorem 2 in terms of  $g$  rather than the preimage, via

$$\int_{g(\Omega)} f(x) dx = \int_{\Omega} \text{stuff}.$$

Next, because  $g$  compresses the domain of integration, if we directly integrate, we will end up with a value that is too small. Again, see Figure ??. We compensate for this compression term by understanding the "compression rate" and compensating for it when computing the integral. I'm going to

<sup>2</sup>We keep  $c > 0$  to make the discussion of "compression" consistent later. But the same reasoning holds for any  $c \neq 0$ , just change the terminology to "dilation" where necessary.

call this the exchange rate. The “exchange rate” at which it compresses can be understood by looking at the preimage of the canonical volume-1 set  $[0, 1]$ . This tells us the volume of the set that maps to  $[0, 1]$ . The exchange rate is given by the ratio.

$$\frac{\text{vol}[0, 1]}{\text{vol } g^{-1}([0, 1])}.$$

Evaluating this expression we get a “compression ratio” (I’m making that term up) of  $c$ . Higher compression ratio means more compressive. Because the function  $g$  is linear, we can also get the exchange rate by just looking at the ratio

$$\frac{\text{vol } g([0, 1])}{\text{vol}[0, 1]}.$$

This tells us the relative size of the set that  $[0, 1]$  is mapped to, and gives us the same information about the compression rate in the previous expression. Evaluating this expression gives us the same compression ratio. This ratio tells us the effect that including  $g$  inside the argument of  $f$  has. In particular, it tells us the rate at which infinitesimal segments are compressed inside the integral. We may now compensate for the infinitesimal compression effect inside the integral by multiplying by the exchange rate

$$\int_{[0,1]} \frac{1}{2}x^2 dx = \int_{[0,1/c]} \frac{1}{2}(cx)^2 c dx = 1$$

This same idea holds more generally if we consider the effect of a bijective linear map  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . The following result is critical.

**Theorem 3.** *Let  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Then the image of the hypercube  $T([0, 1]^n)$  is a parallelipiped. Representing  $T$  as a matrix, the (signed) volume of the parallelipiped is given by  $\det T$ .*

In the previous theorem, the sign on the volume tells us about the orientation of the parallelipiped relative to the original hypercube. We won’t go into that here. Taking  $|\det T|$  gives the unsigned volume.

Example: Consider the matrix

$$T = \begin{pmatrix} 0 & 1 \\ 2 & 3 \end{pmatrix}.$$

The image of the cube  $[0, 1]^2$  is shown in Figure ?? . See also Jupyter notebook.

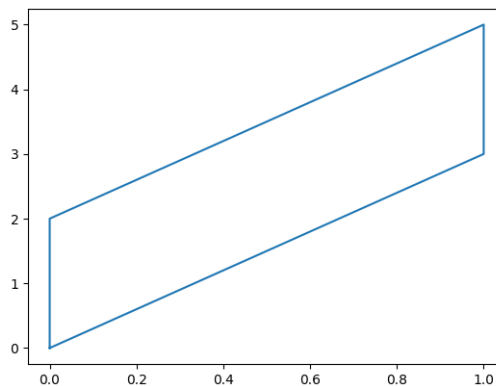


Figure 6.4: The image of the cube  $[0, 1]^2$  under  $T$ .

Note that the vertices are given by the image of the vertices of the hypercube. In particular, two of the vertices are precisely the columns of  $T$ .

Returning to change of variables with a linear function  $T$  inside of  $f$ , we see that  $|\det T|$  gives our exchange rate inside the integral:

$$\int_A f(x) dx = \int_{T^{-1}(A)} f(T(x)) |\det T| dx.$$

More generally, if the function  $g$  is nonlinear, then we must use an instantaneous exchange rate. This is obtained by linearizing  $g$  at each point  $x$ , and computing the volume exchange rate for the linearization; i.e.,  $|\det Dg(x)|$ . This result takes a measure theoretic treatment to prove generally. But this is the basic intuition.

Also, a concrete example with a density function would be nice. So you can see how the map is sort of redistributing the mass in the density function, and how the determinant term helps

## 6.4 High dimensional spaces

We often work in high dimensional spaces. Many datasets have 10 or more features. MNIST data has 728 dimensions. A high resolution image has a million dimensions. If you're working with a kernel machine, even if your data has relatively low dimensionality, it's embedded in a high dimensional space. We develop intuition in 2 and 3 dimensions. But our basic mental picture fails in serious ways to generalize in higher dimensions, where properties are counterintuitive, and sometimes quite bizarre at first sight. The purpose of this discussion is to help build intuition about high dimensional spaces so you can be a better data scientist, ML practitioner, or engineer. Also, understanding the basic geometric properties of high dimensional spaces lays the foundation for better understanding statistical properties in high dimensional spaces.

### 6.4.1 Volume of $d$ -ball

The ball of radius  $r$  in  $\mathbb{R}^d$  is given by

$$B_r := \{x \in \mathbb{R}^d : \|x\| \leq r\}.$$

The volume of the ball in  $d$  dimensions is obtained by integration. There are a few different ways of doing this, some pretty clever.<sup>3</sup> But the brute force way is to just integrate in  $d$  dimensions analogous to how you've done it in 3. Set up the integral you want to solve in cartesian coordinates.

$$\int_{B_r} dx_1 \cdots dx_d.$$

Do a change of variables to spherical coordinates. (Hyper)-spherical coordinates in  $d$  dimensions are defined analogous to the 3d case, and you can look up the conversion and wrap your head around how it works. At the end of the day, you solve the integral

$$\int r^{d-1} \sin^{d-2}(\phi_1) \cdots \sin(\phi_{d-1}) d\phi_{d-1} \cdots d\phi_1 dr,$$

where I've omitted the domain because I'm lazy. But it's a product of intervals. So, a high dimensional rectangle. You chug through this integral and you get that the volume of the  $d$  ball of radius  $r$  is given by

$$V_d(r) = \frac{\pi^{d/2} r^d}{\Gamma(\frac{d}{2} + 1)}. \quad (6.6)$$

Just think of the Gamma function as an extension of the factorial to non integer values. Figure 6.5, taken from wikipedia, plots the volume of a few different balls as the dimension increases. Note that the volume of a ball plummets to zero as  $d \rightarrow \infty$ . For bigger radii, the volume will be a lot bigger in lower dimensions, but no matter what, the volume of any ball goes to zero as  $d \rightarrow \infty$ . That's freaking weird.

---

<sup>3</sup>I was tempted to put in the version that uses the Gaussian integral, since we solved that in the CoV section. But it makes an assumption about a proportionality relationship with the surface area of high dimensional spheres that isn't obvious. And I didn't want to have to prove. And integrating this stuff out isn't really the point of these notes. So, I took the lazy route of just giving a high level idea of the brute force approach, but skipping the details because they're tedious and don't advance the objective of these notes.

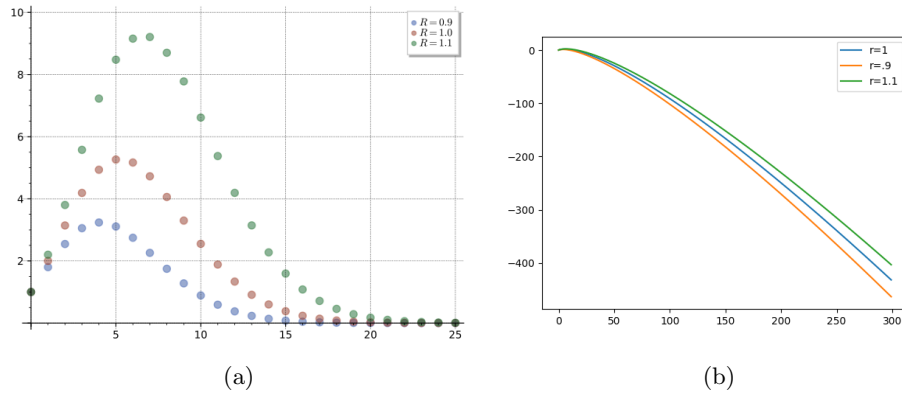


Figure 6.5: Volume of unit ball as a function of the dimension. (a) is the standard plot, taken from Wikipedia (b) plots the log of the volume vs dimension.

### 6.4.2 Balls and Cubes

Consider a ball of radius  $1/2$  inscribed inside a cube of width 1. e.g., 2d example. The ball always touches the surface of the cube. As  $d$  increases, the volume of the cube stays one. But the volume of the sphere inside goes to zero (quickly!). As another example, consider a ball of radius 1 and cube of width 1. Draw in 2d. The cube is contained well within the ball. But when we go to 3 dimensions, the gap closes a bit. The corner of the cube is at  $\sqrt{d/2^2} = 1/2\sqrt{d}$ . At  $d = 4$  the corners of the unit cube touch the surface of the ball of radius 1. For  $d > 5$ , the corners poke through the surface. Moreover, there are  $2^d$  vertices poking through.

### 6.4.3 Most of the mass is near the equator

See section 1.2.3 of the CMU book chapter on geometry of high dimensional spaces.

Consider a small slice through the ball near some equator. That is, consider the portion of the ball that lies "above" the plane where  $x_1 = \epsilon$ . One can show that very little of the volume of the sphere lies in this region. This is accomplished by simply integrating (and using a few approximations). Add this if I have time. But the end result is the following lemma.

**Lemma 1.** *For any  $c > 0$ , the volume of the hemisphere above the plane  $x_1 = \frac{c}{\sqrt{d-1}}$  is less than  $\frac{2}{c}e^{-c^2/2}$ .*

Intuitively, I think this makes a lot of sense. You're fixing a plane by setting  $x_1 = 0$  and considering a  $\epsilon$  padding of that plane, and considering how much mass is in that padded region. While it's true that  $\epsilon$  is small, by considering a plane, you're considering an  $d - 1$  dimensional set. When you move in the orthogonal direction, you only have 1 dimension to move in. In some sense, this dimension can't contain a lot of volume. Most of the volume is contained near the plane.

You can get some intuition for this by going from 2 dimensions to 3 dimensions. If you fix a diametric plane in 2 dimensions, then moving orthogonal to the plane you capture a lot of volume. Now do the same thing in three dimensions. More of the volume is contained near the equatorial plane. As the dimension increases, more and more of the volume is contained near  $d - 1$  dimensional equatorial hyperplane.

While the idea that most of the mass is near the equator, its implication in practice is almost stupidly obvious. Suppose you draw a uniform random sample from the unit ball. Then with high probability, all of the coordinates are going to be small. This makes sense when you think of the fact that the norm of the sample must be contained near within a radius of 1. If any one of the components is close to 1, then the others have to be really, really tiny. Another way to think about it, suppose that  $x_i = \text{big}$ , for some coordinate  $i$ . This means that the other coordinates must all be contained in a very small  $d - 1$  dimensional ball. But the product of a tiny  $d - 1$  dimensional ball with, even the entire diameter of the circle, is going to have really small volume. So, that's not going to happen.

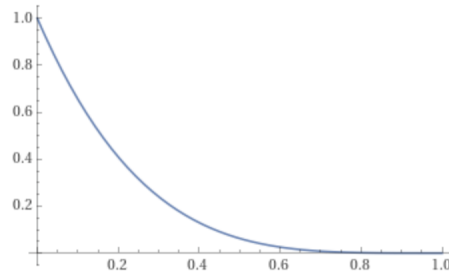


Figure 6.6: Volume near surface of sphere.

#### 6.4.4 Most of the volume is near the surface

(Copied from cmu book.) By (6.6), the ratio of the volume of a sphere of radius  $1 - \varepsilon$  to the volume of a unit sphere in  $d$  dimensions is  $(1 - \varepsilon)^d$ . This is shown in Figure ?? . The higher the dimension, the more the volume is concentrated on the surface of the sphere. This basically follows from the fact that the higher  $d$ , the more  $x^d$  looks flat below 1 and like a wall straight up at 1.

#### 6.4.5 High dimensional cubes

The fact that most of the volume of a high dimensional ball is near the surface is a little more transparent when you look at hypercubes. (Which, I'm going to call a ball, even though by ball here, we really do only mean L2 ball, which is an important distinction. The volume near the surface property is probably true for any convex body?) Consider a unit cube, and the volume outside the inner cube with side length  $1 - \varepsilon$ . Draw this out for 1, 2, and 3d. The volume of the inner region diminishes as  $(1 - \varepsilon)^d$ . (See other notes on this for a compelling presentation.)

#### 6.4.6 Distances between randomly sampled points

Suppose you sample two points uniformly at random from the unit ball. How far apart will they be? Show demo. Answer is that even though the volume of the ball is going to zero, the distance between randomly sampled points converges to a constant. State a general result about distance between furthest and closes points, for iid case. (Technically, this won't cover the ball case. But it tells us something.)

Question: What does this mean about things like knn? Under reasonably general conditions, you can show that as the dimension grows large, the distinction between nearest and furthest point vanishes. An example of this is the following result.

**Theorem 4.** Suppose that  $X_n \in \mathbb{R}^d$  is composed of  $d$  iid random variables and that for some  $k \geq 1$ ,  $\lim_{d \rightarrow \infty} \frac{\|X_d\|_k}{\mathbb{E}\|X_d\|_k} = 0$ . Then

$$\frac{D_{max,d}^k - D_{min,d}^k}{D_{min,d}^k} \rightarrow_p 0 \quad \text{as } d \rightarrow \infty,$$

where  $\rightarrow_p$  indicates convergence in probability to the point mass on zero.

Example: Sampling points from the  $d$ -cube. (TODO: Add computational example. TODO: Is there an interesting two-class classification version of this? Like, can you distinguish between points on opposite sides of a hypercube? Probably best to just look at that paper. TODO: If I don't have time to do this now, then that's OK. Maybe this is good fodder for discussion. When does this tell us about when knn breaks down? How does this impact things like neural networks operating in high dimensional spaces?)

Moral: You need to be careful about the meaning of distance metrics and similarity in high dimensions.

### 6.4.7 Digression: Lp balls

Question: We've implicitly looked at l2 and l-infty balls. You can consider similar properties for other lp balls. Note paper about surprising properties of distances. Add reference. Moral of that paper: Can consider fractional distances, and these tend to be a little more meaningful.

### 6.4.8 Discussion

Mention things like, is data really usually high dimensional. No. usually does something like live in a lower dimensional manifold. Touch on importance of dimensionality reduction for being able to say meaningful things about data in something like knn. Touch on PCA and isomap.

Mention that sometimes you intentionally move up to higher dimensions. Usually for kernel methods. This allows you to do things like linearly separate classes of data. Mention Cover's theorem.

## 6.5 High Dimensional Gaussian

Random variables in high dimensional spaces can behave in counterintuitive ways. We typically develop intuition in low dimensions, but don't realize which aspects of our intuition fail to extend to high dimensional spaces. A good place to start building better intuition is looking at the properties of high-dimensional Gaussian random variables. The main property of high dimensional Gaussians that you'll hear cited is that most of the mass is concentrated around the surface of a sphere. The reason for this bizarre effect has to do with more basic properties of high dimensional geometry than anything else. A high-dimensional ball has most of its mass concentrated on the surface of a sphere. We'll begin by briefly considering properties of high-dimensional spheres and hypercubes, and then characterize the high dimensional Gaussian.

### 6.5.1 Radial density of Gaussian

Suppose  $x \in \mathbb{R}^d$  is sampled from an isotropic Gaussian. We would like to understand the density of  $r = \sqrt{\sum_{i=1}^d x_i^2}$ . To this end, first note that the random variable  $r$  can be viewed as

$$r = \sqrt{w}$$

where  $w$  is a  $\chi^2$  distribution with  $d$  degrees of freedom having density

$$f_w(w) = cw^{\frac{d}{2}-1}e^{-\frac{w}{2}},$$

and  $c$  is a proportionality constant that I'm ignoring. Applying the change of variables expression (6.1), we can compute the density of  $r$ . We have  $h(r) = g^{-1}(r) = r^2$  and  $Dh = 2r$ . Subbing this in we get

$$f_r(r) = cr^{d-2}e^{-\frac{r^2}{2}}r = cr^{d-1}e^{-\frac{r^2}{2}},$$

where we absorb stuff into the constant. We would like to characterize where most of the mass of  $f_r(r)$  lies. A plot of  $f_r(r)$  for various values of  $d$  is shown in Figure ???. In general, the plots suggest that  $f_r(r)$  is log concave, and most of the mass is concentrated in a band that increases roughly as  $\sqrt{d-1}$ . Our plan of attack will be

1. Take the log of the density, so we hopefully deal with a concave function
2. Consider an interval  $I$  of width  $2d$  about the maximum of the function. We'll try and characterize the amount of mass in this interval.
3. Construct a second order approximation of the function. In particular, construct an upper bound on the function and then use this to bound the amount of mass outside of the interval.
4. Because we've been cavalier with the proportionality constant, we don't know what the total mass of  $f$  is. (It's not a probability density.) We'll come up with a lower bound on the total mass of  $f_r$  and use that to estimate the fraction of mass in the interval  $I$ .



5. Continuing with the previous point, let  $M_{\bar{I}}$  be the mass outside the interval  $I$  and let  $TM$  denote the total mass. We are interested in estimating

$$\text{quantity of interest} = \frac{M_{\bar{I}}}{TM}.$$

We'll compute an upper bound on  $M_{\bar{I}}$  and a lower bound on  $TM$ , to get an upper bound estimate on the quantity of interest. If we can say that QOI is small, then it means that most of the mass lies in the interval  $I$ . Through all of this, we'll make our estimate depend on  $c$ , the width parameter of  $I$ , so that we'll be able to estimate the amount of mass in  $I$  as a function of width.

We now proceed along those lines. Disregarding the proportionality constant, let

$$g(r) := r^{d-2} e^{-\frac{r^2}{2}} r = cr^{d-1} e^{-\frac{r^2}{2}} \quad \text{and} \quad f(r) = \log g(r) = (d-1) \log(r) - \frac{r^2}{2}.$$

Note that

$$f'(r) = \frac{d-1}{r} - r \quad \text{and} \quad f''(r) = -\frac{(d-1)}{r^2} - 1 \leq -1.$$

The second inequality implies that  $f(r)$  is concave (and  $g(r)$  log concave). Hence, solving for  $f'(r) = 0$ , we see that  $f$  is maximized at  $\sqrt{d-1}$ . Applying Taylor's theorem about  $r = \sqrt{d-1}$  and we have

$$\begin{aligned} f(r) &= f(\sqrt{d-1}) + \underbrace{f'(\sqrt{d-1})}_{=0} (r - \sqrt{d-1}) + \frac{1}{2} f''(\zeta) (r - \sqrt{d-1})^2 \\ &\leq f(\sqrt{d-1}) - \frac{1}{2} (r - \sqrt{d-1})^2 \end{aligned} \tag{6.7}$$

for some  $\zeta$  between  $r$  and  $\sqrt{d-1}$ . In the first line, we have applied a form of Taylor's theorem where we explicitly handle the remainder. (This can be found on Wikipedia, and would be a good topic for future notes. On the board, maybe set this part up by first stating the standard form of Taylor's theorem. Then backtrack and note that this is a more convenient form that allows us to explicitly handle the remainder. Though, to be honest, could you just use the standard form?) In the last line, we apply the fact that  $f''(\zeta) \leq -1$ . By the definition of  $f$ , this gives

$$\begin{aligned} g(r) &= \exp(f(r)) \\ &\leq \exp(f(\sqrt{d-1}) - \frac{1}{2} (r - \sqrt{d-1})^2) \\ &= g(\sqrt{d-1}) \exp(-\frac{1}{2} (r - \sqrt{d-1})^2). \end{aligned}$$

As desired, we have a simple “quadratic” estimator on  $g$ . Let  $r_d = \sqrt{d-1}$ , because I'm tired of writing this, and for a constant  $c > 0$  consider the interval

$$I = \{r_d - c, r_d + c\}.$$

Let's try to bound the mass outside this interval using our quadratic overestimator.

$$\begin{aligned} \int_{r \notin I} g(r) dr &\leq \int_0^{r_d - c} g(\sqrt{d-1}) \exp(-\frac{1}{2} (r - \sqrt{d-1})^2) dr + \int_{r_d + c}^{\infty} g(\sqrt{d-1}) \exp(-\frac{1}{2} (r - \sqrt{d-1})^2) dr \\ &\leq 2g(\sqrt{d-1}) \int_{r_d + c}^{\infty} \exp(-\frac{1}{2} (r - \sqrt{d-1})^2) dr \\ &= 2g(\sqrt{d-1}) \int_c^{\infty} \exp(-y^2/2) dy \\ &\leq 2g(\sqrt{d-1}) \int_c^{\infty} \frac{y}{c} \exp(-y^2/2) dy \\ &= \frac{2}{c} g(\sqrt{d-1}) \exp(-c^2/2). \end{aligned}$$

In the first line, we use the upper bound on  $g$ . In the second line, we use the fact that the upper bound is symmetric about  $r_d$ , so integrating over the right half interval is at least as large as the left interval (which is truncated at zero). In the fourth line we use the fact that  $y_i \geq 1$  in the interval of integration. (We do this with a change of variables in mind.) In the last line, we apply change of variables. TODO: Actually check the application of COV in the last step. This will be our estimate on  $M_I$ .

Now, let's come up with a lower bound on the total mass. We'll do this by considering only the mass of  $g(r)$  in the subinterval  $[r_d - c, r_d]$ . For  $r$  in this subinterval, we have  $f''(r) \leq -2$ . This is confirmed by noting that  $f''$ , as explicitly computed above, is monotonically increasing for  $r > 0$ , and then evaluating  $f''$  at the left endpoint of the interval. Applying (6.7), we have

$$f(r) \geq f(\sqrt{d-1}) - (r - \sqrt{d-1})^2 \geq f(\sqrt{d-1}) - \frac{c^2}{4}$$

for  $r$  in the designated subinterval. Equivalently, this gives

$$g(r) \geq g(\sqrt{d-1}) \exp(-\frac{c^2}{4}).$$

Applying this (and noting that the bound has no dependence on  $r$ ), we see that the total mass is at least

$$\int_{r_d-c}^{r_d} g(r) dr \geq cg(\sqrt{d-1}) \exp(-\frac{c^2}{4}).$$

Finally, we see that

$$\text{QOI} = \frac{M_I}{TM} \leq \frac{\frac{2}{c}g(\sqrt{d-1}) \exp(-\frac{c^2}{4})}{cg(\sqrt{d-1}) \exp(-\frac{c^2}{4})} = \frac{2}{c^2} \exp(-\frac{c^2}{4}).$$

This proves the following lemma (wording directly copied from Blume/Kannan/Hopcroft book).

**Lemma 2.** *For a  $d$ -dimensional spherical Gaussian of variance 1, all but  $\frac{2}{c^2} \exp(-\frac{c^2}{4})$  fraction of its mass is within the annulus  $\sqrt{d-1} - c \leq r \leq \sqrt{d-1} + c$ .*

Some concrete implications: Independent of  $d$ , .99 fraction of the mass is contained in the annulus with  $c = 3.38$ . and .999 fraction of the mass is contained in the annulus with  $c = 4.32$ .

**Can we do better?** This estimate holds for  $d$  arbitrarily large. One may wonder if this is a conservative estimate. Perhaps the mass actually concentrates onto an arbitrarily narrow annulus as  $d \rightarrow \infty$ ? My experimentation suggests that this is not the case. I think that the width of the high mass region is genuinely “constant” independent of  $d$ .

Example: To generate samples from a traditional VAE, one draws samples from a unit gaussian and then passes these through a neural network. For moderately high dimensional latent space, this means that the samples are roughly drawn from the surface of a unit sphere. A question one might ask is: How does this affect interpolation? One of the most useful properties of a VAE is the ability to interpolate in latent space. But, if you draw two samples from a high dimensional gaussian and consider the euclidean “straight line” interpolation between them, are you passing through regions where there is very little mass, and so you probably have few samples? Or do your interpolations tend to stay close to the surface of the sphere as well? Is there a better way/different geometry to interpolate rather than straight lines?

Example: TODO: Give an information theory example, like from the Shannon paper about communication in the presence of noise?

TODO: Show some simulation examples.

Todos

- Geometry of high dimensional spheres and cubes
- Gaussian in high dimensions (characterize. Maybe there's a reference in one of those old papers about characterizing local minima of neural networks.)
- Law of unconscious statistician (do this and polar and Gaussian constant as applications of change of variables. Then reference all kinds of generative models.)

- Brachistochrone problem
- integration by parts and connections with FTC.
- review of IFT and connection to inverse function theorem? And connections to FTC?