

Far from complete and far from correct. These notes exist because I learn better when I express things in my own words. In some cases I might be assembling material for teaching purposes. Might delete later.

balthasar@rhizomeworks.com

Contents

1	Combinatorics	7
1.1	Sets	7
1.1.1	Set Operations	7
1.1.2	Common Sets	7
1.1.3	DeMorgan's Rules	7
1.1.4	Inclusion - Exclusion Principle	8
1.2	Combinatorial Identities and Expansions	8
1.2.1	Binomial Coefficients and Binomial Expansions	8
1.2.2	Multinomial Expansion	10
1.2.3	Unnamed Polynomial Identity	10
1.2.4	Factorial Expansion	10
1.2.5	Stirling Numbers of the Second Kind	11
1.3	Distributions: The Twentyfold Way	11
1.3.1	Distinct Objects	13
1.3.2	Distinct Objects, Every Recipient Receives At Most One	13
1.3.3	Distinct Objects, Every Recipient Receives at Least One	14
1.3.4	Distinct Objects, Every Recipient Receives Exactly One	15
1.3.5	Distinct Objects, Distributed in Ordered Groups	16
1.3.6	Distinct Objects, Distributed in Ordered Groups of At Least One	17
1.3.7	Identical Objects	17
1.3.8	Identical Objects, Each Receives At Most One	18
1.3.9	Identical Objects, Each Receives At Least One	19
1.3.10	Identical Objects, Each Receives Exactly One	19
1.4	Generating Functions	20
1.4.1	Example: Binomial Coefficients	20
1.4.2	Example: Basket of Goods	20
1.4.3	Example: Dice	21
2	Linear Algebra and Multivariable Calculus	23
2.1	Multi-Index Notation	23
2.1.1	Example: Multinomial Coefficients	24
2.1.2	Example: Taylor Expansion	24
2.2	Linear Systems of Equations	24
2.2.1	$A \in \mathbb{R}^{n \times n}$ Square Matrices	24
2.2.2	$A \in \mathbb{R}^{m \times n}$ Rectangular Matrices, Overdetermined Case	24
2.2.3	$A \in \mathbb{R}^{n \times m}$ Rectangular Matrices, Underdetermined Case	25
2.3	$A = LL^{\dagger}$ Cholesky Decomposition	25
2.3.1	Creating Correlated Random Variables	25
2.4	Generalized Eigenvectors	25
2.5	$A = V\Lambda V^{-1}$ Spectral Theorems, Diagonalization	25
2.5.1	$A = V\Lambda V^T$ Eigendecomposition of Symmetric Matrices	25
2.5.2	$H = U\Lambda U^T$ Eigendecomposition of Hermitian Matrices	27
2.5.3	Eigenvalue Sensitivity and Accuracy	27
2.6	$A = U\Sigma V^{\dagger}$ Singular Value Decomposition	27
2.6.1	Full and Economy SVDs	28
2.6.2	Matrix Approximation	28

2.7	Types of Transformations	29
2.7.1	Similarity Transformations	29
2.7.2	Affine Transformations	29
2.7.3	Unitary Transformations	29
2.7.4	Multilinear Maps	30
2.7.5	Multilinear Forms	30
2.8	Types of Matrices	30
2.8.1	$\text{sgn}(\mathbf{x}^\dagger \mathbf{H} \mathbf{x})$ Definite	30
2.8.2	Triangular	30
2.8.3	$\mathbf{AB} - \mathbf{BA} = 0$ Commuting	30
2.8.4	$\mathbf{AB} + \mathbf{BA} = 0$ Anticommuting	31
2.8.5	$\mathbf{A}^\dagger = \mathbf{A}$ Hermitian, Symmetric	31
2.8.6	$\mathbf{A}^\dagger = -\mathbf{A}$ Skew Hermitian, Skew Symmetric	31
2.8.7	$\mathbf{AA} = \mathbf{I}$ Involutory	31
2.8.8	$\mathbf{U}^\dagger = \mathbf{U}^{-1}$ Unitary, Orthogonal	32
2.8.9	$\mathbf{A} = \mathbf{TBT}^{-1}$ Similarity	32
2.9	Properties of Norms	32
2.10	L^p Lebesgue Vector Norms	32
2.10.1	L^1 Taxicab / Manhattan Norm	32
2.10.2	L^2 Euclidian Norm	32
2.10.3	L^∞ Maximum Norm	32
2.10.4	$L^{-\infty}$ Minimum Norm	32
2.11	Operator and Matrix Norms	33
2.11.1	$\ \mathbf{A}\ _{(\alpha)}$ Operator Norm	33
2.11.2	$\ \mathbf{A}\ _q$ q -Norms	33
2.11.3	$\ \mathbf{A}\ _F$ Frobenius Norm	34
2.11.4	$\ \mathbf{A}\ _{(1)}$ (1)-Norm	34
2.11.5	$\ \mathbf{A}\ _{(\infty)}$ (∞)-Norm	34
2.11.6	$\ \mathbf{A}\ _{(2)}$ (2)-Norm	34
2.12	Vector and Matrix Derivatives	34
2.12.1	Jacobian	34
2.12.2	Inverse Function Theorem	35
2.12.3	Critical Points	35
2.12.4	Differential Volume Element, Change of Variables	35
2.12.5	Hessian	36
3	Probability	37
3.1	Interpretations and Definitions of Probability	37
3.1.1	Naive and Non-Naive Definitions of Probability	37
3.1.2	Frequentist Probability	38
3.1.3	Bayesian Probability	38
3.1.4	Measure Theoretic Probability	38
3.2	Measure Theory	38
3.3	Functions of Random Variables, Derived Distributions	38
3.3.1	Example: Sum of Random Variables	39
3.3.2	Example: Lower Dimensional Random Variable	40
3.4	Indicator Variables	40
3.4.1	Example: The Party Problem	41
3.5	Copulas	44
3.6	Relationships Between Distributions	44
3.7	Large Deviation Theory	44
3.7.1	Gaertner-Ellis Theorem	44
3.7.2	Example: Sum of Uniform Random Variables	44

4	Information Theory	45
4.1	Entropy	45
4.2	Mutual Information	45
4.3	Kullback-Leibler Divergence	45
5	Stochastic Processes and Time Series Analysis	47
5.1	Markov Chains	47
5.2	Martingales	47
5.2.1	Martingale Convergence Theorem	47
5.3	Hidden Markov Models	47
5.4	Ito Calculus	47
6	Statistics	49
6.1	Directional Statistics	49
6.1.1	Mean Direction	49
6.1.2	Dispersion	49
6.2	Features	50
6.2.1	Dense Features, Sparse Features	50
6.3	Multicollinearity	50
6.3.1	Sources	50
6.4	R^2 Value	51
6.5	Regression Diagnostics	51
6.6	t-Statistics	51
6.7	AIC and BIC	51
6.8	Factor Regression	51
6.9	Factor Models	51
7	Linear Regression	53
7.1	Ordinary Least Squares Regression	53
7.1.1	The Normal Equations, Analytical Least Squares Estimator	53
7.1.2	The Quick Way to $\hat{\beta}$	54
7.1.3	The Long Way to $\hat{\beta}$	54
7.1.4	Projection Matrix	54
7.1.5	Bayesian Perspective on Least Squares Regression	55
7.1.6	Q-plots	55
7.1.7	Variance Inflation Factor	55
7.2	Ridge Regression	55
7.2.1	Analytical Ridge Estimator	55
7.2.2	Bayesian Perspective on Ridge Regression	55
8	Bayesian Data Analysis	57
9	Unsupervised Learning	59
9.1	Blind Source Separation	59
9.2	Independent Component Analysis	59
9.3	Eigenanalysis, Singular Value Decomposition of Data Matrices	59
9.3.1	Example: Eigenfaces and Facial Recognition	60
9.3.2	Tracking an Eigensystem over Time	61
9.3.3	Bayesian SVD	61
9.4	Principal Component Analysis (PCA)	61
9.4.1	Relationship between PCA and SVD	63
9.4.2	Tracking Principal Components over Time	63
9.5	$\mathbf{X} = \mathbf{WH}$ Non-Negative Matrix Factorization	64

1 Combinatorics

Contents of this chapter

1.1	Sets	7	1.3.4	Distinct Objects, Every Recipient Receives Exactly One	15
1.1.1	Set Operations	7	1.3.5	Distinct Objects, Distributed in Ordered Groups	16
1.1.2	Common Sets	7	1.3.6	Distinct Objects, Distributed in Ordered Groups of At Least One	17
1.1.3	DeMorgan's Rules	7	1.3.7	Identical Objects	17
1.1.4	Inclusion - Exclusion Principle	8	1.3.8	Identical Objects, Each Receives At Most One	18
1.2	Combinatorial Identities and Expansions	8	1.3.9	Identical Objects, Each Receives At Least One	19
1.2.1	Binomial Coefficients and Binomial Expansions	8	1.3.10	Identical Objects, Each Receives Exactly One	19
1.2.2	Multinomial Expansion	10	1.4	Generating Functions	20
1.2.3	Unnamed Polynomial Identity	10	1.4.1	Example: Binomial Coefficients	20
1.2.4	Factorial Expansion	10	1.4.2	Example: Basket of Goods	20
1.2.5	Stirling Numbers of the Second Kind	11	1.4.3	Example: Dice	21
1.3	Distributions: The Twentyfold Way	11			
1.3.1	Distinct Objects	13			
1.3.2	Distinct Objects, Every Recipient Receives At Most One	13			
1.3.3	Distinct Objects, Every Recipient Receives at Least One	14			

1.1 Sets

Sets are a term to describe collections of things. The things could be countable objects, such as the integers between 1 and 10, or contain a continuum, for example all the real numbers between 1 and 10.

1.1.1 Set Operations

1.1.2 Common Sets

1.1.3 DeMorgan's Rules

DeMorgan's Rules relate the complement of the union to the intersection of the complements, and the complement of the intersection to the union of the complements.

$$\left(\bigcup_{i \in \{i\}} A_i \right)^c = \bigcap_{i \in \{i\}} A_i^c \quad (1.1)$$

$$\left(\bigcap_{i \in \{i\}} A_i \right)^c = \bigcup_{i \in \{i\}} A_i^c \quad (1.2)$$

1.1.4 Inclusion - Exclusion Principle

The inclusion-exclusion principle is used to calculate the size of the union of sets. This requires counting each region of some complicated overlapping Venn diagram exactly once, which, in turn requires accounting for overcounting wherever sets overlap. Let $\{A_i | i \in \{i\}_n\}$ be a collection of n overlapping sets indexed by $i \in \{i\}_n$, then the inclusion-exclusion principle is given by:

$$\left| \bigcup_{i \in \{i\}_n} A_i \right| = \sum_{k=1}^n (-1)^{k-1} \sum_{\{j\}_k \subseteq \{i\}_n} \left| \bigcap_{j \in \{j\}_k} A_j \right| \quad (1.3)$$

Where the sum over $\{j\}_k \subseteq \{i\}_n$ is over all k -element subsets of $\{i\}_n$.

1.1.4.1 Example: n=2 Sets and n=3 Sets

n=2

$$\begin{aligned} \left| \bigcup_{i \in \{1,2\}} A_i \right| &= \sum_{k=1}^2 (-1)^{k-1} \sum_{\{j\}_k \subseteq \{1,2\}} \left| \bigcap_{j \in \{j\}_k} A_j \right| \\ &= (-1)^0 (|A_1| + |A_2|) \\ &\quad + (-1)^1 (|A_1 \cap A_2|) \end{aligned} \quad (1.4)$$

n=3

$$\begin{aligned} \left| \bigcup_{i \in \{1,2,3\}} A_i \right| &= \sum_{k=1}^3 (-1)^{k-1} \sum_{\{j\}_k \subseteq \{1,2,3\}} \left| \bigcap_{j \in \{j\}_k} A_j \right| \\ &= (-1)^0 (|A_1| + |A_2| + |A_3|) \\ &\quad + (-1)^1 (|A_1 \cap A_2| + |A_1 \cap A_3| + |A_2 \cap A_3|) \\ &\quad + (-1)^2 (|A_1 \cap A_2 \cap A_3|) \end{aligned} \quad (1.5)$$

1.1.4.2 Example: Counting Integers

How many integers are there between 1 and 100 that are neither divisible by 3, 5 nor 7?

Let S be the set of all integers between 1 and 100. The size of the set is $|S| = 100$. The subset of S that is numbers divisible by 3 is $A_3 \subseteq S$ with $|A_3| = 33$ because $100/3 = 33.\overline{333}$. Similarly, $|A_5| = 20$ and $|A_7| = 14$. The set of integers that is not divisible by 3, 5 or 7 is:

$$S \setminus \bigcup_{i \in \{3,5,7\}} A_i \quad (1.6)$$

So that the sought after quantity is :

$$\begin{aligned} \left| S \setminus \bigcup_{i \in \{3,5,7\}} A_i \right| &= |S| - [|A_3| + |A_5| + |A_7| \\ &\quad - |A_3 \cap A_5| - |A_3 \cap A_7| - |A_5 \cap A_7| + |A_3 \cap A_5 \cap A_7|] \end{aligned} \quad (1.7)$$

The size of the intersection $|A_3 \cap A_5| = 6$ because 100 is 6 times divisible by $3 \times 5 = 15$. Similarly, $|A_3 \cap A_7| = 4$, $|A_5 \cap A_7| = 2$ and $|A_3 \cap A_5 \cap A_7| = 0$. Hence:

$$\left| S \setminus \bigcup_{i \in \{3,5,7\}} A_i \right| = 100 - 33 - 20 - 14 + 6 + 4 + 2 - 0 = 45 \quad (1.8)$$

There are 45 integers between 1 and 100 that are not divisible by 3, 5 or 7.

1.2 Combinatorial Identities and Expansions

1.2.1 Binomial Coefficients and Binomial Expansions

For two positive integers n and k , the binomial coefficient " n choose k " is:

$$\binom{n}{k} = \begin{cases} \frac{n!}{k!(n-k)!} & \text{for } n \geq k \\ 0 & \text{otherwise} \end{cases} \quad (1.9)$$

The term can be defined for negative arguments, which comes up often when working with generating functions.

$$\binom{-n}{k} = \begin{cases} (-1)^k \binom{n+k-1}{k} & \text{for } n \geq k \\ 0 & \text{otherwise} \end{cases} \quad (1.10)$$

$$\binom{-n}{-k} = \begin{cases} (-1)^{k-n} \binom{k-1}{k-n} & \text{for } n \geq k \\ 0 & \text{otherwise} \end{cases} \quad (1.11)$$

The generalizations can, for example, be derived using symmetry arguments and the Gamma function, which is the generalization of the factorial to non-integers (cf. Kronenburg (2011)).

The binomial expansion can be proven either by expanding the polynomial or by creating the Taylor series for the polynomial.

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^k y^{n-k} \quad (1.12)$$

This holds also for negative integer exponents n , in which case:

$$\frac{1}{(y+x)^n} = \sum_{k=0}^{\infty} \binom{-n}{k} x^k y^{-n-k} = (-1)^k \binom{n+k-1}{k} x^k y^{-(n+k)} \quad (1.13)$$

1.2.1.1 Derivation of the Binomial Theorem for a Negative Exponent

Let $f(x) = \frac{1}{(y+x)^n} = (y+x)^{-n}$. The Taylor expansion about the point $x = 0$ is:

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x=0)}{k!} x^k \quad (1.14)$$

The derivatives of f are:

$$\begin{aligned} f^{(0)}(x) &= (y+x)^{-n} \\ f^{(1)}(x) &= (y+x)^{-(n+1)} (-1)n \\ f^{(2)}(x) &= (y+x)^{-(n+2)} (-1)^2 n(n+1) \\ &\vdots \\ f^{(k)}(x) &= (y+x)^{-(n+k)} (-1)^k n(n+1) \dots (n+k-1) \\ &= (y+x)^{-(n+k)} (-1)^k \frac{(n+k-1)!}{(n-1)!} \end{aligned} \quad (1.15)$$

Combining Eqns. 1.14 and 1.15 gives:

$$f(x) = \sum_{k=0}^{\infty} (-1)^k \binom{n+k-1}{k} x^k y^{-(n+k)} \quad (1.16)$$

1.2.1.2 Derivation of the Binomial Theorem for a Fractional Exponent

Following much the same logic as for a negative exponent, for a fractional exponent the Taylor series is also infinite.

Let $f(x) = (y+x)^m$ with $m \notin \mathbb{Z}$.

The derivatives of f are:

$$\begin{aligned}
f^{(0)}(x) &= (y+x)^m \\
f^{(1)}(x) &= (y+x)^{m-1}m \\
f^{(2)}(x) &= (y+x)^{m-2}m(m-1) \\
&\vdots \\
f^{(k)}(x) &= (y+x)^{m-k}m(m-1)\dots(m-k+1)
\end{aligned} \tag{1.17}$$

Combining Eqns. 1.14 and 1.17 gives:

$$f(x) = \sum_{k=0}^{\infty} \binom{m}{k} x^k y^{m-k} \tag{1.18}$$

Where, boldly, I defined $\binom{m}{k}$ to mean:

$$\binom{m}{k} = \frac{m(m-1)(m-2)\dots(m-k+2)(m-k+1)}{k!} \tag{1.19}$$

1.2.2 Multinomial Expansion

$$(x_1 + x_2 + x_3 + \dots + x_k)^n = \sum_{\substack{i_1, i_2, i_3, \dots, i_k \\ \sum_j i_j = n}} \binom{n}{i_1, i_2, \dots, i_k} x^{i_1} x^{i_2} x^{i_3} \dots x^{i_k} \tag{1.20}$$

with:

$$\binom{n}{i_1, i_2, \dots, i_k} = \frac{n!}{i_1! i_2! i_3! \dots i_k!} \tag{1.21}$$

Where the sum over all possible exponents i_j so that $\sum_j i_j = n$ has $\binom{n+k-1}{n}$ terms.

1.2.3 Unnamed Polynomial Identity

I don't know what this is called, but it's useful.

$$\prod_i^n (1 - x_i) = \sum_{s=0}^n (-1)^s \sum_{\substack{0 \leq i_1, i_2, \dots, i_s \leq n \\ \{i\}_s}} \prod_{i \in \{i\}_s} x_i \tag{1.22}$$

Where $\{i\}_s$ is a set of s indices that range between 0 and n , and the sum is over all possible such sets, of which there are $\binom{n}{s}$.

1.2.4 Factorial Expansion

$$x^n = \frac{x!}{(x-n)!} = \sum_{k=0}^n s(n, k) x^k \tag{1.23}$$

where

$$s(n, k) = (-1)^{n-k} \left[\begin{matrix} n \\ k \end{matrix} \right] \tag{1.24}$$

are the stirling numbers of the first kind.

1.2.5 Stirling Numbers of the Second Kind

Stirling numbers of the second kind $S(k, n)$ measure the amount of ways in which k objects can be divided into n non-empty groups. They give the number of onto functions from a set of k distinct objects to n indistinct recipients. For example: how many ways can a set of k pool balls be put into n bags, so that there is at least one ball in each bag. All the pool balls have numbers on them and have different colors, so that $n = 2$ bags containing $[(1, 2, 3), (4)]$ and $[(1, 2, 4), (3)]$ count as different. This sort of problem is discussed at length in section 1.3

They are given by an explicit formula:

$$S(k, n) = \frac{1}{n!} \sum_{i=1}^n (-1)^{n-i} \binom{n}{i} i^k \quad (1.25)$$

They can also be generated via the recurrence relation:

$$\left\{ \begin{matrix} k+1 \\ n \end{matrix} \right\} = n \left\{ \begin{matrix} k \\ n \end{matrix} \right\} + \left\{ \begin{matrix} k \\ n-1 \end{matrix} \right\} \quad (1.26)$$

The recurrence relation is explained by adding the combinations corresponding to two cases. If the $k + 1$ st object is added to one of the n existing subsets with k objects, then that corresponds to:

$$n \left\{ \begin{matrix} k \\ n \end{matrix} \right\} = 1 \quad (1.27)$$

Possibilities. If the $k + 1$ st object is in a set by itself (a singleton), then the remaining objects are distributed over $n - 1$ set. The combinations arising from this are:

$$\left\{ \begin{matrix} k \\ n-1 \end{matrix} \right\} = 1 \quad (1.28)$$

Furthermore, the following holds:

$$\left\{ \begin{matrix} 0 \\ 0 \end{matrix} \right\} = 1 \quad (1.29)$$

$$\left\{ \begin{matrix} k \\ 0 \end{matrix} \right\} = \left\{ \begin{matrix} 0 \\ n \end{matrix} \right\} = 0 \quad (1.30)$$

And $S(k, n) = 0$ if $n > k$.

1.3 Distributions: The Twentyfold Way

The twentyfold way is a taxonomy of distribution problems developed by Kenneth Bogart in his book *Combinatorics through Guided Discovery* (Bogart 2004). It divides up the way in which k objects may be assigned to n individuals, subject to whether the objects are distinct or identical, and subject to conditions on how the objects are received.

*When we are passing out objects to recipients, we may think of the objects as being either identical or distinct. We may also think of the recipients as being either identical (as in the case of putting fruit into plastic bags in the grocery store) or distinct (as in the case of passing fruit out to children). We may restrict the distributions to those that give at least one object to each recipient, or those that give exactly one object to each recipient, or those that give at most one object to each recipient, or we may have no such restrictions. If the objects are distinct, it may be that the order in which the objects are received is relevant (think about putting books onto the shelves in a bookcase) or that the order in which the objects are received is irrelevant (think about dropping a handful of candy into a child's trick or treat bag). If we ignore the possibility that the order in which objects are received matters, we have created $2 \times 2 \times 4 = 16$ distribution problems. In the cases where a recipient can receive more than one distinct object, we also have four more problems when the order objects are received matters. Thus we have 20 possible distribution problems. - Bogart, *Combinatorics Though Guided Discovery*, Chapter 3.*

Figure 1.1: Bogart's Twentyfold Way

The Twentyfold Way: A Table of Distribution Problems		
k objects and conditions on how they are received	n recipients and mathematical model for distribution	
	Distinct	Identical
1. Distinct no conditions	n^k functions	$\sum_{i=1}^k S(n, i)$ set partitions ($\leq n$ parts)
2. Distinct Each gets at most one	$n^{\underline{k}}$ k -element permutations	1 if $k \leq n$; 0 otherwise
3. Distinct Each gets at least one	$S(k, n)n!$ onto functions	$S(k, n)$ set partitions (n parts)
4. Distinct Each gets exactly one	$k! = n!$ permutations	1 if $k = n$; 0 otherwise
5. Distinct, order matters	$(k + n - 1)^{\underline{k}}$ ordered functions	$\sum_{i=1}^n L(k, i)$ broken permutations ($\leq n$ parts)
6. Distinct, order matters Each gets at least one	$(k)^{\underline{n}}(k - 1)^{\underline{k-n}}$ ordered onto functions	$L(k, n) = \binom{k}{n}(k - 1)^{\underline{k-n}}$ broken permutations (n parts)
7. Identical no conditions	$\binom{n+k-1}{k}$ multisets	$\sum_{i=1}^n P(k, i)$ number partitions ($\leq n$ parts)
8. Identical Each gets at most one	$\binom{k}{n}$ subsets	1 if $k \leq n$; 0 otherwise
9. Identical Each gets at least one	$\binom{k-1}{n-1}$ compositions (n parts)	$P(k, n)$ number partitions (n parts)
10. Identical Each gets exactly one	1 if $k = n$; 0 otherwise	1 if $k = n$; 0 otherwise

Table 3.3.4: The number of ways to distribute k objects to n recipients, with restrictions on how the objects are received

What I like about this approach is that the challenge with most of the basic combinatorics problems is to figure out the right way of counting. For this reason, the idea of having a unified handbook-like taxonomy is very appealing. The weakness (in my opinion) is that the language of "objects" and "recipients" is unclear because in practice it's not obvious which is which: if there are k students and n teachers, do the teachers receive students, or do the students receive a teacher?

A way to resolve this is to say that an object can have only one recipient, but that a recipient might receive more than one object. A more formal path is to think of the act of creating combinations in terms of functions.

- The elements of the domain are the objects.
- The elements of the range are the recipients.
- A function can be many-to-one, but it should not be one-to-many.

Favorite Teachers At a school with k students and n teachers, the students all have a favorite teacher. (They might all like the same one.). How many ways are there for all of the k students to pick a favorite?

Objects: k students. *Recipients:* n teachers. Many students might have one favorite teacher. There are n^k combinations.

Assembling a Team Out of a choice of n athletes, a coach must assemble a team of k . How many ways are there to form a team?

Objects: n athletes. *Recipients:* team, not on the team. Many athletes can be assigned to one outcome of being on the team or not being on the team. There are $\binom{n}{k}$ combinations for the team, which is the same number as the $\binom{n}{n-k}$ selections for the bench.

1.3.1 Distinct Objects

1.3.1.1 Distinct Recipients

The k objects are assigned to n recipients with no conditions as to the number of objects each recipient receives. This is the same as assigning the elements of a k -tuple from a selection of n with replacement.

$$\begin{aligned} S &= \{(a_1, a_2, \dots, a_k) | a_i \in A, |A| = n\} \\ |S| &= n^k \end{aligned} \tag{1.31}$$

Pool Balls into Labeled Buckets All possible ways to put k pool balls, which all have different numbers and colors, into n labeled buckets. Some of the buckets might be empty, and others might contain more than one of the pool balls.

Functions All possible functions $f : x \rightarrow y$ with $\{x | x \in A, |A| = k\}$ and $\{y | y \in B, |B| = n\}$.

Binary Strings of Length k The k distinct positions of a binary string (i_1, i_2, \dots, i_k) of length k are assigned to an element of the set $A \in [0, 1]$. The number of possible binary strings of length k is 2^k .

Subsets of a k -Element Set The subsets of a set of k distinct elements are formed by assigning each of its k distinguishable elements to one of the two labels $A \in [\text{included}, \text{excluded}]$. The number of possible subsets, including the empty subset and the full set, is 2^k .

1.3.1.2 Indistinct Recipients

The k objects are assigned to a recipient that is not distinct.

$$|A| = \sum_{i=1}^k S(n, i) \tag{1.32}$$

Where $S(k, n)$ is the Stirling number of the second kind that gives the number of ways that k objects can be distributed across n non-empty indistinct sets. The sum above takes care of the case where the k objects are divided into up to n collections.

A closed form expression for the Stirling Numbers of the second kind is (c.f. section 1.2.5):

$$S(k, n) = \left\{ \begin{matrix} k \\ n \end{matrix} \right\} = \frac{1}{n!} \sum_{j=0}^n (-1)^{n-j} \binom{n}{j} j^k \tag{1.33}$$

Pool Balls into Unlabeled Bags All possible ways to put k pool balls, which all have different numbers and colors, into n unlabeled bags. Some of the bags might be empty, and others might contain more than one of the pool balls.

1.3.2 Distinct Objects, Every Recipient Receives At Most One

1.3.2.1 Distinct Recipients

At most one of k distinct objects are assigned to one of n distinct recipients. This is the same as assigning the elements of a k -tuple from a selection of n without replacement, so that first there are n choices, then $n - 1$ choices, $n - 2$, etc.

$$\begin{aligned} S &= \{(a_1, a_2, \dots, a_k) | a_i \in A, |A| = n, a_i \neq a_j\} \\ |S| &= \frac{n!}{(n-k)!} = n^k \text{ if } k \leq n, 0 \text{ otherwise.} \end{aligned} \tag{1.34}$$

At Most One Pool Ball into Labeled Buckets All possible ways to put k pool balls, which all have different numbers and colors, into n labeled buckets, but the buckets can have at most one pool ball in them. In other words, you choose any k out of the n labeled buckets and put one pool ball into them.

One-to-One Functions All possible functions $f : x \rightarrow y$ with $\{x|x \in A, |A| = k\}$ and $\{y|y \in B, |B| = n\}$ subject to the constraint that $f(a) = f(b)$ implies $a = b$. That is, the functions are one-to-one, or injective.

k -element Permutations of n elements Each of the k positions in a k -element permutation are distinct objects. These are each assigned to one of n possible values, where each value can only show up once.

Books on a Shelf How many ways are there to order k books on a library shelf when there are n different books available.

1.3.2.2 Indistinct Recipients

At most one of k distinct objects are assigned to one of n indistinct recipients. This is the same as assigning the elements of a k -tuple from a selection of n without replacement. Except, since the recipients are all indistinct, there is only one type of choice of recipient for each object. Either each object finds a recipient if $k \leq n$, or it is impossible to distribute at most one object to each recipient because $n < k$.

$$S = \{(a_1, a_2, \dots, a_k) | a_i \in A, |A| = n, a_i = a_j\} \quad (1.35)$$

$$|S| = 1 \text{ if } k \leq n, 0 \text{ otherwise.}$$

At Most One Pool Ball into Unlabeled Bags All possible ways to put k pool balls, which all have different numbers and colors, into n unlabeled bags, but the bags can have at most one pool ball in them. The bags are identical, so the only result is one where there are k bags with a ball in them and $n - k$ without a ball in them. If there are not enough bags, then there is no possible result.

Distributing Candy There are n pieces of identical candy and k kids. How many ways are there to give each kid a piece of candy? If there is enough candy, the answer is one. Everyone gets candy. If there is not enough candy then the answer is zero. There is no way to give everyone candy if there's not enough candy.

1.3.3 Distinct Objects, Every Recipient Receives at Least One

1.3.3.1 Distinct Recipients

$$|A| = n!S(k, n) = n! \left\{ \begin{matrix} k \\ n \end{matrix} \right\} = n! \frac{1}{n!} \sum_{j=0}^n (-1)^{n-j} \binom{n}{j} j^k \quad (1.36)$$

Where $S(k, n)$ denotes the Stirling function of the second kind.

At Least One Pool Ball into Labeled Buckets All possible ways to put k pool balls, which all have different numbers and colors, into n labeled buckets, so that all of the buckets have at least one ball in them. This is the same number of combinations as if the buckets were unlabeled, but multiplied with $n!$ ways of applying a label to them.

Onto Functions All possible functions $f : x \rightarrow y$ with $\{x|x \in A, |A| = k\}$ and $\{y|y \in B, |B| = n\}$ subject to the constraint that there is an element x in the domain so that $f(x) = y$ for each element y of the codomain. That is the functions are onto, or surjective.

1.3.3.2 Indistinct Recipients

The number of ways to divide k distinct objects into n non-empty subsets is given by the Stirling number of the second kind (c.f. section 1.2.5):

$$|A| = S(k, n) = \left\{ \begin{matrix} k \\ n \end{matrix} \right\} = \frac{1}{n!} \sum_{j=0}^n (-1)^{n-j} \binom{n}{j} j^k \quad (1.37)$$

At Least One Pool Ball into Unlabeled Bags All possible ways to put k pool balls, which all have different numbers and colors, into n unlabeled bags, so that none of the bags are empty.

1.3.4 Distinct Objects, Every Recipient Receives Exactly One

1.3.4.1 Distinct Recipients

One of k distinct objects are assigned to each of n distinct recipients. This is the same as assigning the elements of a k -tuple from a selection of n without replacement and with the requirement that all of the n are selected.

$$S = \{(a_1, a_2, \dots, a_k) | a_i \in A, |A| = k = n, a_i \neq a_j\} \quad (1.38)$$
$$|S| = n! = k! \text{ if } k = n, 0 \text{ otherwise.}$$

Exactly One Pool Ball into Each Labeled Buckets All possible ways to put k pool balls, which all have different numbers and colors, into n labeled buckets, that there is one pool ball in each bucket. This is the same as putting one pool ball into unlabeled buckets and multiplying by the $n!$ ways of attaching a label to the buckets. If there isn't the same amount of balls and buckets then there is no possible way for them to be matched one-for-one.

Bijective Functions All possible functions $f : x \rightarrow y$ with $\{x | x \in A, |A| = k\}$ and $\{y | y \in B, |B| = n\}$ subject to the constraints that $f(a) = f(b)$ implies $a = b$ and that there is an element x in the domain so that $f(x) = y$ for each element y of the codomain. That is, the functions are one-to-one and onto, of bijective.

Permutations Since each of the k objects is given to a different one of the n recipients, there must be as many recipients as there are objects and $k = n$. The number of ways of assigning k objects to n recipients is $k! = n!$.

Unique Identifiers Each of k entries in a database is given one of $n = k$ unique identifiers, so that each identifier leads to an entry and each entry has an identifier.

1.3.4.2 Indistinct Recipients

$$S = \{(a_1, a_2, \dots, a_k) | a_i \in A, |A| = k, a_i = a_j\} \quad (1.39)$$
$$|S| = 1 \text{ if } k = n, 0 \text{ otherwise.}$$

Exactly One Pool Ball into Each Unlabeled Bag All possible ways to put k pool balls, which all have different numbers and colors, into n unlabeled bags, that there is one pool ball in each bag. The result is that you either have n bags with $n = k$ balls in them, if the two have matching numbers, or that you either have a bag or a ball left over and there is no way to match them one-for-one.

Distribute without Leftovers k students are assigned to n identical textbooks. How many ways are there for each child to have a textbook so that there are no textbooks left over?

1.3.5 Distinct Objects, Distributed in Ordered Groups

1.3.5.1 Distinct Recipients

k objects are distributed to n different recipients with an internal ordering, so that each recipient receives an ordered list. This is the same as creating a list of n sequences that are sampled from k without replacement.

$$S = \{(\mathbf{a}_{\{i_1\}} = (a_{i_1,1}, a_{i_1,2}, \dots, a_{i_1,l}), \mathbf{a}_{\{i_2\}}, \dots, \mathbf{a}_{\{i_n\}}) \mid \sum_j |\mathbf{a}_{i_j}| = k, a_{i_j,i} \in A, |A| = n\}$$

$$|S| = \frac{(k+n-1)!}{(n-1)!} = (k+n-1)\underline{k} = k! \binom{k+n-1}{k} \quad (1.40)$$

Books on Labeled Bookshelves k books are distributed across n different bookshelves. The books may all be on the same shelf, or shelves may be empty. The ordering of the books on each of the shelves matters. This is the same as taking all the $k!$ permutations of the k books and multiplying it by the way of dividing that permutation up onto n shelves.

Ordered Functions All functions $f : x \rightarrow y$ with $\{x \mid x \in A, |A| = k\}$ and $\{y \mid y \in B, |B| = n\}$ that assign ordered sequences of elements in x to elements of y .

1.3.5.2 Indistinct Recipients

k objects are broken up into n

$$|S| = \sum_{i=1}^n L(k, i) \quad (1.41)$$

where $L(k, n)$ are the Lah numbers, which describe "How many ways can k objects be distributed to n recipients if order matters and each recipient receives at least 1". They are given by:

$$L(k, n) = \binom{k}{n} (k-1)^{\underline{k-n}} = \frac{k!}{n!} \binom{k-1}{n-1} \quad (1.42)$$

The sum considers the case where all k objects are given to $i \leq n$ recipients, with the remaining recipients receiving none. Note that the expression is more complicated than simply dividing the case for indistinct recipients by $n!$, because more than one of the recipients might receive the same number of objects, in which case i.e. the distribution $[(0), (0), (1, 2)]$ would be counted twice. The expression is analogous to the expression for distinct objects distributed to indistinct recipients, but where the objects did not have an internal ordering. In that case, instead of a sum over the Lah numbers, the sum is over the Stirling numbers of the second kind.

The Lah numbers can be derived by imagining that first one of the k distinct objects is distributed to the n indistinct recipients, to make sure that each recipient receives at least one, and then adding the remaining objects to the first without restrictions. After the first n distinct objects have been distributed to the recipients, the recipients are no longer indistinct, because they each have been labeled by the object they have already received. There are $\binom{k}{n}$ ways of distributing the first n objects and $(k-n)! \binom{k-n+n-1}{k-n}$ ways to add the $k-n$ remaining objects.

Books into Unlabeled Boxes k books are stacked into up to n different unlabeled boxes. Some of the boxes may be empty, and others may contain more than one book. The sequence in which the books are stacked in each box matters. If there are $n-r$ empty boxes and r boxes that have at least one book in them, then there are $\binom{k}{r}$ ways of putting the first book in each of the boxes. Then there are $(k-r)! \binom{k-r+r-1}{k-n}$ ways to stack the remaining books on top of those first books.

Broken Permutations $\leq n$ Parts The permutations of k distinct elements are ordered sequences of length k . If the sequence is cut up into up to n different parts of non-zero length, then what results are *broken permutations*.

Books into Boxes k different books are put into n identical boxes. How many ways are there to pack the boxes if you keep track of the order in which the books in each box are stacked?

1.3.6 Distinct Objects, Distributed in Ordered Groups of At Least One

1.3.6.1 Distinct Recipients

$$|S| = \frac{k!}{(k-n)!} \frac{(k-1)!}{(n-1)!} = k^n (k-1)^{k-n} \quad (1.43)$$

Books on Labeled Bookshelves k books are distributed across n different labeled bookshelves. There is at least one book on each shelf. This is the same as picking n books to go on the shelves first, so that all of the shelves have at least one book on them, and then distributing the remainder without restrictions. There are k^n ways of picking the first n books for each of the n shelves, and $(k-n+n-1)^{k-n}$ ways to add the remaining $k-n$ books.

Ordered Onto Functions All functions $f : x \rightarrow y$ with $\{x|x \in A, |A| = k\}$ and $\{y|y \in B, |B| = n\}$ that assign ordered sequences of elements in x to elements of y , where every element $y \in B$ has an assignment of at least one element.

1.3.6.2 Indistinct Recipients

k elements are divided into n ordered sequences of minimum length 1.

$$L(k, n) = \binom{k}{n} (k-1)^{k-n} = \frac{k!}{n!} \binom{k-1}{n-1} \quad (1.44)$$

Books into Unlabeled Boxes k books are stacked into up to n different unlabeled boxes. All of the n boxes have at least one book in them. The sequence in which the books are stacked in each box matters. There are $\binom{k}{n}$ ways of putting the first book in each of the boxes. Then there are $(k-n)! \binom{k-n+n-1}{k-n}$ ways to stack the remaining $k-n$ books on top of those first books.

Broken Permutations n parts The permutations of k distinct elements are ordered sequences of length k . If the sequence is cut up into up to n different parts, then what results are *broken permutations*.

Books into Boxes k different books are put into n identical boxes, so that there is at least one book in each box. How many ways are there to pack the boxes if you keep track of the order in which the books in each box are stacked?

1.3.7 Identical Objects

1.3.7.1 Distinct Recipients

$$|S| = \binom{k+n-1}{k} \quad (1.45)$$

This coefficient can be easily derived using the "Stars and Bars" concept. A short way of explaining this is as follows: picture putting all the k identical objects in a line (stars). Picture having dividers (bars) to divide the k objects into n groups. For n groups, it is necessary to use $n-1$ dividers. Empty groups result when the dividers sit right next to each other with no object in between them, or if a divider is at the end of a sequence. In total, the objects and the dividers make for a sequence in which $k+n-1$ spots are taken up by either an object or a divider. To pick a particular way of dividing the k objects up, you can either pick the $n-1$ locations of the sequence in which the dividers are located, or, equivalently, the locations in which the objects are located. There are $\binom{k+n-1}{n-1} = \binom{k+n-1}{k}$ to do so.

Ping Pong Balls into Labeled Buckets k identical ping pong balls are put into n labeled buckets. Some of the buckets may be empty, and other buckets might have more than one ball in them.

Multisets Multisets are sets in which identical elements might show up several times. For example $\{a, a, b, b, b\}$. They can also be described in terms of the multiplicity of their elements. For example $[a : 2, b : 3, c : 0]$. How many multisets can be formed with k different elements of n different classes?

Integer Sums How many different configurations of the n integers $\{x_i\}_n$ satisfy $x_1 + x_2 + \dots + x_n = k$?

Bosons in Degenerate States In how many ways might k Bosons populate n degenerate states?

1.3.7.2 Indistinct Recipients

$$|S| = \sum_{i=1}^n P(k, i) \quad (1.46)$$

It turns out that there is no known formula for $P(k, n)$.

Ping Pong Balls into Unlabeled Bags k identical ping pong balls are put into n unlabeled bags. Some of the bags might be empty, and other bags might have more than one ball in them.

Number Partitions How many ways are there to divide k objects across up to n piles. How many ways are there to divide an integer k into a sum of n integers (including zeros). For example: for $k = 5$, $n = 3$, the partitions are $5 + 0 + 0, 4 + 1 + 0, 3 + 2 + 0, 3 + 1 + 1, 2 + 2 + 1$.

Unlabeled Multiplicities of Multisets For multisets of k elements with n different classes, what is the number of possible multiplicities? For example, a multiset of $k = 3$ elements from $n = 2$ classes could have multiplicities $[a : 3, b : 0], [a : 2, b : 1], [a : 1, b : 2], [a : 0, b : 3]$. If we do not care about the labels a, b , then the ways that the k elements might be distributed are $[3, 0]$ and $[2, 1]$.

Boxes of Marbles k marbles are randomly put into n boxes. How many ways are there for the weight to be distributed among the boxes?

1.3.8 Identical Objects, Each Receives At Most One

1.3.8.1 Distinct Recipients

k identical objects are distributed across n recipients so that each recipient receives at most one. That amounts to choosing k out of the n recipients who will receive an object.

$$|S| = \binom{n}{k} \quad (1.47)$$

At Most One Ping Pong Balls into Labeled Buckets k identical ping pong balls are distributed across n different buckets, so that k buckets have one ball in them and $n - k$ buckets are empty. This is the same as choosing k out of the n buckets, for which there are $n \times (n - 1) \times (n - 2) \times \dots \times (n - k + 1)$ choices for the k balls, and correcting for the internal orderings of the k balls by dividing by $k!$ because the balls are identical.

Subsets k element subsets of a set of size n . The subsets are formed by either choosing the k elements that are included or the $n - k$ elements that are excluded.

Set Binary Labels The problem can also be thought of as assigning a binary label to n elements, where there are k times 1 and $n - k$ times 0.

1.3.8.2 Indistinct Recipients

$$|S| = 1 \text{ if } k \leq n, 0 \text{ otherwise} \quad (1.48)$$

At Most One Ping Pong Ball into Unlabeled Bags k identical ping pong balls are put into n identical bags. The result is either k bags with ping pong balls and $n - k$ empty bags, or there is no possible result if there are not enough bags (i.e. if $k > n$).

Boxes How many ways are there to put k marbles in n boxes, if each box is only big enough for one marble. One, if there are enough boxes, or zero, if there aren't enough boxes.

1.3.9 Identical Objects, Each Receives At Least One

1.3.9.1 Distinct Recipients

This problem is the same as giving each of n recipients one of k objects, and then distributing the remaining $k - n$ objects arbitrarily.

$$|S| = \binom{k + n - 1 - k}{k - n} = \binom{n - 1}{k - 1} \quad (1.49)$$

This can be derived by picturing first distributing one object to each of the n recipients, ensuring that each recipient has at least one, and then distributing the remaining $k - n$ objects arbitrarily. There is only one way to give each recipient one of the identical objects, and then there are $\binom{k - n + n - 1}{k - n}$ ways to distribute the remaining $k - n$ objects arbitrarily across the recipients.

At Least One Ping Pong Ball into Labeled Buckets k ping pong balls are distributed into n labeled buckets so that there is at least one ping pong ball in each bucket. The labeled buckets may have more than one ping pong ball in them.

Compositions in n Parts How many ways are there to assign k identical objects to n labeled sets of at least one object?

1.3.9.2 Indistinct Recipients

$$|S| = P(k, n) \quad (1.50)$$

It turns out that there is no known formula for $P(k, n)$.

At Least One Ping Pong Ball in Unlabeled Bags k identical ping pong balls are distributed across n unlabeled bags, so that none of the bags are empty.

Partitions in n Parts How many ways are there to make n piles from k objects.

1.3.10 Identical Objects, Each Receives Exactly One

1.3.10.1 Distinct Recipients

$$|S| = 1 \text{ if } k = n, 0 \text{ otherwise} \quad (1.51)$$

One Ping Pong Ball into Each Labeled Bucket k identical ping pong balls are distributed across n labeled buckets so that each bucket has one ping pong ball in it. If $k = n$, then there is one way to put one ball in each bucket. If the numbers do not match up, then there is no way to match them one-for-one.

1.3.10.2 Indistinct Recipients

$$|S| = 1 \text{ if } k = n, 0 \text{ otherwise} \quad (1.52)$$

One Ping Pong Ball into Each Unlabeled Bag k identical ping pong balls are distributed across n identical bags so that each bag has one ping pong ball in it. If $k = n$, then there is one way to put one ball in each bag. If the numbers do not match up, then there is no way to match them one-for-one.

1.4 Generating Functions

Generating functions are functions that encode sequences of numbers as the coefficients of power series. One example are the moment generating functions in probability theory, though they are generally extremely useful in combinatorics problems and, almost equivalently, discrete probability problems.

Bogart (2004) approaches the concept in terms of *Picture Functions*. For each element $s \in S$, there is a picture function $P(s)$, so that, for example, the multiset $\{1, 1, 2\}$ can be written as $P(1)^1 P(2)$. Collections of combinations can be rewritten in terms of sums and products, which enables factorization and overall easier accounting. Combinations with particular properties can be filtered by looking at exponents.

The picture function enables writing down combinations of elements as an enumerating function $E_P(s)$. For example, the enumerating function for all multisets that include either one or two times some element a and between zero and two times some element b is written:

$$E_P(s) = \frac{P(a) + P(a)^2 + P(a)P(b) + P(a)^2P(b) + P(a)P(b)^2 + P(a)^2P(b)^2}{(P(a) + P(a)^2)(1 + P(b) + P(b)^2)} \quad (1.53)$$

Generating functions get much more complicated. ? is a good resource.

1.4.1 Example: Binomial Coefficients

Consider the case of a collection of n indistinguishable objects s , and write $P(s) = x$. Then the enumerator for selecting any subset of those n objects is given by:

$$E_P(s) = \prod_i^n (x^0 + x^1) = (1 + x)^n = \sum_i^n \binom{n}{i} x^i \quad (1.54)$$

Where each term $(x^0 + x^1)$ corresponds to the two options of either excluding a particular element (x^0) or including a particular element (x^1). The exponent of the expanded product encodes how many objects were included into a particular subset. This is one way of "proving" the binomial coefficients, and one can say $(1+x)^n$ is the generating function for the binomial coefficients $\binom{n}{i}$.

1.4.2 Example: Basket of Goods

An apple costs $20c$, a pear costs $25c$ and a banana costs $30c$. How many different fruit baskets can be bought for $100c$?

By replacing the picture function $P(s)$ with x , it was possible to identify subsets of n objects by looking at the exponent of x^n in the enumerating function. In this case, the exponent is supposed to show the price. This can be done by writing $P(\text{apple}) = x^{20}$, $P(\text{pear}) = x^{25}$ and $P(\text{banana}) = x^{30}$.

$$E_P(s) = \left(\sum_{i=0}^5 x^{20} \right) \left(\sum_{i=0}^4 x^{25} \right) \left(\sum_{i=0}^3 x^{30} \right) \quad (1.55)$$

Which results in some power series of the form:

$$E_P(s) = 1x^0 + 1x^{20} + 1x^{25} + 1x^{30} + 1x^{40} + \dots + 2x^{60} + \dots + 1x^{290} \quad (1.56)$$

To obtain the number of combinations that correspond to a cost of exactly $100c$, one can apply the operator $\frac{1}{n!} \frac{d^n}{dx^n}$ and set $x = 0$ to obtain the desired term. This is what is done with moment generating functions in statistics.

But actually it is easier to think through what the coefficients will be so that:

$$\sum_{l=0}^{n+m+h} d_l x^l = \left(\sum_{i=0}^n a_i x^i \right) \left(\sum_{j=0}^m b_j x^j \right) \left(\sum_{k=0}^h c_k x^k \right) \quad (1.57)$$

$$d_l = \sum_{\substack{i, j, k \\ i+j+k=l}} a_i b_j c_k \quad (1.58)$$

If $a_i = b_j = c_k = 1$, then:

$$d_l = \sum_{\substack{i, j, k \\ i+j+k=l}} 1 = \binom{l+3-1}{l} \quad (1.59)$$

Which is the number of all multisets of size l and 3 classes.

In this case, however, the coefficients are equal to 1 only for $i = a20$, $j = b25$, and $k = c30$ for $a, b, c \in \mathbb{N}_0$, and zero otherwise, so that there are 4 combinations of a, b, c for which $20a + 25b + 30c = 100$. Hence, $d_{100} = 4$ for the basket of goods above.

The sum can also be rewritten:

$$d_l = \sum_{i=0}^l \sum_{j=0}^{l-i} a_i b_j c_{l-i-j} \quad (1.60)$$

Of course, this is the discrete version of a convolution. So, it is no surprise that the addition of random variables winds up being a convolution.

1.4.3 Example: Dice

How many ways are there for n dice with k faces to show s eyes?

$$E_P = \left(\sum_{i=0}^{\infty} a_i x^i \right)^n = \sum_{i=0}^{\infty} d_i x^i \quad (1.61)$$

Where $a_i = 1 \forall i \in (1, k)$ and $a_i = 0$ otherwise.

$$E_P = \left(\sum_{i=1}^k x^i \right)^n = \left(x(1-x^k) \sum_{i=0}^{\infty} x^i \right)^n = x^n \left(\frac{1-x^k}{1-x} \right)^n = x^n \left(\sum_{i=0}^n (-1)^i \binom{n}{i} x^{ik} \right) \left(\sum_{j=0}^{\infty} (-1)^j \binom{-n}{j} x^j \right) \quad (1.62)$$

The coefficient for x^s is given the sum:

$$d_s = \sum_{ki+j=s-n} (-1)^{i+j} \binom{n}{i} \binom{-n}{j} \quad i \in [0, n], j \in [0, \infty] \quad (1.63)$$

Where the indices i and j satisfy $ki + j = s - n$. For example, for $s = 7$, $n = 2$ and $k = 6$:

$$6i + j = 7 - 2 = 5s \quad (1.64)$$

Holds for $i = 0, j = 5$:

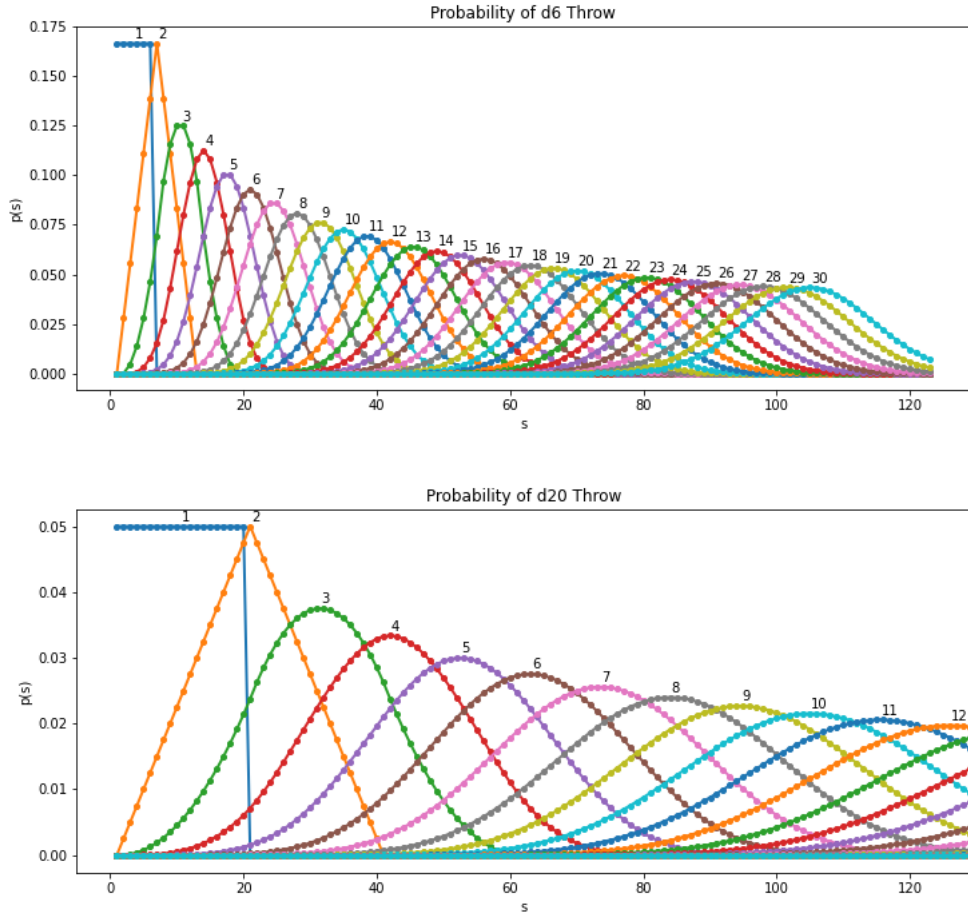


Figure 1.2: Probability of the sum of 6-sided and 20-sided dice, calculated with Eqn. 1.65. Note that these are simply the convolutions of $n = 1, 2, 3, \dots$ square waves, which rapidly adopts the shape of a Bell curve.

$$d_s = (-1)^5 \binom{2}{0} \binom{-2}{5} = (-1)^{10} 1 \binom{2+5-1}{5} = \binom{6}{5} = \frac{6!}{5!1!} = 6 \quad (1.65)$$

Indeed, there are 6 ways for two d6 to add to 7:

$$[6, 1], [5, 2], [4, 3], [3, 4], [2, 5], [1, 6] \quad (1.66)$$

2 Linear Algebra and Multivariable Calculus

Contents of this chapter

2.1	Multi-Index Notation	23	2.8	Types of Matrices	30
2.1.1	Example: Multinomial Coefficients	24	2.8.1	$\text{sgn}(\mathbf{x}^\dagger \mathbf{H} \mathbf{x})$ Definite	30
2.1.2	Example: Taylor Expansion	24	2.8.2	Triangular	30
2.2	Linear Systems of Equations	24	2.8.3	$\mathbf{AB} - \mathbf{BA} = 0$ Commuting	30
2.2.1	$A \in \mathbb{R}^{n \times n}$ Square Matrices	24	2.8.4	$\mathbf{AB} + \mathbf{BA} = 0$ Anticommuting	31
2.2.2	$A \in \mathbb{R}^{m \times n}$ Rectangular Matrices, Overdetermined Case	24	2.8.5	$\mathbf{A}^\dagger = \mathbf{A}$ Hermitian, Symmetric	31
2.2.3	$A \in \mathbb{R}^{n \times m}$ Rectangular Matrices, Underdetermined Case	25	2.8.6	$\mathbf{A}^\dagger = -\mathbf{A}$ Skew Hermitian, Skew Symmetric	31
2.3	$\mathbf{A} = \mathbf{LL}^\dagger$ Cholesky Decomposition	25	2.8.7	$\mathbf{AA} = \mathbf{I}$ Involutory	31
2.3.1	Creating Correlated Random Variables	25	2.8.8	$\mathbf{U}^\dagger = \mathbf{U}^{-1}$ Unitary, Orthogonal	32
2.4	Generalized Eigenvectors	25	2.8.9	$\mathbf{A} = \mathbf{TBT}^{-1}$ Similarity	32
2.5	$\mathbf{A} = \mathbf{VAV}^{-1}$ Spectral Theorems, Diagonalization	25	2.9	Properties of Norms	32
2.5.1	$\mathbf{A} = \mathbf{VAV}^T$ Eigendecomposition of Symmetric Matrices	25	2.10	L^p Lebesgue Vector Norms	32
2.5.2	$\mathbf{H} = \mathbf{U\Lambda U}^T$ Eigendecomposition of Hermitian Matrices	27	2.10.1	L^1 Taxicab / Manhattan Norm	32
2.5.3	Eigenvalue Sensitivity and Accuracy	27	2.10.2	L^2 Euclidian Norm	32
2.6	$\mathbf{A} = \mathbf{U\Sigma V}^\dagger$ Singular Value Decomposition	27	2.10.3	L^∞ Maximum Norm	32
2.6.1	Full and Economy SVDs	28	2.10.4	$L^{-\infty}$ Minimum Norm	32
2.6.2	Matrix Approximation	28	2.11	Operator and Matrix Norms	33
2.7	Types of Transformations	29	2.11.1	$\ \mathbf{A}\ _{(\alpha)}$ Operator Norm	33
2.7.1	Similarity Transformations	29	2.11.2	$\ \mathbf{A}\ _q$ q -Norms	33
2.7.2	Affine Transformations	29	2.11.3	$\ \mathbf{A}\ _F$ Frobenius Norm	34
2.7.3	Unitary Transformations	29	2.11.4	$\ \mathbf{A}\ _{(1)}$ (1)-Norm	34
2.7.4	Multilinear Maps	30	2.11.5	$\ \mathbf{A}\ _{(\infty)}$ (∞)-Norm	34
2.7.5	Multilinear Forms	30	2.11.6	$\ \mathbf{A}\ _{(2)}$ (2)-Norm	34
			2.12	Vector and Matrix Derivatives	34
			2.12.1	Jacobian	34
			2.12.2	Inverse Function Theorem	35
			2.12.3	Critical Points	35
			2.12.4	Differential Volume Element, Change of Variables	35
			2.12.5	Hessian	36

2.1 Multi-Index Notation

Multi-index notation makes high-dimensional things faster and easier. A collection of indices is represented by a tuple $\alpha = (\alpha_1, \alpha_2, \alpha_3, \dots)$. The absolute value $|\alpha| = \sum_i \alpha_i$, partial derivatives $\partial^\alpha = \prod \partial^{\alpha_i}$, powers $\mathbf{x}^\alpha = \prod_i x_i^{\alpha_i}$.

2.1.1 Example: Multinomial Coefficients

Instead of:

$$\sum_{0 \leq i_1, i_2, i_3, \dots, i_k \leq n} \binom{n}{i_1, i_2, i_3, \dots, i_k} \quad (2.1)$$

Write:

$$\sum_{0 \leq |\alpha| \leq n} \binom{n}{\alpha} \quad (2.2)$$

2.1.2 Example: Taylor Expansion

For a vector valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ that is analytical in a neighborhood of the point \mathbf{a} :

$$f(\mathbf{x}) = \sum_{|\alpha| \geq 0} \frac{(\mathbf{x} - \mathbf{a})^\alpha}{\alpha!} (\partial^\alpha) f \quad (2.3)$$

2.2 Linear Systems of Equations

(Real numbers only this time.)

Linear equations are of the form $Ax = b$ where A is a matrix and x and b are vectors. The rows of A and b form a system of equations that must be simultaneously satisfied by the entries of x . If $x, b \in \mathbb{R}^n$, then the solutions to the equation of a single row corresponds to an $n - 1$ -dimensional hyperplane. If the rows of A are linearly independent, then solutions that simultaneously satisfy the equations in k -rows correspond to the $n - k$ -dimensional intersection of k n -dimensional hyperplanes. The solution of x that satisfies all n equations is a $n - n = 0$ -dimensional point, and so x is uniquely determined. If any two rows of A are not linearly independent, then the hyperplanes that correspond to values of x that satisfy them overlap exactly, and their intersection is n dimensional, rather than $n - 1$ dimensional. In this case, the value of x that satisfies all rows of A is not narrowed down to a single point. The system of equations is said to be *underdetermined*. This is equivalently the case when A has $m < n$ rows.

2.2.1 $A \in \mathbb{R}^{n \times n}$ Square Matrices

If $A \in \mathbb{R}^{n \times n}$, then the solution to the system is formally $x = A^{-1}b$, where A^{-1} is the matrix inverse, satisfying $A^{-1}A = I$, where I is the identity matrix.

Since $Ax = b$ is the same as expressing b in terms of a linear combination of the columns of A , the entries of x can be interpreted as the coefficients resulting from the projection of b into the column space of A . Therefore, for an orthonormal matrix, the A^{-1} is simply A^T .

2.2.2 $A \in \mathbb{R}^{m \times n}$ Rectangular Matrices, Overdetermined Case

If $A \in \mathbb{R}^{m \times n}$ with $m > n$ rows, then there need not be any point $x \in \mathbb{R}^n$ in which the m hyperplanes all intersect. In that case, the system does not have a solution $x \in \mathbb{R}^n$, and the system is considered *overdetermined*. (The intersection of m distinct hyperplanes in n dimensional space would have negative dimension $(n - m) < 0$ if $m > n$, which my feeble brain can't make sense of.)

In the overdetermined case $A^{m \times n}$ with $m > n$, the columns of A do not span \mathbb{R}^m and therefore $b \in \mathbb{R}^m$ may have some component ϵ that lies outside of the column space of A . In that case, no linear combination x of the columns of A can express b perfectly, but we might look for approximate solutions \hat{x} so that:

$$A\hat{x} + \epsilon = b \quad (2.4)$$

So that the error $\|\epsilon\|_\alpha$ is minimized. This is the starting point for linear regression from the linear algebra perspective. In practice, the approximation is usually approximated by applying an iterative gradient descent algorithm to minimize the *loss function* $\|\epsilon\|_\alpha$. The choice of metric $\|\cdot\|_\alpha$ is essentially a design choice. For $\alpha = 2$, the metric is the L^2 norm (cf. section ??) and an analytic solution exists, named the *normal equations*. The solution minimizes the least squares error and corresponds to the projection of b into the column space of a . The procedure is better known as ordinary least squares regression and therefore I will move a more elaborate discussion to chapter 7.

2.2.3 $A \in \mathbb{R}^{n \times m}$ Rectangular Matrices, Underdetermined Case

For an underdetermined system with $m < n$, there is either no solution (if the m hyperplanes don't intersect), or there are infinitely many possible solutions that lie on an $n - m$ dimensional hyperplane. One idea is to pick the solution that minimizes $\|\hat{x}\|_2$ based on the idea that it might generalize better. The least norm solution is $\hat{x} = A^T (A A^T)^{-1} b$, which is the projection of $\vec{0}$ on the solution set.

2.3 $A = LL^\dagger$ Cholesky Decomposition

The Cholesky Decomposition exists when a matrix is hermitian and positive-definite. It expresses the matrix A as:

$$A = LL^\dagger \quad (2.5)$$

Where L is a lower-triangular matrix with positive, real diagonal entries. When A is real, then so is L . The Cholesky decomposition enables fast solution of a linear system, but it can also be used to create correlated random variables in Monte Carlo simulations.

2.3.1 Creating Correlated Random Variables

Let u_t be a vector of uncorrelated samples with mean 0 and standard deviation 1. If the covariance matrix of the system to be simulated is Σ with Cholesky decomposition $\Sigma = LL^\dagger$, then the vector $v_t = Lu_t$ has the desired covariance.

2.4 Generalized Eigenvectors

2.5 $A = V\Lambda V^{-1}$ Spectral Theorems, Diagonalization

Spectral theorems deal with diagonalizable linear operators.

A diagonalization of a matrix A is always possible when a matrix is square, and refers to a decomposition of the matrix into the matrix of eigenvectors V and eigenvalues Λ as

$$A = V\Lambda V^{-1} \quad (2.6)$$

2.5.1 $A = V\Lambda V^T$ Eigendecomposition of Symmetric Matrices

A symmetric matrix A has orthogonal eigenvectors so that $V^{-1} = V^T$ and real eigenvalues. In that case, the diagonalization is:

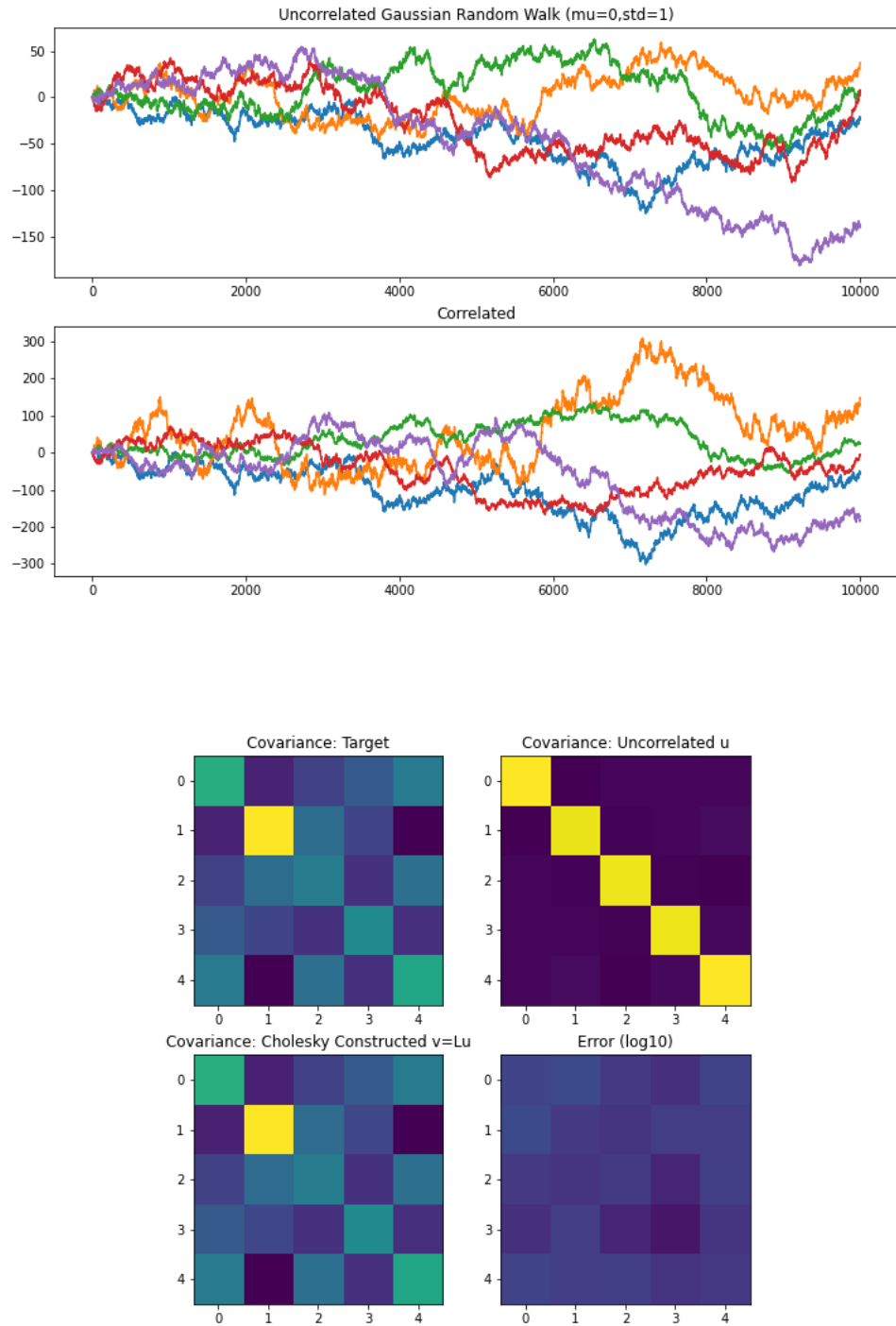


Figure 2.1: Creating correlated random variables from uncorrelated random variables using the Cholesky decomposition of the covariance matrix. The 5 uncorrelated random variables are sampled from a standard normal distribution. It is difficult to see a difference between the correlated and uncorrelated random walks.

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \begin{bmatrix} | & | & \cdots & | \\ v_1 & v_2 & \ddots & v_n \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \lambda_n \end{bmatrix} \begin{bmatrix} \text{---} & v_1 & \text{---} \\ \text{---} & v_2 & \text{---} \\ \vdots & \ddots & \vdots \\ \text{---} & v_n & \text{---} \end{bmatrix} \quad (2.7)$$

This means that enables the expression of \mathbf{A} in terms of projections on the eigenvectors:

$$\mathbf{A} = \sum_i \lambda_i (v_i \otimes v_i) \quad (2.8)$$

Where \otimes is the outer product. Since the eigenvectors are an orthonormal basis, $\sum_i v_i \otimes v_i = \mathbb{I}$.

2.5.2 $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ Eigendecomposition of Hermitian Matrices

Similarly, a hermitian matrix $\mathbf{H} \in \mathbb{C}^{n \times n}$ (the complex equivalent to a symmetric matrix) has real eigenvalues and the matrix of eigenvectors is unitary, so that $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$.

2.5.3 Eigenvalue Sensitivity and Accuracy

2.5.3.1 General Case

In general, the values of the eigenvalues of a square matrix \mathbf{A} may vary wildly under a slight change of $\mathbf{A} \rightarrow \mathbf{A} + \delta\mathbf{A}$. The sensitivity of the eigenvalues to a change in \mathbf{A} can be investigated using matrix norms (?). Let $\|\cdot\|$ denote a submultiplicative matrix norm, then:

$$\begin{aligned} \mathbf{\Lambda} + \delta\mathbf{\Lambda} &= \mathbf{X}^{-1} (\mathbf{A} + \delta\mathbf{A}) \mathbf{X} \\ \delta\mathbf{\Lambda} &= \mathbf{X}^{-1} \delta\mathbf{A} \mathbf{X} \\ \|\delta\mathbf{\Lambda}\| &= \|\mathbf{X}^{-1} \delta\mathbf{A} \mathbf{X}\| \leq \|\mathbf{X}^{-1}\| \|\mathbf{X}\| \|\delta\mathbf{A}\| \end{aligned} \quad (2.9)$$

When $\|\cdot\|$ is chosen to be the operator norm with respect to L^2 , $\|\cdot\|_{(2)}$, then $\|\mathbf{X}^{-1}\| = \sigma_1$ and $\|\mathbf{X}\| = \frac{1}{\sigma_n}$ where σ_1 and σ_2 are the square roots of the largest and the smallest eigenvalue of $\mathbf{X}^\dagger \mathbf{X}$ respectively (cf. section 2.11 on matrix norms). In that case, the sensitivity of the eigenvalues to a change in \mathbf{A} is:

$$\|\delta\mathbf{\Lambda}\|_{(2)} \leq \frac{\sigma_1}{\sigma_n} \|\delta\mathbf{A}\|_{(2)} = \kappa(\mathbf{X}) \|\delta\mathbf{A}\|_{(2)} \quad (2.10)$$

Where $\kappa(\mathbf{X})$ is the conditioning number of the matrix \mathbf{X} . Upper bounds on the error on individual eigenvalues can also be derived quite easily, which is shown in ?? pp 10-12.

2.5.3.2 Hermitian Matrices

For hermitian (or orthogonal) matrices, the conditioning number for the individual eigenvalues $\kappa(\lambda_i, \mathbf{H}) = 1$, so that the error on an individual eigenvalue $\|\lambda_i\|_{(2)} \leq \kappa(\lambda_i, \mathbf{H}) \|\mathbf{H}\|_{(2)} = 1 \times \|\mathbf{H}\|_{(2)}$.

2.6 $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\dagger$ Singular Value Decomposition

The Singular Value Decomposition (SVD) exists for any matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$, and is a closely related alternative to the eigendecomposition (cf. section 2.5) that works for non-square matrices. The contrast to diagonalization is that the eigenvectors and eigenvalues of $\mathbf{A}^\dagger \mathbf{A}$ and $\mathbf{A} \mathbf{A}^\dagger$ found, rather than to look for the eigensystem of \mathbf{A} itself. The advantage is that $\mathbf{A}^\dagger \mathbf{A}$ and $\mathbf{A} \mathbf{A}^\dagger$ always have the convenient properties of being square, hermitian and positive semidefinite.

SVD comes up incessantly in the context of data analysis. In general, the decomposition has the form:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\dagger \quad (2.11)$$

Where $\mathbf{U}^\dagger\mathbf{U} = \mathbf{I}$, \mathbf{V} is unitary, and $\mathbf{\Sigma}$ is a diagonal matrix with real and positive entries σ_i^2 along the diagonal, so that $\sigma_1 \geq \sigma_2^2 \geq \dots \geq \sigma_n^2$. \mathbf{U} is the matrix of left singular vectors, which are the eigenvectors of $\mathbf{A}\mathbf{A}^\dagger$. \mathbf{V} is the matrix of right singular vectors, which are the eigenvectors of $\mathbf{A}^\dagger\mathbf{A}$. The singular values are the square roots of the eigenvalues of $\mathbf{A}^\dagger\mathbf{A}$ or, equivalently, $\mathbf{A}\mathbf{A}^\dagger$. If the \mathbf{A} happens to be square and symmetric ($\mathbf{A} = \mathbf{A}^\mathbf{T}$, cf. section 2.8.5), then the singular values are simply the absolute values of the eigenvalues of \mathbf{A} . The right and left singular vectors of the matrix are orthonormal bases that respectively span the column and row space of \mathbf{A} . The number of singular values of \mathbf{A} is the rank of \mathbf{A} .

The singular values are always real and positive or zero, because $\mathbf{A}^\dagger\mathbf{A}$ and $\mathbf{A}\mathbf{A}^\dagger$ are hermitian and positive semidefinite (cf. sections 2.8.5, 2.8.1). When the matrix \mathbf{A} has purely real entries, then the left and right singular vectors are also real, which is not necessarily the case for complex \mathbf{A} .

The SVD of a matrix is unique up to the sign columns in \mathbf{U} and \mathbf{V} . That is, the SVD is valid under transforming $\mathbf{u}_i, \mathbf{v}_i \rightarrow -\mathbf{u}_i, -\mathbf{v}_i$ if \mathbf{u}_i and \mathbf{v}_i are the i th vectors in \mathbf{U} and \mathbf{V} . This has the consequence that the basis of singular vectors that is found using SVD does not remain consistently oriented in the presence of noise or the slow evolution of a system. This is addressed in greater length in section 9.3 on performing SVD specifically on data matrices.

2.6.1 Full and Economy SVDs

While these properties are always true, unfortunately people use a range of conventions when it comes to the size of \mathbf{U} , $\mathbf{\Sigma}$ and \mathbf{V} .

The first convention is for \mathbf{U} and \mathbf{V} to be square, in which case they contain the full set of left and right singular vectors, and $\mathbf{\Sigma}$ has dimension $m \times n$, with 0 entries in rows $i > n$. This is known as the *full SVD*.

Friedman et al. (2001) uses the convention where $\mathbf{\Sigma}$ is square, i.e. $\mathbf{U} : m \times n$, $\mathbf{\Sigma} : n \times n$ and $\mathbf{V} : n \times n$. This means that \mathbf{U} does not contain the full set of m left singular vectors. Note that, in this case, $\mathbf{U}^\dagger\mathbf{U} = \mathbf{I}$ but $\mathbf{U}\mathbf{U}^\dagger \neq \mathbf{I}$. Friedman et al.'s (2001) convention is of advantage in the context of data analysis, where the data matrix tends to be "tall and skinny" (i.e. $m \gg n$), and only the first n left singular vectors are relevant. Also, letting $\mathbf{\Sigma}$ be a square matrix significantly simplifies calculations. Friedman et al.'s (2001) is known as the *economy SVD*.

The two different layouts are illustrated in Figure 2.2, which I brazenly copied from Mathworks (n/a).

2.6.2 Matrix Approximation

The SVD enables the decomposition of the matrix into a sum:

$$\mathbf{A} = \sum_{i=1}^n \sigma_i (\mathbf{u}_i \otimes \mathbf{v}_i) \quad (2.12)$$

Where \mathbf{u}_i is the i th left singular vector and \mathbf{v}_i is the i th right singular vector. Taking the first $r \leq n$ terms of this series is the best rank- r approximation to \mathbf{A} under the Frobenius norm of the error, i.e. $\arg \min_{\text{rank}(\mathbf{A}^{(r)})=r \leq n} \|\mathbf{A}^{(r)} - \mathbf{A}\|_F$. (The Frobenius norm essentially measures the elementwise least-squares error, cf. section 2.11.3). In the context of data analysis, the reduced rank representation of a data matrix \mathbf{A} acts as a filter in which the components that explain less of the variance in the dataset are removed (the hope is that these correspond to noise). For a more elaborate discussion of SVD in data analysis, see section 9.3

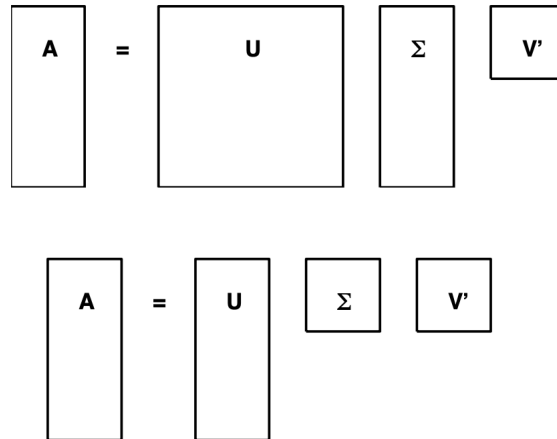


Figure 2.2: Dimensions for Full vs. Economy SVDs

2.7 Types of Transformations

2.7.1 Similarity Transformations

If T is a nonsingular matrix, then a similarity transformation is defined as:

$$A = TBT^{-1} \quad (2.13)$$

And A and B are said to be *similar*.

2.7.2 Affine Transformations

Affine transformations are the combination of a linear map and a translation, which has the form $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$.

$$f : V \rightarrow W \quad (2.14)$$

Where V and W are vector spaces. Affine transformations can be expressed as matrices by adding an entry with a constant to the vectors that describe a point in space. For example, for $\mathbf{x} \in \mathbb{R}^n$, the affine transform $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ with $A \in \mathbb{R}^{n,n}$ and $\mathbf{b} \in \mathbb{R}^n$ can be expressed as the product of a rectangular matrix M and a vector \mathbf{c} as:

$$A\mathbf{x} + \mathbf{b} = \underbrace{\begin{bmatrix} A & \mathbf{b} \end{bmatrix}}_M \underbrace{\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}}_{\mathbf{c}} \quad (2.15)$$

Where $\mathbf{c}^T = [x_1, x_2, x_3, \dots, x_n, 1]$ and $M \in \mathbb{R}^{n,n+1}$.

2.7.3 Unitary Transformations

Unitary transformations are transformations that preserve the inner product, i.e. $\hat{U}x \cdot \hat{U}y = x \cdot y$. As linear transformations, they are represented by unitary matrices (cf. section ??). Unitary transformations include translations, reflections and rotations.

2.7.4 Multilinear Maps

A multilinear map acts on several vectors in a way that is linear in each of its arguments. A k -linear map acts on k vectors, where $k = 2$ are bilinear maps and $k = 1$ are linear maps.

$$f : V_1 \times V_2 \times \dots \times V_n \rightarrow W \quad (2.16)$$

Where V_1, V_2, \dots, V_n and W are vector spaces. An example would be the addition or subtraction of two or more vectors.

2.7.5 Multilinear Forms

Multilinear forms are multilinear maps that have a scalar output. An example is the dot product between two vectors, or summing over the elements of one or more vectors.

$$f : V_1 \times V_2 \times \dots \times V_n \rightarrow K \quad (2.17)$$

Where V_1, V_2, \dots, V_n and K is a scalar field.

2.8 Types of Matrices

2.8.1 $\text{sgn}(\mathbf{x}^\dagger \mathbf{H} \mathbf{x})$ Definite

A hermitian matrix $\mathbf{H} \in \mathbb{C}^n$ is positive definite, if for any non-zero column vector $\mathbf{x} \in \mathbb{C}^n$, the quadratic form $\mathbf{x}^\dagger \mathbf{H} \mathbf{x} > 0$, and negative definite if $\mathbf{x}^\dagger \mathbf{H} \mathbf{x} < 0$. The matrix is positive or negative *semidefinite* if $\mathbf{x}^\dagger \mathbf{H} \mathbf{x} \geq 0$ or $\mathbf{x}^\dagger \mathbf{H} \mathbf{x} \leq 0$, respectively. Definiteness plays a role in investigating the convexity of a function by looking at the Hessian (section 2.12.5). Sometimes notation with curly comparison symbols are used. The relationship between definiteness and eigenvalues is intuitive:

$\mathbf{H} \succ 0$	positive definite	all eigenvalues are positive
$\mathbf{H} \prec 0$	negative definite	all eigenvalues are negative
$\mathbf{H} \succeq 0$	positive semidefinite	all eigenvalues are positive or 0
$\mathbf{H} \preceq 0$	negative semidefinite	all eigenvalues are negative or 0

The curly comparison symbols can mean other stuff though, for example in the context of partially ordered sets (order theory) or comparisons between multidimensional arrays.

2.8.2 Triangular

A lower triangular matrix is a matrix that has all-zero entries above the diagonal.

$$\mathbf{L} = \begin{bmatrix} l_{1,1} & & & & & 0 \\ l_{2,1} & l_{2,2} & & & & \\ l_{3,1} & l_{3,2} & \ddots & & & \\ \vdots & \vdots & \ddots & \ddots & & \\ \vdots & \vdots & & \ddots & \ddots & \\ l_{n,1} & l_{n,2} & \dots & \dots & l_{n,n-1} & l_{n,n} \end{bmatrix} \quad (2.18)$$

Upper triangular matrices are matrices that have all-zero entries below the diagonal.

2.8.3 $\mathbf{AB} - \mathbf{BA} = 0$ Commuting

Two matrices commute if $\mathbf{AB} = \mathbf{BA}$, or, equivalently, their *commutator* $[\mathbf{A}, \mathbf{B}] = \mathbf{AB} - \mathbf{BA}$ is zero. This means that \mathbf{A} and \mathbf{B} both have to be square. Matrices commute when they have the same eigenspace, i.e.

they have the same eigenvectors. This can be seen by considering the diagonal representations of \mathbf{A} and \mathbf{B} .

Let \mathbf{A} and \mathbf{B} be two square matrices with the same eigenvectors \mathbf{V} , then they can be diagonalized as:

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}_\mathbf{A}\mathbf{V}^{-1}\mathbf{B} = \mathbf{V}\mathbf{\Lambda}_\mathbf{B}\mathbf{V}^{-1} \quad (2.19)$$

They commute, because:

$$\mathbf{AB} = \mathbf{V}\mathbf{\Lambda}_\mathbf{A}\mathbf{V}^{-1}\mathbf{V}\mathbf{\Lambda}_\mathbf{B}\mathbf{V}^{-1} = \mathbf{V}\mathbf{\Lambda}_\mathbf{B}\mathbf{\Lambda}_\mathbf{A}\mathbf{V}^{-1} = \mathbf{V}\mathbf{\Lambda}_\mathbf{B}\mathbf{V}^{-1}\mathbf{V}\mathbf{\Lambda}_\mathbf{A}\mathbf{V}^{-1} = \mathbf{BA} \quad (2.20)$$

2.8.4 $\mathbf{AB} + \mathbf{BA} = 0$ Anticommuting

Two matrices anticommute if $\mathbf{AB} = -\mathbf{BA}$, or, equivalently, their *anticommutator* $\{\mathbf{A}, \mathbf{B}\} = \mathbf{AB} + \mathbf{BA}$ is zero.

2.8.5 $\mathbf{A}^\dagger = \mathbf{A}$ Hermitian, Symmetric

Hermitian matrices are matrices that are equal to their complex transpose. That is:

$$\mathbf{A}^{*T} = \mathbf{A}^\dagger = \mathbf{A} \quad (2.21)$$

2.8.5.1 Properties

There are many properties of Hermitian matrices.

- By definition: $\mathbf{A} = \mathbf{A}^\dagger$
- Diagonal Entries are all real, since $a_{i,i} = a_{i,i}^*$, but not necessarily positive (physics will mislead you there...)
- Inverse is also hermitian: $\mathbf{A}^{-1} = \mathbf{A}^{-1\dagger}$
- Diagonalizable with real eigenvalues and orthogonal eigenvectors $\in \mathbb{C}^n$.

Hermitian matrices with only real entries are called symmetric matrices. In that case $\mathbf{A}^T = \mathbf{A}$.
Hermitian matrices can only have real elements along their diagonal.

2.8.6 $\mathbf{A}^\dagger = -\mathbf{A}$ Skew Hermitian, Skew Symmetric

Skew Hermitian matrices that are equal to the negative of their complex transpose. That is:

$$\mathbf{A}^{*T} = \mathbf{A}^\dagger = -\mathbf{A} \quad (2.22)$$

Real matrices that are skew Hermitian are called skew symmetric. In that case:

$$\mathbf{A}^T = -\mathbf{A} \quad (2.23)$$

Skew Hermitian matrices can only have complex values on their diagonal, and skew symmetric matrices can only have zeros as diagonal elements.

2.8.7 $\mathbf{AA} = \mathbf{I}$ Involutory

Involutory matrixes are matrices that are their own inverse, so that:

$$\mathbf{AA} = \mathbf{I} \quad (2.24)$$

Involutory matrices are all square roots of the identity matrix. A famous example are the 2×2 Pauli matrices.

2.8.8 $U^\dagger = U^{-1}$ Unitary, Orthogonal

Unitary matrices satisfy $U^\dagger U = U U^\dagger = \mathbf{I}$, and they have $\det(U) = 1$. They are diagonalizable and can be expressed as $e^{i\mathbf{H}}$ where \mathbf{H} is a Hermitian matrix.

Unitary matrices that are real are called orthogonal. Orthogonal matrices satisfy $\mathbf{A}^{-1} = \mathbf{A}^T$. The rows (and columns) of \mathbf{A} are an orthonormal basis in \mathbb{R}^n .

Unitary matrices are necessarily invertible, and have determinant $|U| = 1$ or $|U| = -1$. They represent unitary transformations, which means that they preserve the inner product between two vectors.

The set of $n \times n$ orthogonal matrices is known as the orthogonal group $O(n)$ and the subgroup of orthogonal matrices with determinant 1 is known as the special orthogonal group $SO(n)$. The elements of $SO(n)$ are rotations, and the elements of $O(n)$ represent translations, reflections or rotations. Similarly, the group of $n \times n$ unitary matrices is the unitary group $U(n)$ and the subgroup of $U(n)$ that has determinant 1 is the special unitary group $SU(n)$.

2.8.9 $\mathbf{A} = \mathbf{T}\mathbf{B}\mathbf{T}^{-1}$ Similarity

Two matrices are said to be similar if they can be related through a similarity transformation $\mathbf{A} = \mathbf{T}\mathbf{B}\mathbf{T}^{-1}$ where \mathbf{T} is some nonsingular matrix (cf. section 2.7.1). An important example is that square matrices are similar to diagonal matrices, see section 2.5.

2.9 Properties of Norms

2.10 L^p Lebesgue Vector Norms

Let $p \geq 1$ be a real number, then the p -norm or L^p norm of a vector $\mathbf{x} \in \mathbb{C}^n$ is defined as:

$$\|\mathbf{x}\|^p = \left[\sum_i^n |x_i|^p \right]^{1/p} \quad (2.25)$$

The expression can still be useful for $0 < p < 1$, but in that case the result is not a proper norm, because it is not subadditive (does not satisfy $f(x+y) \leq f(x) + f(y)$). p -norms are closely related to expressions for the generalized mean.

2.10.1 L^1 Taxicab / Manhattan Norm

$$\|\mathbf{x}\|^1 = \sum_i |x_i| \quad (2.26)$$

2.10.2 L^2 Euclidian Norm

$$\|\mathbf{x}\|^2 = \sum_i |x_i|^2 \quad (2.27)$$

2.10.3 L^∞ Maximum Norm

$$\|\mathbf{x}\|^\infty = \max(x_1, x_2, \dots, x_n) \quad (2.28)$$

2.10.4 $L^{-\infty}$ Minimum Norm

Formally, I only came across values $0 < p$, but it is my opinion that $-\infty$ picks out the minimum value:

$$\|\mathbf{x}\|^{-\infty} = \min(x_1, x_2, \dots, x_n) \quad (2.29)$$

2.11 Operator and Matrix Norms

Matrix norms are functions $\|\cdot\| : K^{m \times n} \rightarrow \mathbb{R}$ where K is a field of real or complex numbers. They satisfy:

- $\|\alpha A\| = |\alpha| \|A\|$ (absolutely homogenous)
- $\|A + B\| \leq \|A\| + \|B\|$ (triangle inequality, subadditivity)
- $\|A\| \geq 0$ (positive valued)
- $\|A\| = 0 \implies A_{n,m} = 0$ (definiteness)

A norm is submultiplicative if it satisfies $\|AB\| \leq \|A\| \|B\|$, which Gera (2009) calls a requirement of "useful matrix norms".

The main risk of confusion is that norms for operators and vectors are different animals. Norms for operators normally measure some relationship between input and output. Norms for vectors are normally some kind of size, length or distance metric. In as far as matrices can be thought of as both operators and multidimensional vectors, norms of either type may be applied to them. People's notation and language is all over the place. Below I've used $\|\cdot\|_{(\alpha)}$ to denote norms in the operator sense and $\|\cdot\|_{\alpha}$ in the vector sense.

2.11.1 $\|A\|_{(\alpha)}$ Operator Norm

The operator norm describes the largest change in size that it may impart on any of its inputs. That means that the operator norm is defined with respect to a definition of size in both domain and codomain. I.e., for an operator \mathbf{A} and a given way of measuring size $\|\cdot\|_{\alpha}$:

$$\|\mathbf{A}\|_{(\alpha)} = \sup \left\{ \frac{\|\mathbf{A}\mathbf{v}\|_{\alpha}}{\|\mathbf{x}\|_{\alpha}} : \mathbf{v} \in V \right\} \quad (2.30)$$

When the operator is given by a matrix \mathbf{A} , and the length of the vector \mathbf{x} is measured using the usual euclidian 2-norm ($\|\cdot\|_2$), then the operator norm is given by the square root of the largest eigenvalue of $\mathbf{A}^T \mathbf{A}$. In that case, the operator norm is the same as the 2-norm (cf. section 2.11.6).

Note that $\|A\|_{(\alpha)}$ and $\|A\|_{\alpha}$ are two different things. The former measures the change in input size, where the size of the input is measured according to the latter. That is the reason for why the 1-Norm and 2-Norms are so different from the vector norms L_1 and L_2 .

2.11.2 $\|A\|_q$ q -Norms

The q norms for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with entries $a_{i,j}$ in row i and column j are defined:

$$\|\mathbf{A}\|_q = \left(\sum_i \sum_j a_{i,j}^q \right)^{1/q} \quad (2.31)$$

For $q = 2$, this becomes the Frobenius norm (section 2.11.3). For vectors $\mathbf{v} \in \mathbb{R}^n$, the q -norm is more known as p -norm or L^p norm (cf. section 2.10).

2.11.3 $\|\mathbf{A}\|_F$ Frobenius Norm

The Frobenius Norm is the sum of the squares of all entries of a matrix. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a matrix with entries $a_{i,j}$ in row i and column j , then:

$$\|\mathbf{A}\|_F = \sqrt{\sum_i \sum_j a_{i,j}^2} \quad (2.32)$$

The Frobenius norm is invariant under rotations, and $\|\mathbf{A}\|_F = \sqrt{\sum_i \sigma_i^2}$ where σ_i are the singular values of \mathbf{A} .

2.11.4 $\|\mathbf{A}\|_{(1)}$ (1)-Norm

Let \mathbf{A} be a matrix with entries $a_{i,j}$ in row i and column j , then:

$$\|\mathbf{A}\|_{(1)} = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{i,j}| \quad (2.33)$$

That is, it is the maximum of the sums of the absolute values of any of the columns of \mathbf{A} .

2.11.5 $\|\mathbf{A}\|_{(\infty)}$ (∞)-Norm

Let \mathbf{A} be a matrix with entries $a_{i,j}$ in row i and column j , then:

$$\|\mathbf{A}\|_{(\infty)} = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{i,j}| \quad (2.34)$$

That is, it is the maximum of the sums of the absolute values of any of the rows of \mathbf{A} .

2.11.6 $\|\mathbf{A}\|_{(2)}$ (2)-Norm

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a matrix with entries $a_{i,j}$ in row i and column j , then:

$$\|\mathbf{A}\|_{(2)} = \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2} \quad (2.35)$$

Which is the square root of the largest eigenvalue of $\mathbf{A}^T \mathbf{A}$. Or, equivalently, $\|\mathbf{A}\|_{(2)} = \sigma_1$, where σ_1 is the largest singular value of the SVD of $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$. That means that $\|\mathbf{A}^{-1}\| = \frac{1}{\sigma_n}$, where σ_n is the smallest singular value of the SVD of \mathbf{A} .

2.12 Vector and Matrix Derivatives

Derivatives involving matrices and vectors can look nonintuitive when the usual symbolic matrix notation is used, but can be derived handily when index notation is used. A very concise and helpful resource for this is Barnes (n/a).

2.12.1 Jacobian

It is particularly helpful to remember the Jacobian, which is the derivative of a function with respect of a vector. The Jacobian of some function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is:

$$\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \left[\frac{\partial \mathbf{f}}{\partial x_1}, \dots, \frac{\partial \mathbf{f}}{\partial x_n} \right] = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad (2.36)$$

I enjoy writing the gradient $\frac{d}{d\mathbf{x}}$ as $\nabla_{\mathbf{x}}$. The relationships below can all be derived as applications of the Jacobian.

$$\begin{aligned} \nabla_{\mathbf{x}} (\mathbf{u}^T \mathbf{x}) &= \left[\frac{\partial}{\partial x_1} (\sum_i u_i x_i), \dots, \frac{\partial}{\partial x_n} (\sum_i u_i x_i) \right] = \mathbf{u}^T \\ \nabla_{\mathbf{x}} (\mathbf{x}^T \mathbf{u}) &= \left[\frac{\partial}{\partial x_1} (\sum_i u_i x_i), \dots, \frac{\partial}{\partial x_n} (\sum_i u_i x_i) \right] = \mathbf{u}^T \\ \nabla_{\mathbf{x}} (\mathbf{x}^T \mathbf{x}) &= \left[\frac{\partial}{\partial x_1} (\sum_i x_i^2), \dots, \frac{\partial}{\partial x_n} (\sum_i x_i^2) \right] = 2\mathbf{x}^T \\ \nabla_{\mathbf{x}} (\mathbf{A}\mathbf{x}) &= \begin{bmatrix} \frac{\partial}{\partial x_1} \left(\underbrace{\sum_i A_{1i} x_i}_{A_{11}} \right) & \dots & \frac{\partial}{\partial x_n} \left(\underbrace{\sum_i A_{1i} x_i}_{A_{1n}} \right) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} \left(\underbrace{\sum_i A_{ni} x_i}_{A_{n1}} \right) & \dots & \frac{\partial}{\partial x_n} \left(\underbrace{\sum_i A_{ni} x_i}_{A_{nn}} \right) \end{bmatrix} = \mathbf{A} \end{aligned} \quad (2.37)$$

2.12.2 Inverse Function Theorem

The inverse function theorem gives a sufficient condition for the invertibility of a function near some point in its domain. If the derivative f' of a function f is continuous and non-zero near some point a within its domain, then the function is invertible near that point. If $b = f(a)$, then:

$$\frac{d[f^{-1}(b)]}{dx} = \frac{1}{\frac{df(a)}{dx}} \quad (2.38)$$

That is, the derivative of the inverse function at a point $b = f(a)$ of the range, is the reciprocal of the derivative of the function near the point a in the domain. This extends to multivariable calculus. Given a function $\mathbf{f} : \mathbf{x} \rightarrow \mathbf{y}$:

$$\nabla_{\mathbf{y}} [\mathbf{f}^{-1}] = [\nabla_{\mathbf{x}} \mathbf{f}]^{-1} \quad (2.39)$$

In words: the Jacobian of the inverse function at the point $\mathbf{b} = \mathbf{f}(\mathbf{a})$ is the matrix inverse of the Jacobian of the function at the point \mathbf{a} . The sufficient condition is that the Jacobian $\nabla_{\mathbf{x}} \mathbf{f}$ is continuous and *nonsingular* near \mathbf{a} .

2.12.3 Critical Points

Critical points are points where the Jacobian does not have maximal rank. In case of a square Jacobian, this means that the Jacobian is singular.

2.12.4 Differential Volume Element, Change of Variables

The Jacobian is used when transforming between different coordinate systems. Consider a transformation $\mathbf{x} = \mathbf{H}(\mathbf{y})$, then:

$$d^n x = |\nabla_y \mathbf{H}| d^n y \quad (2.40)$$

And:

$$\int_{\mathbf{x}} d^n \mathbf{x} f(\mathbf{x}) = \int_{\mathbf{y}} d^n \mathbf{y} |\nabla_y \mathbf{H}| f(\mathbf{H}(\mathbf{y})) \quad (2.41)$$

Alternatively, if $\mathbf{y} = \mathbf{H}^{-1}(\mathbf{x})$:

$$\begin{aligned} d^n y &= |\nabla_x \mathbf{H}^{-1}(\mathbf{x})| d^n x \\ &= |\nabla_y \mathbf{H}(\mathbf{y})|^{-1} d^n x \\ d^n x &= \frac{1}{|\nabla_y \mathbf{H}(\mathbf{y})|^{-1}} d^n y \end{aligned} \quad (2.42)$$

The Jacobian has to be nonsingular within the domain of integration. This implies that \mathbf{x} and \mathbf{y} have to have the same dimension. In the context of probability theory that sometimes requires artificially defining additional variables so that \mathbf{H} is bijective because the quantity of interest has lower dimension (for example, if you calculate the mean of a random variable).

2.12.5 Hessian

The Hessian is the second derivative of a scalar valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with respect to a vector, i.e. $\nabla \cdot \nabla f$. The elements are $\mathbf{H}(f)_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j}$.

$$\mathbf{H}(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad (2.43)$$

The Hessian of a vector valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a third order tensor with elements $\mathbf{H}(\mathbf{f})_{i,j,k} = \frac{\partial^2 f_k}{\partial x_i \partial x_j}$.

2.12.5.1 Testing Convexity

The definiteness (cf. section ??) of the Hessian is used to test convexity.

- $\mathbf{H} \succeq 0$ convex
- $\mathbf{H} \succ 0$ strictly convex
- $\mathbf{H} \preceq 0$ concave
- $\mathbf{H} \prec 0$ strictly concave

If this holds at a point, the property is local (for example at a local maximum or minimum), and if it holds everywhere on the domain, then the property is global.

3 Probability

Contents of this chapter

3.1	Interpretations and Definitions of Probability	37	3.3.2	Example: Lower Dimensional Random Variable . . .	40
3.1.1	Naive and Non-Naive Definitions of Probability	37	3.4	Indicator Variables	40
3.1.2	Frequentist Probability . . .	38	3.4.1	Example: The Party Problem	41
3.1.3	Bayesian Probability	38	3.5	Copulas	44
3.1.4	Measure Theoretic Probability	38	3.6	Relationships Between Distributions	44
3.2	Measure Theory	38	3.7	Large Deviation Theory	44
3.3	Functions of Random Variables, Derived Distributions	38	3.7.1	Gaertner-Ellis Theorem . . .	44
3.3.1	Example: Sum of Random Variables	39	3.7.2	Example: Sum of Uniform Random Variables	44

3.1 Interpretations and Definitions of Probability

3.1.1 Naive and Non-Naive Definitions of Probability

Probability courses typically start with what ? calls the *naive* definition of probability, which is to look at the fraction of the event space that corresponds to a particular event. This works well for combinatorial type of probability problems, such as the likely outcome of dice throws. The dichotomy below is from ?.

3.1.1.1 Naive Probability

The event space S consists of a collection of equally likely outcomes. The actual outcome $s_{actual} \in S$. The probability of an outcome in a subset $A \subseteq S$, $P(s_{actual} \in A)$ corresponds to the fraction of events in A out of S .

$$P_{naive}(A) = \frac{|A|}{|S|} = \frac{\text{number of outcomes favorable to } A}{\text{total number of outcomes in } S} \quad (3.1)$$

This approach runs into trouble when the outcomes in $|A|$ are not equally likely, or if the sample space is infinite, i.e. $|S| = \infty$.

3.1.1.2 Non-Naive Definition of Probability

A probability space consists of a sample space S in addition to a *probability function*. The job of the probability function is to take an event $A \subseteq S$ and map it to a number between 0 and 1. I.e. $P : S \rightarrow [0, 1]$.

The function P must satisfy:

- $P(\emptyset) = 0$, $P(S) = 1$

- If A_1, A_2, \dots are *disjoint* events, then:

$$P\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} P(A_j) \quad (3.2)$$

In other words, all you need to work with probability is a set of possible outcomes S and some function to apply to parts of that set.

3.1.2 Frequentist Probability

The *frequentist* interpretation of probability is that it represents the frequency of an outcome of a particular experiment upon running the experiment many times. On one hand, this view on probability is arguably empirically grounded. On the other hand it relies on the impractical notion of identically repeating something a large number of times.

3.1.3 Bayesian Probability

The *Bayesian* view on probability is that it represents a degree of belief in a particular outcome. On one hand, this implies a degree of subjectivity. On the other hand, this notion of probability is much more broadly applicable, because it does not require repeatable experiments. It also permits the introduction of subjective biases through priors.

I'm not sure whether the supposed "battle" between frequentists and Bayesians has ever been of any consequence for me. Also, the frequentist perspective strikes me as inherently contradictory, because supposed empirical proof can strictly speaking never be practically delivered.

In practice, it seems that frequentist approaches tend to be structured around parametrization of solutions in terms of summary statistics, while Bayesians have a tendency to approach problems in terms of full probability distributions. Social scientists seem to be leaning towards frequentist methods, while physicists and engineers seem to be more interested in Bayesian data analysis.

I suspect that an additional source of bias is that frequentist methods tend to be computationally lighter and that many standard accepted recipes exist for common problems. In contrast, Bayesian approaches are usually derived ad-hoc in the context of a particular problem.

3.1.4 Measure Theoretic Probability

The modern approach to probability is rooted in measure theory. In that context:

- The *sample space* is a measurable set
- An *event* is a measurable set; a subset of the sample space.
- A *random variable* is a measurable function on the sample space.
- The *expectation* of a random variable is its integral with respect to the probability measure.

3.2 Measure Theory

3.3 Functions of Random Variables, Derived Distributions

(This section needs revision)

I find changes of variables to be the easiest to understand by writing down a joint probability distribution using a delta function to express the conditional probability for the new variables based on the old variables. The next step is to perform a change of variables in the delta function, and integrate. The

delta function means that the relationship between the old and the new variables is deterministic. If the relationship is non-deterministic you can use whatever expression for the conditional probability is appropriate.

Let's say that you want to know the probability distribution of the function $\vec{u} = \mathbf{H}(\vec{x})$ of some random variable \vec{x} for which the probability distribution is known.

I personally find it least confusing to approach this problem by thinking about the joint probability distribution, and then obtaining $p(\vec{u})$ through marginalization:

$$p(\vec{u}) = \int dx^n p(\vec{x}, \vec{u}) = \int dx^n p_x(\vec{x}) p_u(\vec{u}|\vec{x}) \quad (3.3)$$

Since \vec{u} is a deterministic function of \vec{x} , the conditional probability $p_u(\vec{u}|\vec{x}) = \delta(\vec{u} - \mathbf{H}(\vec{x}))$. Explicitly:

$$p(\vec{u}) = \int dx^n p_x(\vec{x}) \delta(\vec{u} - \mathbf{H}(\vec{x})) \quad (3.4)$$

What this does, is to integrate over all the points $p_x(\vec{x})$ where the argument of the delta function is zero. The important thing is that care needs to be taken when the argument of the delta function is itself a function. In that case, a change of variables has to be performed, so that this is no longer the case. On wikipedia, this is done by defining the new variable $du = |\frac{d}{dx}g(x)|dx$, from which follows:

$$\delta(g(x)) = \sum_{x_0} \frac{\delta(x - x_0)}{|g'(x_0)|} \quad (3.5)$$

Where it is necessary to sum over each point x_0 for which $g(x_0) = 0$.

Which makes sense. Except, if you rewrite in terms of the new variable, you get an expression that does not strike me as necessarily the same:

$$\int f(x) \delta(g(x)) dx = \sum_n \int f(g_n^{-1}(u)) \delta(u) \left| \frac{d}{du} g_n^{-1}(u) \right| du \quad (3.6)$$

Where the sum n is over all functions g_n^{-1} that satisfy $g(g_n^{-1}(u)) = u$. For example, if $u = g(x) = \sin(x)$ then $g_n^{-1} = \arcsin(u) + n2\pi$, for any integer n . I feel more comfortable with the second path. The answer to my confusion is most likely the inverse function theorem (cf. section 2.12.2). In other words: this section needs revision! But hey, I flagged it.

Note that a delta function with a vector argument can be written as a product of the delta functions along each dimension.

Define a new set of variables $\vec{a} = \vec{u} - \mathbf{H}(\vec{x})$. It follows that $\vec{x} = \mathbf{H}_n^{-1}(\vec{u} - \vec{a})$ and the volume element $dx^n = \left| \frac{d}{d\vec{a}} \mathbf{H}_n^{-1}(\vec{u} - \vec{a}) \right| d\vec{a}$. Here, \mathbf{H}_n^{-1} are all the functions that satisfy $\mathbf{H}(\mathbf{H}_n^{-1}(\vec{u})) = \vec{u}$. The integral becomes:

$$p(\vec{u}) = \int d\vec{a} p_x(\mathbf{H}_n^{-1}(\vec{u} - \vec{a})) \left| \frac{d}{d\vec{a}} \mathbf{H}_n^{-1}(\vec{u} - \vec{a}) \right| \delta(\vec{a}) \quad (3.7)$$

At this point it is safe to evaluate the delta function integral. The result is:

$$p(\vec{u}) = p_x(\mathbf{H}_n^{-1}(\vec{u})) \left| \frac{d}{d\vec{a}} \mathbf{H}_n^{-1}(\vec{u} - \vec{a}) \right|_{(\vec{a}=0)} \quad (3.8)$$

3.3.1 Example: Sum of Random Variables

If the map \mathbf{H} is not bijective, for example because \vec{u} has lower dimensionality than \vec{x} , then a bijective map can be artificially constructed by introducing additional variables that are then also marginalized out. For example, if the goal is to calculate the probability of measuring some sum of random variables s , then you can define:

$$u_0 = s - \sum (x_i) u_1 = x_1 u_2 = x_2 u_{(n-1)} = x_{n-1} \quad (3.9)$$

The inverse is:

$$x_1 = u_1, x_2 = u_2, \dots, x_n = s - u_0 - \sum_{i < n} x_i \quad (3.10)$$

The argument in the delta function is transformed $\delta(s - \sum(x_i)) \rightarrow \delta(u_0)$. The determinant of the Jacobian $|J| = |\frac{d}{d\vec{u}} H^{-1}(\vec{u})|$ is given by $[\frac{d}{du_1} \mathbf{H}^{-1}, \frac{d}{du_2} \mathbf{H}^{-1}, \dots, \frac{d}{du_{n-1}} \mathbf{H}^{-1}]$. That's an nxn matrix where the first row is all -1 , the lower left is an $(n-1)x(n-1)$ identity matrix, and the lower right is a $(n-1)x1$ vector of 0 s. The determinant is one. Consequently the probability distribution of measuring a sum s is given by:

$$p(s) = \int dx^{n-1} p_x(x_1, x_2, x_3, \dots, s - \sum_{x < n} x_i) \quad (3.11)$$

Which turns out to be the convolution when the variables are independent.

3.3.2 Example: Lower Dimensional Random Variable

This was already the case for the sum of several random variables, in which case the dimensionality of the problem was reduced from many to one.

Consider the case many to fewer. Again, you just need to perform a change of variables. The new set of variables needs to have the same dimensionality as the old set of variables. The rest should follow pretty obviously.

Linear example: $\vec{y} = \mathbf{M}\vec{x}$ where y has 2 dimensions and x has 3.

Transform:

$$u_1 = y_1 - \vec{M}_1 \cdot \vec{x}, u_2 = y_2 - \vec{M}_2 \cdot \vec{x}, u_3 = x_3 \quad (3.12)$$

You can write this in terms of some invertible matrix $\mathbf{W} = [\mathbf{M}, [0, 0, 1]]$ as $\vec{u} = \mathbf{W}\vec{x}$, so that $dx^n = |\mathbf{W}^{-1}| du^n$. The integral is then:

$$p(y) = \int du^n p_x(\mathbf{W}^{-1}(\vec{y} - \vec{u})) \delta(u_1) \delta(u_2) |\mathbf{W}^{-1}| \quad (3.13)$$

3.4 Indicator Variables

Indicator variables are useful devices that simplify probability calculations involving binary outcomes, namely whether the outcome lies in a certain region of event space or not. Specifically, they allow one to translate set expressions into algebraic expressions.

Let $A \subseteq S$ be a subset of event space S (ex. $A \equiv$ "it rains tomorrow"). Then the indicator variable I_A :

$$I_A = \begin{cases} 1 & \text{if outcome in } A \\ 0 & \text{if outcome in } A^c \end{cases} \quad (3.14)$$

Which means that:

$$\mathbb{E}(I_A) = 1 \times p(A) + 0 \times p(A^c) = p(A) \quad (3.15)$$

To indicator variables for different events A, B, C, \dots can also be combined:

$$I_{A \cap B \cap C \dots} = I_A I_B I_C \dots \quad (3.16)$$

And indicator variables for the complement can be constructed trivially:

$$I_{A^c} = 1 - I_A \quad (3.17)$$

Whenever accounting for complicated combinations of events becomes overwhelming, indicator variables are often a good approach.

3.4.1 Example: The Party Problem

The party problem is a typical interview question, so I will include the two easier problems that do not make use of indicator variables. (This version of my notes also still has bitterness included.)

There are n drunk kids at a party that is presumably getting shut down by the fun police in Cambridge, MA. They (meaning the kids, presumably) grab their coats at random, and the problem is built around thinking about how many people wind up with the correct coat.

3.4.1.1 Lame Interview Question 1: Every one finds their coat

The common version of this problem asks "what is the probability that all of the kids wind up with the right coat". This is much easier than the general case. Imagine the list of party guests as a sequence $(1, 2, 3, 4, \dots, n-1, n)$ and the coats they grab as a random permutation of that sequence, such as $\alpha = (6, 1, 40, 21, 9, \dots)_n$. There are $n!$ such permutations, and they are all equally likely. The probability of all guests getting the correct coat is the probability that the permutation happens to be the single correct one, i.e. $p(\alpha = (1, 2, 3, 4, \dots, n-1, n)_n)$, which is the probability of one particular permutation, i.e. $p(\alpha = (1, 2, 3, 4, \dots, n-1, n)_n) = \frac{1}{n!}$.

3.4.1.2 Lame Interview Question 2: At least r people find their coat

This is still pretty easy, because the permutations that are correct are easily counted. If at least r coats are correctly assigned, then there are $\binom{n}{r}$ ways of choosing which of the r people wind up with the right coat. Then, while the location of r indices in the sequence is fixed, $n-r$ indices can be assigned arbitrarily.

The number of permutations where at least r are assigned correctly are then:

$$\binom{n}{r} (n-r)! \quad (3.18)$$

And the probability of at least r people finding their coat is the number of permutations multiplied with the probability of an individual permutation (that is, $\frac{1}{n!}$):

$$p(\# \text{ correct} \geq r) = \binom{n}{r} \frac{(n-r)!}{n!} = \frac{1}{r!} \quad (3.19)$$

Where $r \leq n$.

3.4.1.3 Not lame: Exactly r people find their coat

So, then, what's the probability that exactly nobody finds their coat? What's the probability that 2 people find their coat but nobody else does? What's the probability that r out of n people find their coat? After thinking quickly on your feet for two seconds, you realize that the answer is obviously:

$$p(\# \text{ correct} = r) = \frac{1}{r!} \sum_{s=0}^{n-r} \frac{(-1)^s}{s!} \quad (3.20)$$

The interviewer grunts ambiguously. They never call you back. You never find out why. You can't sleep. You can't eat. You become an anarchist and you declare war on the system.

The first two versions here are what I've come across in interview prep-type materials. I find them a bit annoying, because they represent particular cases that are much simpler than the general case. Applicants in the habit of studying stupid interview questions are rewarded because those answers are easily memorized (false positive). Applicants who intuit the complexity of the general problem, and who don't know the answer beforehand, might become overwhelmed during an interview and fail (false-ish negative).

/ rant

My approach here follows the extraordinary lecture notes <https://mast.queensu.ca/stat455/> by Glen Takahara at the University of Queensland.

As so often, indicator variables help:

$$I_A = \begin{cases} 1 & \text{if outcome in } A \\ 0 & \text{if outcome in } A^c \end{cases} \quad (3.21)$$

Which means that:

$$\mathbb{E}(I_A) = 1 \times p(A) + 0 \times p(A^c) = p(A) \quad (3.22)$$

To indicator variables for different events can also be combined:

$$I_{A \cap B \cap C \dots} = I_A I_B I_C \dots \quad (3.23)$$

And indicator variables for the complement can be constructed trivially:

$$I_{A^c} = 1 - I_A \quad (3.24)$$

Let A_i be the region of state space in which guest i grabbed the right coat. Let's say a *particular* subset $\{i\}_r$ of r guests grabs their correct coats (for example, $\{i\}_r = \{5, 9, 11, 24, \dots\}_r$), and that the set of remaining $n - r$ guests $\{j\}_{n-r} = \{1, 2, 3, \dots\}_n \setminus \{i\}_r$ grab the wrong coat. The area of state space that corresponds to this outcome is:

$$A_{\{i\}_r, \{j\}_{n-r}} = \bigcap_{i \in \{i\}_r} A_i \bigcap_{j \in \{j\}_{n-r}} A_j^c \quad (3.25)$$

The event that the outcome lies within that region of configuration space can be described with an indicator function:

$$I_{A_{\{i\}_r, \{j\}_{n-r}}} = \prod_{\{i\}_r} I_{A_i} \prod_{\{j\}_{n-r}} (1 - I_{A_j}) \quad (3.26)$$

And the probability of those *particular* r people finding their coat is it's expectation value, $p(A_{\{i\}_r, \{j\}_{n-r}}) = \mathbb{E}(I_{A_{\{i\}_r, \{j\}_{n-r}}})$. Good stuff.

The expression above will consist of a bunch of products of indicator variables that describe whether a particular guest wound up with the correct coat. We know that the product of indicator variables corresponds to an indicator variable for the *intersection* of the corresponding subsets of the state space. Explicitly, if it is a product of s indicator variables for some particular set of coats $\{k\}_s$:

$$\prod_k I_{A_k} = I_{\bigcap_k A_k} \quad (3.27)$$

And $\mathbb{E}(I_{\bigcap_k A_k}) = p(\bigcap_k A_k)$ is the probability of a particular s coats being picked up correctly, which is $(n-s)!/n!$. (There is no binomial factor, because it's one *specific* set of s coats).

Products of the sort $\prod_{i=1}^n (1 - x_i)$ can be expanded:

$$\prod_{i=1}^n (1 - x_i) = \sum_{s=0}^n (-1)^s \sum_{1 \leq i_1, \dots, i_s \leq n} x_{i_1} x_{i_2} \dots x_{i_s} \quad (3.28)$$

The sum $\sum_{1 \leq i_1, \dots, i_s \leq n}$ is over all possible sets of up to s indices that can be drawn from $\{1, 2, 3, \dots, n\}$. A simple example with $n = 3$:

$$(1 - x_1)(1 - x_2)(1 - x_3) = \underbrace{1}_{s=0} - \underbrace{(x_1 + x_2 + x_3)}_{s=1} + \underbrace{(x_1 x_2 + x_2 x_3 + x_1 x_3)}_{s=2} - \underbrace{(x_1 x_2 x_3)}_{s=3} \quad (3.29)$$

The number of terms of order s is the amount of ways that s indices can be sampled from n indices, $\binom{n}{s}$.

Returning to the original problem, then:

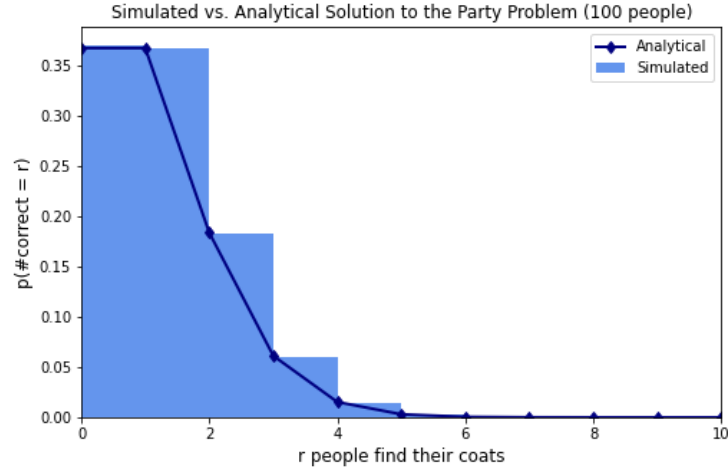


Figure 3.1: Simulated and analytically calculated probabilities that exactly r people at a $n = 100$ people party randomly pick up their coat.

$$\begin{aligned}
 I_{A_{\{i\}_r, \{j\}_{n-r}}} &= \prod_{\{i\}_r} I_{A_i} \prod_{\{j\}_{n-r}} (1 - I_{A_j}) \\
 &= \sum_{s=0}^{n-r} (-1)^s \sum_{\substack{n-r \leq j_1, \dots, j_s \leq n \\ \text{sum of } \binom{n-r}{s} \text{ terms}}} \underbrace{\prod_{\{i\}_r} I_{A_i} \prod_{\{j\}_s} I_{A_j}}_{\text{product of } r+s \text{ terms}}
 \end{aligned} \tag{3.30}$$

By linearity of expected value, and using the relationship of the expected value of indicator variables to their probabilities:

$$\begin{aligned}
 \mathbb{E}(I_{A_{\{i\}_r, \{j\}_{n-r}}}) &= \sum_{s=0}^{n-r} (-1)^s \sum_{\substack{n-r \leq j_1, \dots, j_s \leq n \\ \text{sum of } \binom{n-r}{s} \text{ terms}}} p \left(\underbrace{I_{\cap_{\{i\}_r} A_i \cap_{\{j\}_s} A_j}}_{r+s \text{ coats picked up correctly}} \right) \\
 &= \sum_{s=0}^{n-r} (-1)^s \binom{n-r}{s} \frac{(n-r-s)!}{n!}
 \end{aligned} \tag{3.31}$$

Now, this is already a pretty neat expression for the probability of a *particular* r coats being picked up correctly (i.e. Joe, Mary, and "Hans" picked up the right coat). We don't really care which of the r guests got lucky, though, so since there are $\binom{n}{r}$ ways of r coats having been picked out correctly, you sum over all of them:

$$\begin{aligned}
 p(\# \text{ correct} = r) &= \binom{n}{r} \sum_{s=0}^{n-r} (-1)^s \binom{n-r}{s} \frac{(n-r-s)!}{n!} \\
 &= \sum_{s=0}^{n-r} (-1)^s \frac{n!}{r!(n-r)!} \frac{(n-r)!}{s!(n-r-s)!} \frac{(n-r-s)!}{n!} \\
 &= \frac{1}{r!} \sum_{s=0}^{n-r} \frac{(-1)^s}{s!}
 \end{aligned} \tag{3.32}$$

3.5 Copulas

3.6 Relationships Between Distributions

3.7 Large Deviation Theory

3.7.1 Gaertner-Ellis Theorem

3.7.2 Example: Sum of Uniform Random Variables

4 Information Theory

Contents of this chapter

4.1	Entropy	45	4.3	Kullback-Leibler Divergence	45
4.2	Mutual Information	45			

4.1 Entropy

4.2 Mutual Information

4.3 Kullback-Leibler Divergence

5 Stochastic Processes and Time Series Analysis

Contents of this chapter

5.1	Markov Chains	47	5.3	Hidden Markov Models	47
5.2	Martingales	47	5.4	Ito Calculus	47
5.2.1	Martingale Convergence Theorem	47			

5.1 Markov Chains

5.2 Martingales

5.2.1 Martingale Convergence Theorem

5.3 Hidden Markov Models

5.4 Ito Calculus

6 Statistics

Contents of this chapter

6.1	Directional Statistics	49	6.3.1	Sources	50
6.1.1	Mean Direction	49	6.4	R^2 Value	51
6.1.2	Dispersion	49	6.5	Regression Diagnostics	51
6.2	Features	50	6.6	t-Statistics	51
6.2.1	Dense Features, Sparse Features	50	6.7	AIC and BIC	51
6.3	Multicollinearity	50	6.8	Factor Regression	51
			6.9	Factor Models	51

6.1 Directional Statistics

6.1.1 Mean Direction

The mean direction of a collection of i vectors $\{\mathbf{x}\}_i$ is (?):

$$\langle \mathbf{x} \rangle = \frac{\mathbf{x}_s}{\|\mathbf{x}_s\|_2} \quad (6.1)$$

Where

$$\mathbf{x}_s = \sum_i \mathbf{x}_i \quad (6.2)$$

The mean vector is not defined in case $\|\mathbf{x}_s\|_2 = 0$. The alternative way to calculate the mean direction might make use of angles but apparently that creates ambiguity with respect to the choice of a "zero" angle (cf. footnote in Damask (2019)).

6.1.2 Dispersion

Dispersion is a measure of the variance on the direction of a set of vectors $\{\mathbf{x}\}_i$. For a system of vectors, for example an eigenbasis, both common and differential modes of estimation exist. One way of measuring directional dispersion is to look at the *mean resultant length*:

$$\mu_r = \frac{\|\mathbf{x}_s\|}{N} \quad 0 \leq \mu_r \leq 1 \quad (6.3)$$

Circular variance may be defined as $\sigma_c = 1 - \mu_r$, but apparently there is an issue with generalization to higher dimensions. (TO DO)

An alternative is a model based approach, for example based on the von Mises - Fisher Distribution.

In the case of a basis, Damask (2019) states that, in order to understand the drivers for variation, dispersion parameters should be calculated for each descending subspace of the basis, first including all of the eigenvectors, then excluding the first eigenvector, etc.

6.2 Features

Features are sources of information that hopefully allow conclusions towards some kind of quantity of interest. Some people like to call them independent variables.

6.2.1 Dense Features, Sparse Features

Sparsity and density refer to the fraction of a matrix that is zeros. In the context of features, sparse features typically refer to feature vectors that have many zeros in them. I.e., $[1, 2, 5, 2, 6, 3, 0]$ would be a dense feature vector and $[2, 0, 0, 0, 0, 3, 0, 0, 1]$ would be a sparse feature vector. Typically, sparsity is touted as an advantage because it allows for lossless compressed representations of high-dimensional data, which has computational advantages. Another advantage, though, is that sparse representations can be more interpretable by containing information in an inherently more condensed fashion. They are also thought to prevent overfitting: many regularization techniques aim to minimize the number of parameters or features used in a predictive model, which amounts to biasing an algorithm towards learning sparse coefficients. An example is ridge regression.

6.3 Multicollinearity

So far this section is based on Software (n/a).

Multicollinearity, or collinearity, occurs when features are linearly correlated with each other. The effects are horrible. They include inaccurate estimates of the regression coefficients, higher standard errors of the regression coefficients, lower partial t-tests for the regression coefficients, falsely insignificant p-values and decreased predictive power of the model. And other things.

6.3.0.1 Detection

Scatter Plots Scatterplots provide a visual test for collinearity by hopefully exposing relationships between independent variables. This is subjective and unreliable, but people love plots.

Variance Inflation Factors (VIF) A VIF over 10 it said to indicate collinear variables.

Eigenvalues of the Correlation Matrix Linear relationships between two or more variables cause the corresponding rows of the correlation matrix to be identical or very similar. Correspondingly, the matrix will be singular or near-singular, which will manifest itself through zero or near-zero eigenvalues. The conditioning number, given by the largest eigenvalue divided by the smallest eigenvalue, are a quick way to test for this. A large conditioning number indicates collinearity.

Regression Coefficients Collinearity increases the standard error of the regression coefficients because it allows for the variation of the dependent variable to be explained in terms of a greater variety of different weights assigned to the collinear variables. Counterintuitive results for the regression coefficients may be the result of collinearity.

6.3.1 Sources

Data Collection

Physical Constraints

Over-defined Model

Model Choice or Specification

Outliers

6.3.1.1 Remedies

Dimensionality Reduction SVD, PCA, NNMF and other dimensionality reduction techniques allow for the feature space to be shrunk in a way that aims to optimally preserve information. Any technique worth its salt will either collapse or filter collinear variables into a reduced-rank representation of the information in the dataset.

Regularization Certain forms of regularization are similar in spirit to dimensionality reduction, except that, rather than addressing the dataset, they reduce the parameters used by a model to fit to the data. The canonical example is ridge regression. Ridge regression penalizes the use of a larger number of parameters and, if two variables are collinear, will tend to push the weight of one of them towards zero. It is important to standardize the variables before fitting the model, so that the regression weights for different variables are on the same scale. Ridge regression remains controversial (Software n/a).

6.4 R^2 Value

6.5 Regression Diagnostics

6.6 t-Statistics

6.7 AIC and BIC

6.8 Factor Regression

6.9 Factor Models

7 Linear Regression

Contents of this chapter

7.1	Ordinary Least Squares Regression .	53	7.1.6	Q-plots	55
7.1.1	The Normal Equations, Analytical Least Squares Estimator	53	7.1.7	Variance Inflation Factor . .	55
7.1.2	The Quick Way to $\hat{\beta}$	54	7.2	Ridge Regression	55
7.1.3	The Long Way to $\hat{\beta}$	54	7.2.1	Analytical Ridge Estimator .	55
7.1.4	Projection Matrix	54	7.2.2	Bayesian Perspective on Ridge Regression	55
7.1.5	Bayesian Perspective on Least Squares Regression . .	55			

There is no escaping from linear regression but hardly anyone seems to agree on how to do it properly. Even worse, people seem to expect you to know things like the normal equations, which you'll never, ever need outside of a job interview.

7.1 Ordinary Least Squares Regression

7.1.1 The Normal Equations, Analytical Least Squares Estimator

Section 2.2 discussed the case of the linear system of equations:

$$\mathbf{A}\mathbf{x} = \mathbf{b} \quad (7.1)$$

When $\mathbf{A}^{m \times n}$ with $m > n$, so that the system is overdetermined. This is, of course, the starting point for least squares regression, only that the convention is to use different letters:

$$\mathbf{X}\beta = \mathbf{y} \quad (7.2)$$

And that, seeing that the system is overdetermined, one looks for an approximate solution $\hat{\beta}$, so that

$$\mathbf{X}\hat{\beta} + \epsilon = \mathbf{y} \quad (7.3)$$

A natural approach for picking an approximate solution $\hat{\beta}$ is to look for the projection of \mathbf{y} in the column space of \mathbf{X} . That is, since the column rank $\leq n$ of \mathbf{X} is insufficient to express m -dimensional \mathbf{y} exactly in terms of only m coefficients β , we look for the n -dimensional shadow $\hat{\beta}$ of some hypothetical higher dimensional exact solution.

The projection has the property that it maximizes the dot product $(\mathbf{X}\hat{\beta}) \cdot \mathbf{y}$, and hence minimizes the length of the difference vector ϵ . In turn, the length of the difference vector ϵ is $\sqrt{\epsilon \cdot \epsilon}$, which is monotonic to $\epsilon \cdot \epsilon = \sum_i^m \epsilon_i^2$. That means that finding the projection of \mathbf{y} in the column space of \mathbf{X} minimizes the L_2 norm of ϵ , also known as *least squares error*.

There are two ways to go about finding $\hat{\beta}$.

7.1.2 The Quick Way to $\hat{\beta}$

By construction, the vector ϵ is orthogonal to the column space of \mathbf{X} . Which means:

$$\begin{aligned} \mathbf{X}^T \epsilon &= 0 \\ \mathbf{X}^T (\mathbf{X} \hat{\beta} - \mathbf{y}) &= 0 \\ \mathbf{X}^T \mathbf{X} \hat{\beta} &= \mathbf{X}^T \mathbf{y} \\ \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned} \tag{7.4}$$

Making use of the fact that $\mathbf{X}^T \mathbf{X}$ is square and therefore hopefully invertible.

7.1.3 The Long Way to $\hat{\beta}$

Loss functions play a central role in computational statistics (for example when regularization is introduced), and therefore it is of interest to approach finding $\hat{\beta}$ by instead minimizing the least square error. This requires:

$$\frac{d}{d\hat{\beta}} L_2(\epsilon) = 0 \tag{7.5}$$

where

$$\begin{aligned} L_2(\epsilon) &= (\mathbf{X} \hat{\beta} - \mathbf{y})^T (\mathbf{X} \hat{\beta} - \mathbf{y}) \\ &= \hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - \hat{\beta}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\beta} + \mathbf{y}^T \mathbf{y} \end{aligned} \tag{7.6}$$

Taking derivatives with respect to a vector is covered in section 2.12.

It follows:

$$\begin{aligned} \frac{d}{d\hat{\beta}} \left(\mathbf{X}^T \mathbf{X}^T \underbrace{\mathbf{X} \hat{\beta}}_{u(\hat{\beta})} \right) &= \frac{d}{du} (u^T u) \frac{d}{d\hat{\beta}} u = 2u^T \frac{d}{d\hat{\beta}} u = 2\hat{\beta}^T \mathbf{X}^T \mathbf{X} \\ \frac{d}{d\hat{\beta}} \hat{\beta}^T \mathbf{X}^T \mathbf{y} &= \mathbf{y}^T \mathbf{X} \\ \frac{d}{d\hat{\beta}} \mathbf{y}^T \mathbf{X} \hat{\beta} &= \mathbf{y}^T \mathbf{X} \\ \frac{d}{d\hat{\beta}} \mathbf{y}^T \mathbf{y} &= 0 \end{aligned} \tag{7.7}$$

So that

$$\begin{aligned} \frac{d}{dx} L_2(\epsilon) = 0 &= 2\hat{\beta}^T \mathbf{X}^T \mathbf{X} - 2\mathbf{y}^T \mathbf{X} \\ \hat{\beta}^T \mathbf{X}^T \mathbf{X} &= \mathbf{y}^T \mathbf{X} \\ \mathbf{X}^T \mathbf{X} \hat{\beta} &= \mathbf{X}^T \mathbf{y} \\ \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned} \tag{7.8}$$

7.1.4 Projection Matrix

If $\hat{\mathbf{y}} = \mathbf{X} \hat{\beta}$ is the projection of \mathbf{y} in the column space of \mathbf{X} , then, based on the result for $\hat{\beta}$, the projection matrix is $\mathbf{P} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$. In a fully determined system, $\mathbf{P} = \mathbf{I}$. Projection matrices have eigenvalues that are either 1 or 0, corresponding to dimensions that are kept or discarded during the projection operation.

7.1.5 Bayesian Perspective on Least Squares Regression

7.1.6 Q-plots

7.1.7 Variance Inflation Factor

7.2 Ridge Regression

7.2.1 Analytical Ridge Estimator

7.2.2 Bayesian Perspective on Ridge Regression

8 Bayesian Data Analysis

Contents of this chapter

9 Unsupervised Learning

Contents of this chapter

9.1	Blind Source Separation	59	9.3.3	Bayesian SVD	61
9.2	Independent Component Analysis	59	9.4	Principal Component Analysis (PCA)	61
9.3	Eigenanalysis, Singular Value Decomposition of Data Matrices	59	9.4.1	Relationship between PCA and SVD	63
9.3.1	Example: Eigenfaces and Facial Recognition	60	9.4.2	Tracking Principal Components over Time	63
9.3.2	Tracking an Eigensystem over Time	61	9.5	$\mathbf{X} = \mathbf{WH}$ Non-Negative Matrix Factorization	64

9.1 Blind Source Separation

Blind Source Separation or Blind Signal Separation (BSS) is the separation of mixed signals. An example is the cocktail party problem, where a microphone is recording the conversations of many people speaking at the same time, making the recording unintelligible. Successful BSS would be able to extract the voice of only a single person.

9.2 Independent Component Analysis

Independent Component Analysis (ICA) is a decomposition of data into components that have vanishing mutual information.

9.3 Eigenanalysis, Singular Value Decomposition of Data Matrices

Singular Value Decomposition (SVD) is a matrix decomposition of the form $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, where \mathbf{U} and \mathbf{V} are both unitary (cf. section 2.6). The decomposition always exists for a general complex matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$.

If $\mathbf{X} \in \mathbb{R}^{N \times p}$ is a data matrix with N samples in the rows and p features in the columns, then SVD allows for the decomposition of the data into p linearly independent components, ranked by their explained variance (i.e. their strength in the dataset). The basis of the decomposition turns out the same as for PCA (cf. section ??) except that it is arrived at slightly differently, because PCA involves diagonalizing the covariance matrix. The decomposition of the dataset can be understood as follows. In the SVD of the data matrix:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (9.1)$$

The individual rows of \mathbf{X} , i.e. the individual data points, are expressed as:

$$\mathbf{x}_i^T = \sum_{j=1}^p u_{i,j} \sigma_j \mathbf{v}_j^T \quad (9.2)$$

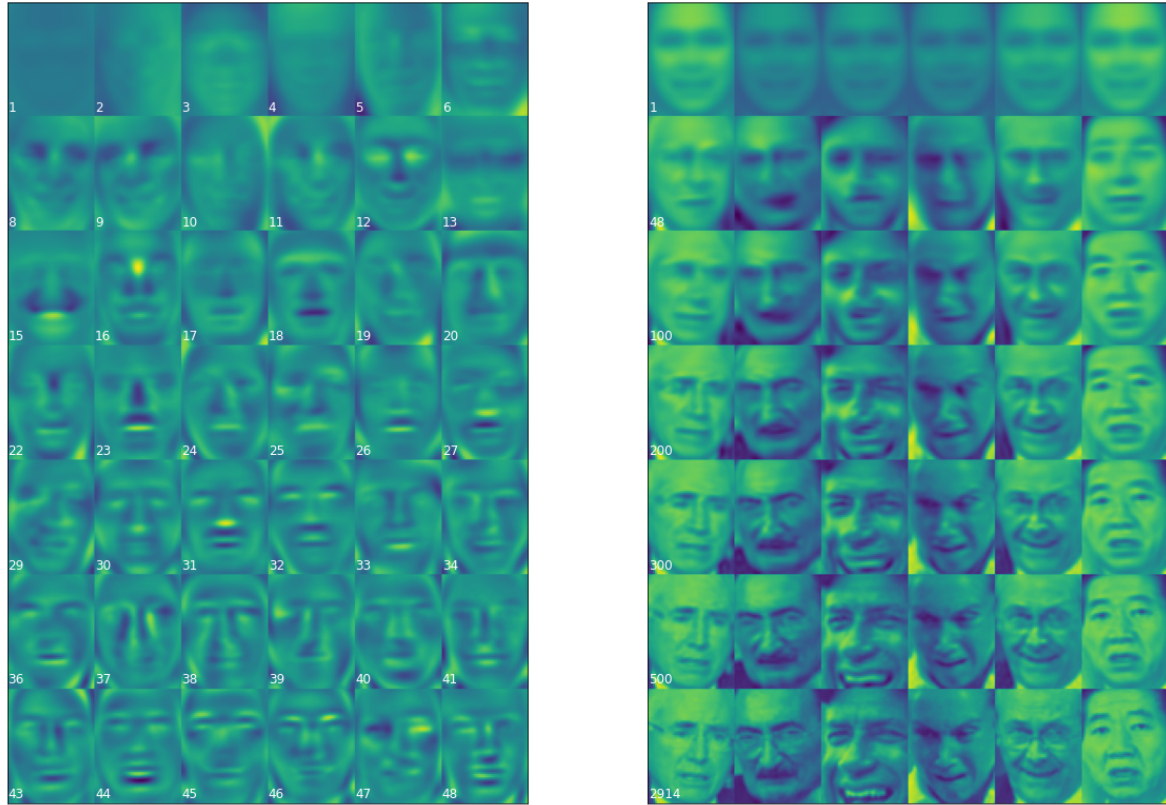


Figure 9.1: Left: The first 48 eigenfaces. The colorscale is consistent across the images. As can be expected, the eigenfaces seem to have an ordering from more general features, that are highly prevalent in the dataset, towards more specific features. Right: Five sample portraits from the dataset approximated using different numbers of eigenvectors. 2914 corresponds to the original image, which had 2914 pixels (degrees of freedom).

Where \mathbf{x}_i is the i th data point. Hence the right singular vectors \mathbf{v}_i are the normalized basis vectors, and the columns of \mathbf{U} give the coefficients the basis expansion, together with the singular values σ_i , which give a measure of the global strength of the corresponding basis vector.

SVD is scalable to very large datasets and finds many applications in the wild, including page rank, facial recognition, recommendation algorithms etc.

9.3.1 Example: Eigenfaces and Facial Recognition

One famous result are the so-called eigenfaces. The data matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$ consists of N pictures of faces that each have p pixels. The right singular eigenvectors yield p eigenfaces in terms of which any of the N pictures can be expressed.

Below are the first eigenfaces extracted from the "Labeled Faces in the Wild" dataset, which includes 13233 portraits with $62 \times 47 = 2914$ pixels each. Running SVD on the data matrix $\mathbf{X} \in \mathbb{R}^{13233 \times 2914}$ yields the eigenfaces shown in Figure 9.1.

Facial recognition may be performed by projecting a new face into the space of eigenfaces and matching to the coefficients of a known subject. This could be done using the euclidean distance, or it could be a classifier. In case of classification, the dimensionality reduction that is enabled by approximating images with a smaller set of eigenvectors may be critical to making the problem tractable by overcoming the curse of dimensionality.

9.3.2 Tracking an Eigensystem over Time

As discussed in section 2.6, the sign of the basis vectors can be flipped without affecting the validity of the singular value decomposition. That means that if the SVD is performed on some system may equally well return, say, a left-handed or a right-handed coordinate system. This becomes a problem when the results of repeated SVDs are supposed to be compared, for example to study the evolution of a system over time. Figure 9.2 shows the eigenbasis of a bivariate Gaussian with principal axes slowly rotating over a 180-degree angle. The top shows how the eigenbases that is found flips back and forth, so that there is a double-trajectory corresponding to results with positive and negative sign on an eigenvector. The trajectory is not described by an injective function, which complicates analysis.

9.3.2.1 Heuristic Method

Say an SVD is performed on a slowly-evolving system at time $t = 0$ and at time $t = 1$. To compensate the sign flips, a heuristic method is to first match the basis vectors at $t = 0$ to the basis vectors at $t = 1$ using a distance metric that is immune to the sign of the vectors, for example the absolute value of the dot product. (This assumes that the vectors are similar enough at time $t = 0$ and $t = 1$ that the matching can be done unambiguously. For the higher-order singular vectors this might be a problem, because they have lower numerical certainty.) Once the vectors are matched, the sign of the inner product between the vectors at $t = 0$ and $t = 1$ may be compared, and the sign flipped accordingly.

This method works, but the problem is that the sign of the vectors is somewhat arbitrarily pinned relative to the result at $t = 0$. That is, if, for example, the coordinate system found at $t = 0$ was left-handed, then the time series of eigenbases will be adjusted to be left-handed. If the point of comparison had instead been the SVD performed at some other time $t = t'$, then one might have wound up with a right-handed coordinate system instead. It is desirable to find a consistent orientation.

9.3.2.2 Consistent Method

Damask (2019) recently developed a method that can be used to find a consistently oriented basis. The method relies on reconstructing the rotations and reflections necessary to transform the eigenbasis in question to the natural basis \mathbf{I} as reference. While rotations preserve the orientation of a basis, but reflections do not. When the reflections of an eigenbasis with respect to the natural basis are known, then they can be undone by flipping them back in place.

9.3.2.3 Rank Order Changes

The singular vectors or eigenvectors of a system are labeled only in terms of their associated singular values or eigenvalues. In order to track a singular vector throughout rank order changes, it is necessary to figure out a way to attach a separate label, for example by looking at the "content" of a particular singular vector. How well that works has to be figured out in context, I am not currently aware of a generally valid solution.

9.3.3 Bayesian SVD

To do! Cf. <https://ieeexplore.ieee.org/document/7336426>

It's awful to be tied to the assumption of normality and it's awful to not know the uncertainty of my singular vectors!

9.4 Principal Component Analysis (PCA)

Principal Component Analysis (PCA) is a decomposition of data into linearly uncorrelated components, ordered by their explained variance. It turns out that these axes are the right-singular eigenvectors \mathbf{V} that are found from SVD. SVD and PCA have correspondence (cf. section 9.4.1). My guess is that the advantage of PCA is that its results are interpretable in terms of commonly used summary statistics of the

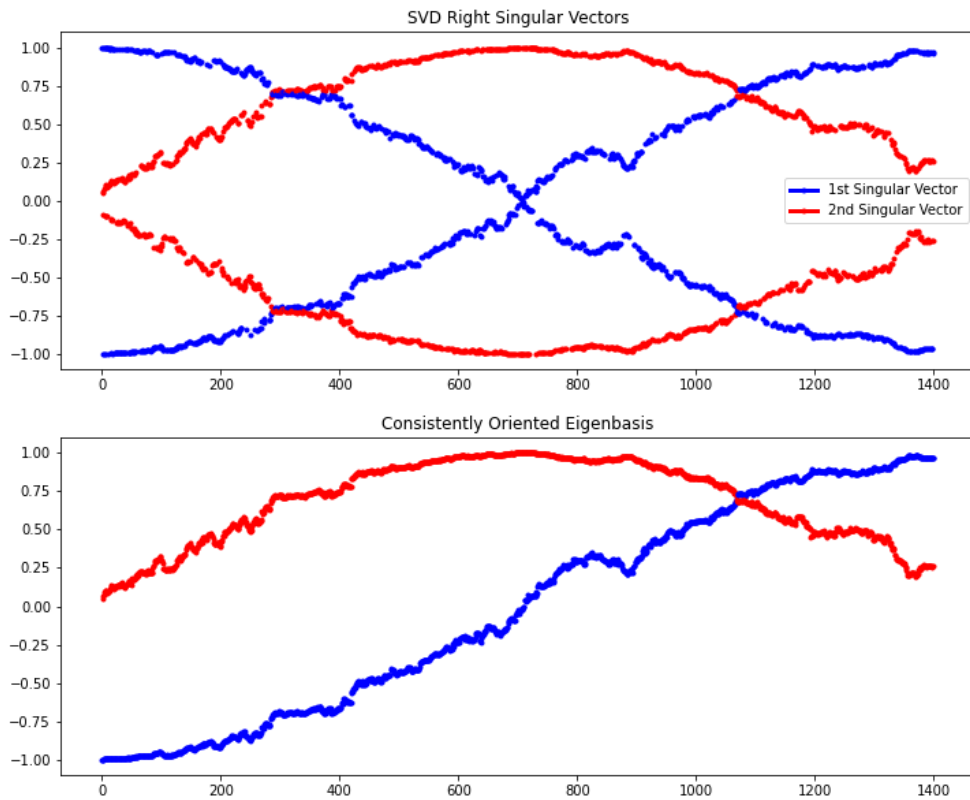


Figure 9.2: Right singular vectors extracted using SVD, showing sign flips (top) and with a consistently oriented basis (bottom). The underlying data is a bivariate gaussian with principal axes undergoing a 180 degree rotation. With the consistently oriented basis, the singular vectors trace out a one-to-one trajectory that can be analyzed.

distribution (namely the variances or Pearson correlations of the data).

More elaborately, for a data matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$, PCA extracts the ordered, rank- p , orthonormal basis in which the $p \times p$ covariance matrix of \mathbf{X} is diagonal. The basis vectors are called the principal axes or principal directions of the data, and their ordering is by the magnitude of the variance that they explain.

The property that the covariance matrix is diagonal in the principal component basis means that projecting the data onto any of the basis vectors extracts a linearly uncorrelated component of the data that has variance corresponding to the corresponding eigenvalue of the covariance matrix.

Just like with SVD, the ordering of the principal component basis in terms of their explained variance allows for lower-rank approximations of the data matrix to be constructed (section 2.6). Just like SVD, the lower rank approximations \mathbf{X} are the best possible approximations with respect to the Frobenius Norm (cf. section 2.11.3).

The $p \times p$ covariance matrix \mathbf{C} of \mathbf{X} is:

$$\mathbf{C} = \frac{(\mathbf{X} - \langle \mathbf{X} \rangle)^T (\mathbf{X} - \langle \mathbf{X} \rangle)}{n - 1} \quad (9.3)$$

Where $(\mathbf{X} - \langle \mathbf{X} \rangle)^T$ is often referred to as the *centered* data matrix. The covariance matrix is a symmetric, positive-definite matrix that can be diagonalized with orthonormal eigenvectors \mathbf{V} and positive (or vanishing) eigenvalues λ_i :

$$\mathbf{C} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \quad (9.4)$$

Where the eigenvectors in \mathbf{V} are ordered so that the eigenvectors along the diagonal of $\mathbf{\Lambda}$ have decreasing magnitude. The eigenvectors are the principal axes or principal directions of the data.

9.4.1 Relationship between PCA and SVD

This is based on a great Stack Exchange answer (amoeba 2015).

Let the singular value decomposition of the centered data matrix be:

$$(\mathbf{X} - \langle \mathbf{X} \rangle) = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \quad (9.5)$$

Then:

$$\mathbf{C} = \frac{\mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T}{n - 1} = \mathbf{V} \frac{\mathbf{\Sigma}^2}{n - 1} \mathbf{V}^T \quad (9.6)$$

That means that:

- The principal axes are the right-singular vectors \mathbf{V} that are obtained during SVD.
- The singular values and the eigenvalues of the covariance matrix are related via $\lambda_i = \frac{\sigma_i^2}{n-1}$.

9.4.2 Tracking Principal Components over Time

Given the direct parallel between PCA and SVD, the identical issue with sign flips emerges with the principal component basis, and the remedy is the same. See section 9.3.2.

9.5 $X = WH$ Non-Negative Matrix Factorization

Non-Negative Matrix Factorization (NNMF) is an approximate decomposition technique for positive data. It can be used, for example, for dimensionality reduction, blind source detection or topic extraction. The reputation is that it can often generate decompositions that are "sparse and meaningful".

Following the convention in Gillis (2014), it works by expressing a positive data matrix $\mathbf{X} \in \mathbb{R}_+^{p \times n}$ in terms of two smaller positive matrices $\mathbf{W} \in \mathbb{R}_+^{p \times r}$ and $\mathbf{H} \in \mathbb{R}_+^{r \times n}$. Here the n columns of the data matrix are data points and the rows are the p features. An excellent introduction to NNMF is found in Colyer's (2019) blog post and in Gillis (2014). Dimensionality reduction derives from the ability to adjust the rank of the decomposition via the dimension r . The decomposition is simply:

$$\mathbf{X} \approx \mathbf{WH} \tag{9.7}$$

Where the error $\|\mathbf{X} - \mathbf{WH}\|_\alpha$ is minimized with respect to some norm (normally, the Frobenius norm with $\alpha = 2$). The decomposition allows the interpretation of $\mathbf{W} \in \mathbb{R}^{p \times r}$ as a matrix with a number of r p -dimensional strictly positive basis vectors and of $\mathbf{H} \in \mathbb{R}^{r \times n}$ as a matrix of coefficients that express the n data points in terms of linear sums of the r basis vectors with strictly positive coefficients.

One might ask: why NNMF when I can create linear decompositions of my data with SVD and PCA? After all, these techniques are well established and are known to give the optimal approximation under the Frobenius norm. In my perception, the key difference is that NNMF solves the constrained problem of forcing both coefficients and basis vectors to be strictly positive. For many data types, negative values (word frequencies, pixel values) are unnatural, and so the results of NNMF are immediately more interpretable. NNMF is promising in particular in situations where a signal arises from the addition of a number of positive signals. A canonical example is hyperspectral imaging, where, assuming incoherent light, the measured spectrum is the linear sum of the spectra of the individual light sources. If all goes well, the basis vectors from NNMF of a hyperspectral image will be the source spectra. In the context of natural language processing, NNMF on a collection of documents can be interpreted as *topics*. Figure 9.3 shows the NNMF-based decomposition and approximation of the "labeled faces in the wild" dataset, analogous to the SVD of the dataset shown in Figure 9.1. The two are quite different. While for SVD, the faces look as if they are gradually coming into focus, for NNMF the faces are at first not recognizable. The portraits also become brighter in the plot as more basis components are added, because, in contrast to the SVD case, the composition is purely additive.

The drawback is that NNMF is computationally more difficult (possibly NP hard), nor is the matrix approximation optimal, nor is the decomposition unique.



Figure 9.3: Left: The first 34 basis faces stored in the \mathbf{W} matrix extracted using NNMF under Frobenius norm from the "Labeled Faces in the Wild" Dataset. The basis images are normalized and shown on a log scale because they have very different contrast and brightness. The normalized images are then all shown on the same color scale. Right: Five sample portraits from the dataset approximated using different numbers of basis faces. The NNMF was done for a value of $r = 300$. The original portraits had 2914 pixels. In contrast to SVD, the basis vectors obtained through NNMF do not necessarily have an internal ordering in terms of how much of the variance in the dataset they explain.



Figure 9.4: Left: The first 6 basis faces stored in the \mathbf{W} matrix extracted using NMF under Frobenius norm from the "Labeled Faces in the Wild" Dataset. Right: Sample images expressed in the reduced basis.

Bibliography

- amoeba (2015), 'Relationship between svd and pca.'
URL: <https://stats.stackexchange.com/questions/134282/relationship-between-svd-and-pca-how-to-use-svd-to-perform-pca>
- Barnes, R. J. (n/a), 'Matrix differentiation (and some other stuff)'.
URL: <https://atmos.washington.edu/~dennis/MatrixCalculus.pdf>
- Bogart, K. P. (2004), *Combinatorics through guided discovery*, Kenneth P. Bogart.
URL: <https://bogart.openmathbooks.org/>
- Colyer, A. (2019), 'The why and how of nonnegative matrix factorization'.
URL: <https://blog.acolyer.org/2019/02/18/the-why-and-how-of-nonnegative-matrix-factorization/>
- Damask, J. (2019), 'A consistently oriented basis for eigenanalysis', *arXiv preprint arXiv:1912.12983*.
URL: <https://arxiv.org/abs/1912.12983>
- Friedman, J., Hastie, T. & Tibshirani, R. (2001), *The elements of statistical learning*, Vol. 1, Springer series in statistics New York.
- Gera, R. (2009), 'Numerical linear algebra lecture notes, chapter 7'.
URL: <http://faculty.nps.edu/rgera/MA3042/2009/ch7.4.pdf>
- Gillis, N. (2014), 'The why and how of nonnegative matrix factorization'.
URL: <https://arxiv.org/abs/1401.5226>
- Kronenburg, M. (2011), 'The binomial coefficient for negative arguments', *arXiv preprint arXiv:1105.3689*.
URL: <https://arxiv.org/pdf/1105.3689.pdf>
- Mathworks (n/a), 'Eigenvalue and singular value decompositions'.
URL: <https://www.mathworks.com/content/dam/mathworks/mathworks-dot-com/moler/eigs.pdf>
- Software, N. S. (n/a), 'Ridge regression'.
URL: https://ncss-wpengine.netdna-ssl.com/wp-content/themes/ncss/pdf/Procedures/NCSS/Ridge_Regression.pdf