



25 Key Equations in Machine Learning

Insight Hatch Machine Learning Reviews - Vol 1

to-ohru iwanami {the hatch keeper}

2024.11.27

version 0.10.0

InsightHatch: “25MLEqs”

IH-25MLEqs-v0.10.0

Contents

Prefice and Introduction	i
The 25 (uhh 24)	iii
1 Gradient Descent	1
1.1 What is behind the equation	3
1.2 Structure of the Parameter Vector θ	3
1.3 Typical Value of the Learning Rate α	4
1.4 Math Behind the Gradient $\nabla J(\theta_j)$	4
1.5 Gradient Descent in R	6
1.6 Linking Gradient Descent to Our Polynomial Example	9
1.7 Additonal ideas to explore	11
2 Normal Distribution	13
2.1 What is Behind the Equation	15
2.2 The Players	16
2.3 Historical Context	17
2.4 Relative Standard Deviation (RSD)	17
2.5 Understanding the Notation $f(x \mu, \sigma^2)$	18
2.6 The Nature of the Exponential Function in Normal Distribution	18
2.7 Non-existence of an Analytical Anti-Derivative	18
2.8 Central Limit Theorem	19
2.9 Applications in Machine Learning	20
2.10 Suggested Additional Content for the Chapter	21
3 Z-Score	23
3.1 What is Behind the Equation	25
3.2 Applications in Machine Learning	25
3.3 Typical Range for Z-Scores	26
3.4 Mathematical Notation and Concept	26
3.5 Mathematical Insights: Area Under the Normal Curve	26
3.6 Example: Z-Score in R	27
4 Sigmoid Function	29
4.1 What is Behind the Equation	31
4.2 General Overview of Applications of Sigmoid Functions	31
4.3 Application to Machine Learning: Sigmoid Function in Binary Classification	32

4.4	How Sigmoid Works in Logistic Regression	32
4.5	Example Walkthrough	33
4.6	Why Sigmoid is Useful for Optimization	34
4.7	Deriving the Derivative of the Sigmoid Function	35
4.8	Typical Ranges or Values for the Sigmoid Function	36
4.9	Important Mathematical Identities	36
4.10	Comparative Functions	37
5	Correlation	41
5.1	What is Behind the Equation	43
5.2	Historical Context	45
5.3	Understanding the Notation	45
5.4	Applications in Machine Learning	45
5.5	Correlation in R	45
5.6	Scatter Plots with Linear Overlays	46
6	Cosine Similarity	49
6.1	What is Behind the Equation	51
6.2	Historical Context	54
6.3	Cosine Similarity - The basics and the usage	56
6.4	Cosine Similarity and Linear Separability	58
6.5	Relation to Spherical Geometry:	62
6.6	Cosine Similarity and Kernel Methods:	62
6.7	Cosine Similarity and Information Retrieval:	62
6.8	Cosine Similarity and the Norms of Vectors:	63
6.9	Can Cosine Similarity vs What can be used as a Norm?	63
6.10	Cosine Distance:	64
6.11	Triangle inequality for cosine similarity	67
6.12	Cosine Similarity Applications in Machine Learning	69
6.13	Cosine Similarity in R	69
6.14	Cosine Similarity in R, 3D-plot	70
7	Naive Bayes	73
7.1	Naive Baye's - overview	75
7.2	Introduction to Bayesian Theory	76
7.3	Fundamental Bayes' Probability Space	79
7.4	Kolmogorov Axioms of Probability	80
7.5	Probability Space Graph over Ω	84
7.6	Maringal and Bayes' probability space sums	88
7.7	bayes' by example	89
7.8	Sum and Product Rules of Probability	91
7.9	Connectinge Bayes' to Naive Bayes'	93
7.10	Raw Bayes Text Here First	94
7.11	Modifications in Venn Diagram Representation	95
7.12	Modifications to Your Original Bayesian Framework	95
7.13	Naive Bayes' as Filter	96
7.14	Example to classical Bayes' vs. Naive Bayes'	98
7.15	USing Bayes and the tie to physics	99

7.16 Using Naive Bayes	101
7.17 after area	103
7.18 Bayes start building	103
8 Maximum Likelihood	105
8.1 What is Behind the Equation	107
8.2 Historical Context	108
8.3 MLE General Introduction	110
8.4 Connection to Other Estimation Methods	113
8.5 Theoretical Properties of MLE	113
8.6 The Role of the Fisher Information	115
8.7 The Role of the Cramer-Rao Lower Bound	115
8.8 2. Intuition Behind the Cramer-Rao Lower Bound	115
8.9 3. Efficiency of Estimators	116
8.10 4. Role in Statistical Inference	116
8.11 5. Practical Implications	116
8.12 6. Examples	117
8.13 Estimating Parameters with MLE	118
8.14 Properties of MLE	119
8.15 Challenges and Limitations of MLE	119
8.16 GPT Dump Conclusion	119
8.17 Claude Dump	120
8.18 Historical Context and Development	120
8.19 Theoretical Foundation	120
8.20 The Maximum Likelihood Principle	120
8.21 Properties of MLE	121
8.22 Applications in Machine Learning	121
8.23 MLE Section	123
8.24 MLE Section	124
8.25 MLE Section	125
8.26 MLE Section	126
8.27 MLE Section	127
8.28 Determining MLE: A Generalized Approach	128
8.29 STEP By STEP Description of MLE (worth it ?)	131
9 Linear Regression	133
9.1 Introduction	135
9.2 Mathematical Formulation	135
9.3 Understanding the Components	135
9.4 Estimation using Ordinary Least Squares (OLS)	135
9.5 Implementation in R	136
9.6 Limitations or Pitfalls of Linear Regression	137
9.7 Impact of Outliers on Goodness of Fit	138
9.8 Equation Summary	141
10 Ordinary Least Squares (OLS)	143
10.1 Linear Model Intro	145

11 R-squared (R²) Score	155
12 Mean Squared Error (MSE)	157
13 Mean Squared Error with L2 Regularization (MSE + L2 Reg)	159
14 Eigenvectors and Eigenvalues	161
15 Entropy	163
16 K-Means Clustering	165
17 Kullback-Leibler (KL) Divergence	167
17.1 Introduction to Kullback-Leibler (KL) Divergence	168
17.2 Example Calculation	168
17.3 Results Table	168
17.4 Observations	169
18 Log Loss	171
19 Support Vector Machine (SVM) Objective	173
20 Singular Value Decomposition (SVD)	175
21 Lagrange Multiplier	177
22 The Human Equation	179
Appendix	181
Matrix Formulation for Linear Separability	181
Epilogue - Guidance and Test Materials	183
Applying Musk Rules to Writing Your Book on the “25 Key Equations in Machine Learning”	183
22.1 TEMP _ R code - Python code - test area	186
22.2 testing Python Integration (temp section - to remove later)	186
22.3 Define Numbers in R	187
22.4 R Code Block	187
22.5 Python Code To Pass *	187
22.6 Python Code Block	187
Apendix II - test area	189
Galley Sheet	195
22.7 GALLEY Sheet info	195
22.8 GALLEY Sheet info - ends	197
References	199

Revision	Date	Author(s)	Description
0.00.0	2024.10.25	DP	Created!
0.01.0	2024.10.25	DP	Thereom boxes changes, layout changes
0.03.0	2024.10.26	DP	Buildout Chapter Frames
0.04.0-0.06.0	2024.11.07	DP	Chapter Builds
0.07.0-0.07.6	2024.11.24	DP	Naive Bayes chapter drafted and updates
0.07.7	2024.11.24	DP	Naive Bayes chapter update - prob space
0.08.0	2024.11.25	DP	Start MLE
0.09.0	2024.11.27	DP	Start Linear Regression
0.10.0	2024.11.27	DP	Continue MLE, start OLS

Prefice and Introduction

This book is a journey through the 25 most important mathematical equations that underpin modern data science and machine learning. The idea is not only to learn about each of these equations but also to deeply understand their applications, background, and practical examples. We will explore each concept visually and instructively, shedding light on how these fundamental equations contribute to building intelligent systems.

Mathematics is the language of the universe, and it is the foundation of the incredible advances we see today in artificial intelligence and machine learning. The purpose of this book is to develop a learning guide that is both informative and visually compelling, bringing out the beauty and utility of these powerful mathematical tools. From optimization algorithms like Gradient Descent to classification techniques like Naive Bayes, and from measures of information like Entropy to clustering algorithms like K-Means, this book aims to provide an accessible yet thorough explanation of how these concepts work and why they are essential.

Each chapter focuses on one equation, starting with an introduction, followed by a detailed description, and ending with examples of how it is used in practice. Whether you're an aspiring data scientist, a seasoned engineer, or just curious about the mathematics that drives intelligent technology, this book will provide you with a solid understanding of these essential tools and how they interact to solve complex problems. Visual aids, illustrative examples, and in-depth explanations will help demystify each topic, making learning both engaging and enjoyable.

Each section will be followed by questions and homework problems. These problems are selected based on their popularity and effectiveness in enhancing understanding. The goal is to provide exercises that solidify the reader's comprehension of each topic. An answer key is provided at the end of the book to facilitate learning and self-assessment.

Many of the examples and problems use Python and R to provide practical insights. Code blocks in Python and R (and in some instances Mathematica) are used throughout to illustrate the concepts discussed. All the example codes, as well as the entire book, are posted on the author's GitHub repository for easy access and reproducibility.

Whenever possible, all examples that originate from a particular source are acknowledged by citation and recognition of the original author. If any citation or attribution is missed, the author kindly requests feedback so that proper corrections can be made.

Some notes on the typesetting framework and tools

With the current advent of ML based LLMs and the availability of typesetting \LaTeX , we combine tools like **Mathpix** (the \LaTeX equation generating snipping tool) and R Markdown driven by **Pandoc** and **knitr** to essentially “GET IT DONE”. Obviously your mileage may vary; we find this tool set perfect for typesetting combined with computational analysis. Why? Well you have nearly full \LaTeX capability with access to R & Python (via **reticulate**) programming code blocks and also the ability to bridge to **MATHEMATICA[®]** if needed. On a powerful laptop, this tool suite (under RStudio) leaves the door open to a wider world of analysis and computation. While RStudio is not perfect, it does allow for an IDE approach which leverages typesetting and computation framework (through R, Python and **MATHEMATICA[®]**) that makes it a worthy Swiss Army Knife. The alternative is also VS Code, which is great, however the mixed mode R Markdown & \LaTeX computational environment is pretty flexible.

The complimentary nod to ML LLM tools

Most of this compilation and summary work would not be possible without the advent of LLMs and new found ability to build concise summaries of difficult mathematical concepts. The time savings is basically a gift which save more grey hairs and time flipping pages of numerous cross references. Actually we want to spend time “learning” and not “flipping and skimming” to find the properly defined statement which is of utility in our endeavor. So for that - we lean on LLMs to be more productive, and help clearly convey underlying concepts. This is not to be lazy, each section must be scrutinized for accuracy and consistency. The speed gains allow for a very large group of examples and highly detailed presentations.

Notes to usage

We give full credit for any original works that are used herein to reach the stated objective. At this point this is not a financial pursuit we do not provide any license related respect to any previous publisher. Go figure. Why would we take such a stand - this is educational open source material. We strongly recommend you take the time to do some background research on any referenced authors.

We truly hope you enjoy the tour, as much as we have enjoyed our journey to write it...

to-ohru iwanami, 2024, somewhere in asia

The 25 (uhh 24)

How did we arrive at the 25 (uhh 24) most important equations? Well, it's partly a matter of educated opinion, partly a rough statistical average of what people talk about most when they're exhilarated by the magic of machine learning, and partly—let's be real, you should just listen to me (fully joking here). These equations are more than just mathematical tools; they represent milestones that have shaped the evolution of artificial intelligence, each one contributing a brick or building block in the formidable wall of progress we've built over the decades.

Take, for instance, Gradient Descent — a cornerstone of optimization, whose principles date back to the early 19th century and the work of Adrien-Marie Legendre and Carl Friedrich Gauss in minimizing residuals. The modern incarnation of Gradient Descent, essential for training deep neural networks, only emerged in the mid-20th century, revitalized by researchers trying to teach machines to “learn.” Fast forward to the 1960s, when Frank Rosenblatt was developing the perceptron, it became evident that iterative optimization methods like Gradient Descent could unlock the potential of early neural networks.

Then we have Bayes’ theorem, after Reverend Thomas Bayes, who developed his idea of conditional probability back in the 18th century. Although it lay relatively dormant for years, it became a key breakthrough in the 1950s when Alan Turing and others applied Bayesian methods to codebreaking during World War II. Today, its cousin—Naive Bayes — is still a powerful tool for classification, especially for natural language processing, allowing us to create chatbots and email spam filters. This dormant theorem found new life, propelling us toward a world of probabilistic understanding in machines.

Singular Value Decomposition (SVD), another member of our elite list, found its fame in the 1990s when it was used for information retrieval in the famous Latent Semantic Analysis (LSA) algorithm—paving the way for how we handle and understand textual data. The underlying mathematics, discovered by Eugenio Beltrami and others in the 19th century, helps us not only reduce dimensionality but also reveal the hidden relationships between concepts in data.

Fast forward to today to the ReLU (Rectified Linear Unit) function. While it seems so simple,

this piece of mathematical elegance emerged as a game-changer for deep learning in the 2010s, thanks to the work of Geoffrey Hinton and his colleagues. Before ReLU, neural networks struggled with the vanishing gradient problem, limiting their ability to learn effectively. By simply transforming negative inputs to zero and retaining positive ones, ReLU helped neural networks go deeper and deeper (pun-ch!), giving birth to the deep learning revolution we know today.

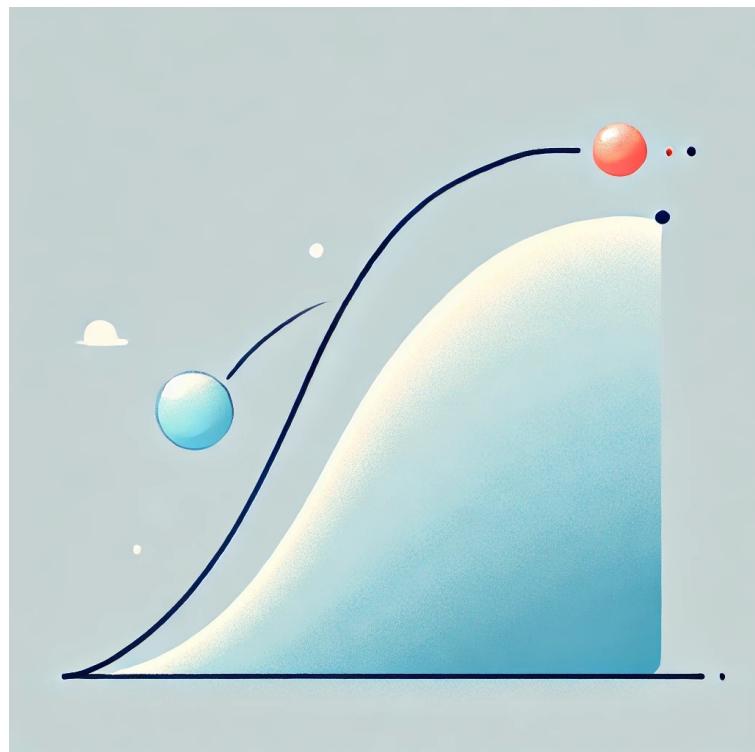
Each of these equations is a testament to the evolution of knowledge, a story of breakthroughs spanning centuries. Together, they form a bridge from the deterministic calculus of Newton and Leibniz to the probabilistic dreams of Bayes and the practical algorithms of today's AI pioneers. Think of them as the greatest hits of mathematics for data scientists—carefully curated, occasionally debated, and ultimately distilled into the essence of what drives machine learning today. This collection captures a human legacy of curiosity and ingenuity, guiding you on your intellectual journey through the ever-growing forest of machine learning—one equation, one insight at a time.

As we reflect on these 1 through 24 equations, and on to 25, one cannot ignore the human spirit that threads through each mathematical symbol and formula. Behind every breakthrough is a mind driven by curiosity, a determination to see beyond the obvious, to turn abstraction into discovery. Each equation, in its own way, captures a moment where human thought transcended barriers—whether it was a problem of optimization, uncertainty, or understanding the nature of learning itself. Yet, even as these 24 pillars of insight stand tall, there remains one more — a keystone that binds them all, born not out of calculation alone but from the essence of what it means to explore, to question, and to create. This final entry goes beyond numbers and symbols, embracing the very source of every theorem, every insight, and every spark of ingenuity. So as you journey onward, remember—there's always one more equation to unveil, one that's been with us all along.

Equation 1

Gradient Descent

$$\theta_{j+1} = \theta_j - \alpha \nabla J(\theta_j)$$



Finding a local minimum with Gradient Descent

Key ML Equation 1: Gradient Descent

$$\theta_{j+1} = \theta_j - \alpha \nabla J(\theta_j) \quad (1.0.1)$$

θ_j	The current value of the parameter vector at iteration j , which represents the current estimate of the model parameters.
θ_{j+1}	The updated value of the parameter vector for the next iteration, which results from applying the gradient step to the current parameters.
α	The learning rate , which is a positive scalar determining the size of the step to take in the direction of the negative gradient.
J	The cost function , which is the function being minimized by adjusting the parameter vector θ . It measures the difference between predicted values and actual values.

Introduction: Gradient Descent is an optimization algorithm used to minimize the cost function by iteratively moving in the direction of the steepest descent as defined by the negative of the gradient.

Description: In machine learning, gradient descent is used to update the parameters of the model, θ , to reduce the difference between the predicted and actual outcomes.

Importance in ML: Gradient Descent is foundational for training machine learning models, particularly in neural networks and linear regression. It helps in finding optimal parameters by iteratively reducing the error, making it crucial for model accuracy.

1.1 What is behind the equation

1.2 Structure of the Parameter Vector θ

The parameter vector, typically denoted as θ , contains all the adjustable parameters or weights of the model that you want to optimize in order to minimize the cost function $J(\theta)$. Depending on the type of model, the structure of θ can vary:

- **Linear Regression:**

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad (1.2.1)$$

In linear regression, θ is often a column vector of weights where each θ_i corresponds to the weight associated with feature x_i . θ_0 is often referred to as the bias term or intercept.

- **Logistic Regression / Neural Networks:** In these models, θ can have more dimensions. For a neural network, the parameter vector could be a collection of weight matrices for different layers, such as:

$$\theta = \{W_1, W_2, \dots, W_L, b_1, b_2, \dots, b_L\} \quad (1.2.2)$$

where W_i and b_i are weights and biases associated with layer i of the network.

- **Deep Learning Models:** The parameter vector is much more complex, often containing multiple matrices and vectors representing weights and biases for each layer in a deep neural network.

In general, θ is a vector that can be expressed as:

$$\theta = (\theta_1, \theta_2, \dots, \theta_n)^T \quad (1.2.3)$$

where n is the number of features or neurons, depending on the model type.

1.3 Typical Value of the Learning Rate α

The learning rate α controls the step size when updating the parameters during gradient descent. Choosing a proper value for α is crucial for the convergence of the algorithm. Here are some general guidelines:

- **Typical Range:** A typical value for the learning rate lies between 0.001 and 0.1. Values outside this range can either lead to a slow convergence or an unstable training process.
- **Considerations:**
 - **Too Small:** If α is too small, gradient descent will take very small steps towards the optimum, resulting in a very slow convergence process.
 - **Too Large:** If α is too large, gradient descent might overshoot the minimum or even diverge, causing the cost function to oscillate or increase.
 - **Adaptive Learning Rates:** Some advanced algorithms like **Adam** use adaptive learning rates which adjust automatically as training progresses.

In practice, tuning the learning rate often involves trial and error or the use of techniques like **learning rate schedules** or **grid search** to find the best value for a specific problem.

1.4 Math Behind the Gradient $\nabla J(\theta_j)$

The term $\nabla J(\theta_j)$ is the **gradient** of the cost function $J(\theta)$ evaluated at θ_j . Let's break it down:

- **Gradient Definition:** The gradient of a function is a vector of partial derivatives with respect to each parameter in θ . In the case of $J(\theta)$, the gradient $\nabla J(\theta_j)$ tells us how much the cost function changes when we make an infinitesimally small change to each component of θ . Mathematically, for a parameter vector $\theta_j = (\theta_1, \theta_2, \dots, \theta_n)^T$:

$$\nabla J(\theta_j) = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \Big|_{\theta_j} \\ \frac{\partial J}{\partial \theta_2} \Big|_{\theta_j} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \Big|_{\theta_j} \end{bmatrix}$$

Each partial derivative $\frac{\partial J}{\partial \theta_i}$ represents the rate of change of the cost function with respect to parameter θ_i .

- **Intuition:** The gradient $\nabla J(\theta_j)$ points in the direction of the **steepest ascent** of the cost function J . In gradient descent, we want to **minimize** $J(\theta)$, so we move in the opposite direction, which is why the update rule is:

$$\theta_{j+1} = \theta_j - \alpha \nabla J(\theta_j)$$

- $\nabla J(\theta_j)$ represents the slope or direction in which $J(\theta)$ increases most quickly.
- By subtracting $\alpha \nabla J(\theta_j)$, we effectively move in the direction of steepest **descent**, hence reducing $J(\theta)$.
- **Computing the Gradient:** In practice, the gradient $\nabla J(\theta_j)$ is computed using **differentiation**. For different cost functions, the gradient takes different forms:
 - For **linear regression** with a **mean squared error (MSE)** cost function, the gradient is relatively simple and involves the residuals (errors) between predictions and actual values.
 - For **neural networks**, computing the gradient involves **backpropagation**, which is a process of applying the chain rule of calculus to calculate the gradients efficiently for each layer.

In short, $\nabla J(\theta_j)$ gives us the necessary information to adjust θ in a way that reduces the error. The learning rate α then determines how big the step should be in this direction.

1.5 Gradient Descent in R

Gradient Descent is a foundational optimization algorithm used to iteratively minimize cost functions. Here, we apply the gradient descent method to find the local minimum of a specific polynomial function. The function used is a quartic polynomial $P(x) = x^4 - 6x^3 + 11x^2 - 6x$. Our goal is to observe how the gradient descent algorithm updates our parameter over successive iterations to converge towards a minimum point of the polynomial.

The following R code demonstrates how to implement gradient descent for this polynomial, including plotting the descent path and function value across iterations. This practical example aims to give you a clear understanding of how gradient descent operates on real functions, using R for illustration. Note that gd_result is the dataframe result fed to ggplot.

Code Listing 1.1: Gradient Descent: Over a Polynomial

```
# Set seed for reproducibility
set.seed(42)

# Define the new polynomial function and its derivative
polynomial_function <- function(x) {
  return(x^4 - 8 * x^3 + 18 * x^2 - 11 * x + 2)
}

# Derivative of the new polynomial function
polynomial_derivative <- function(x) {
  return(4 * x^3 - 24 * x^2 + 36 * x - 11)
}

# Gradient Descent Algorithm
gradient_descent <- function(learning_rate = 0.01, iterations = 1000, start = 3) {
  x <- start # Starting point

  # Store x values and function values for plotting
  x_values <- numeric(iterations)
  function_values <- numeric(iterations)

  for (i in 1:iterations) {
    grad <- polynomial_derivative(x)
    x <- x - learning_rate * grad
    x_values[i] <- x
    function_values[i] <- polynomial_function(x)
  }
  return(data.frame(Iteration = 1:iterations, x_values, function_values))
}

# Run the Gradient Descent with specific parameters
learning_rate <- 0.01
iterations <- 100
start_point <- 3
gd_result <- gradient_descent(learning_rate, iterations, start_point)
```

Code Listing 1.2: Gradient Descent: Objective Function vs. Iteration

```
library(ggplot2)
ggplot(gd_result, aes(x = Iteration, y = function_values)) +
  geom_line(color = "blue", size = 1.2) +
  ggtitle("Gradient Descent on Polynomial Function") +
  xlab("Iteration") +
  ylab("Objective Function Value: f(x)") +
  theme_minimal() +
  theme(
    plot.title = element_text(size = 11) # Adjust the size value as needed
  )
```

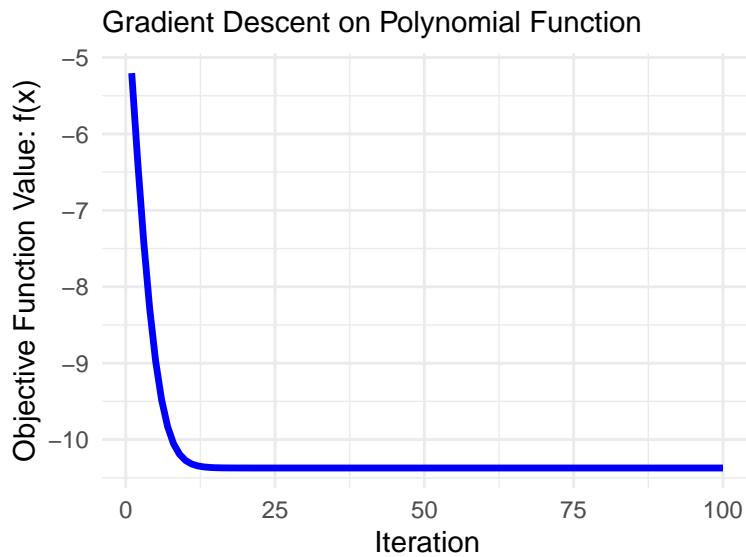


Figure 1.1: Gradient Descent on Polynomial Function - Iteration vs Function Value

Code Listing 1.3: Gradient Descent: Plot the Path of Descent

```
# Additional Plot: Trace the Path of Gradient Descent on the Polynomial
x_range <- seq(min(gd_result$x_values) - 1, max(gd_result$x_values) + 1, length.out = 500)
polynomial_values <- polynomial_function(x_range)

ggplot() +
  geom_line(aes(x = x_range, y = polynomial_values), color = "black", size = 1) +
  geom_point(data = gd_result, aes(x = x_values, y = function_values), color = "red", size = 1.5) +
  ggtitle("Gradient Descent Path on Polynomial Function") +
  xlab("x") +
  ylab("f(x)") +
  theme_minimal() +
  theme(
    plot.title = element_text(size = 11) # Adjust the size value as needed
  )
```

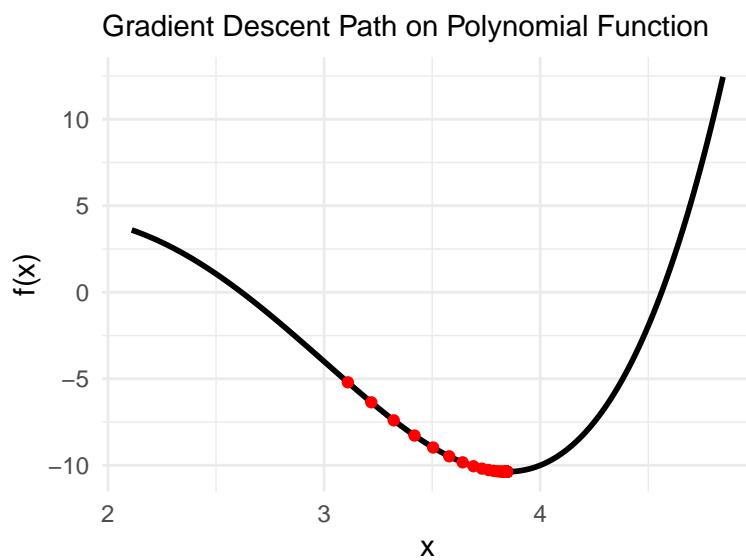


Figure 1.2: Gradient Descent on Polynomial Function - Trace the Path on the Polynomial

1.6 Linking Gradient Descent to Our Polynomial Example

The equation we used to explain gradient descent is:

$$\theta_{j+1} = \theta_j - \alpha \nabla J(\theta_j)$$

1.6.1 Step-by-Step Breakdown

The fundamental idea behind gradient descent is to iteratively adjust the model's parameters, θ , in a direction that reduces the cost function, $J(\theta)$. The parameter update is governed by the **gradient** of the cost function (or loss function), $\nabla J(\theta)$.

In the above equation:

- θ_j represents the parameter(s) at step j , and θ_{j+1} is the updated parameter at step $j + 1$.
- α is the **learning rate**, which controls how large each step is in the descent.
- $\nabla J(\theta_j)$ is the **gradient** of the cost function with respect to the parameters, evaluated at θ_j . It gives the direction of the steepest ascent, and since we want to minimize the function, we move in the opposite direction, hence the negative sign.

1.6.2 Gradient Descent Applied to the Polynomial

In our R code example, we have a **polynomial function** defined as:

$$P(x) = x^4 - 6x^3 + 11x^2 - 6x$$

Our goal is to minimize this polynomial function, which represents our **cost function**, $J(x)$. Here, x plays the role of our parameter, θ , that we want to adjust iteratively using gradient descent.

The **derivative** of the polynomial is:

$$P'(x) = 4x^3 - 18x^2 + 22x - 6$$

In the gradient descent algorithm, this derivative represents the **gradient** of the cost function with respect to our parameter, x . This means that $\nabla J(x) = P'(x)$. The iterative update step becomes:

$$x_{j+1} = x_j - \alpha P'(x_j)$$

Where: - x_j is the current value of the parameter (similar to θ_j). - α is the learning rate that we set in the R code. - $P'(x_j)$ is the gradient of the polynomial at the current point.

1.6.3 Bringing It Full Circle

In our R code, we used **gradient descent** to minimize the polynomial function by starting at an initial point ($x = 3$) and repeatedly updating x using the gradient descent rule:

```
x <- x - learning_rate * grad
```

This corresponds exactly to the equation:

$$x_{j+1} = x_j - \alpha P'(x_j)$$

As each iteration proceeds, x gets closer and closer to a point where the gradient is zero (i.e., a local minimum or a stationary point). In this particular case, the polynomial $P(x)$ has one real minimum, and the gradient descent algorithm helps us converge towards it.

1.6.4 Visual Connection

The **first plot** in our example (function value vs. iteration) shows how the value of the polynomial decreases over time as gradient descent proceeds. The **second plot** (gradient descent path on the polynomial curve) visually shows how x moves along the polynomial curve, getting closer and closer to the minimum.

This linkage between the theoretical equation of gradient descent and the practical implementation of minimizing the polynomial demonstrates how gradient descent serves as a universal tool to solve optimization problems, whether in machine learning or even in simple polynomial functions like the one used in our example.

Gradient descent works by always taking steps in the direction that most rapidly reduces the cost function—allowing us to find optimal parameters, minimize errors, or reach local minima. It is an essential part of training many machine learning models, and the concept remains consistent regardless of the specific problem context.

1.7 Additional ideas to explore

Stochastic Gradient Descent (SGD), Mini-Batch Gradient Descent, or concepts like Momentum and Learning Rate Scheduling.

We could also explore more code examples, such as:

Adding stopping criteria to our gradient descent R code, like a tolerance for change in cost function.

Visualizing the convergence of the gradient descent with multiple starting points.

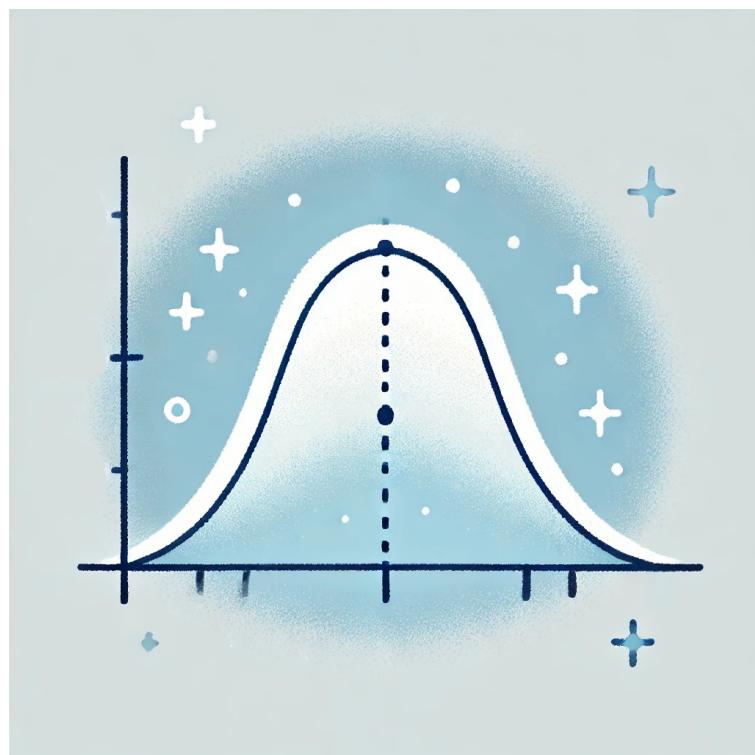
Implementing different cost functions and comparing their optimization trajectories. Let me know which direction you'd like to explore, and I'll get started!

Equation 2

Normal Distribution

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$f(x | \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Normal Distribution, Variance and Mean

Key ML Equation 2: Normal Distribution

$$f(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (2.0.1)$$

x	The random variable for which the probability density function is being calculated. It represents the value within the distribution.
μ	The mean of the distribution, which represents the center or "average" value around which the data clusters.
σ	The standard deviation of the distribution, indicating how spread out the values are around the mean. A larger σ means more spread, while a smaller σ means values are more tightly clustered.
σ^2	The variance of the distribution, which is the square of the standard deviation. It provides a measure of the dispersion of the distribution.
$f(x \mu, \sigma^2)$	The probability density function (PDF) for the normal distribution, which provides the likelihood of x occurring given the parameters μ and σ^2 .
\exp	The exponential function , which ensures that the PDF value falls off symmetrically from the mean μ . It plays a key role in modeling the bell-shaped curve characteristic of the normal distribution.

Introduction: The normal distribution, also called the Gaussian distribution, is a probability distribution that is symmetric about the mean.

Description: It represents how data tends to cluster around a central point. The parameters μ and σ^2 represent the mean and variance, respectively.

Importance in ML: The normal distribution is used in many ML algorithms, especially in probabilistic models and hypothesis testing. Assumptions of normality often simplify the mathematics of learning models and are vital in Bayesian networks.

2.1 What is Behind the Equation

The normal distribution equation represents a probability density function (PDF) that models how data points are distributed around a central value (the mean, μ). This equation is essential in statistics because it describes a common pattern found in natural phenomena, from heights of people to errors in measurements. The bell-shaped curve is a striking visual, representing how values tend to cluster around the mean, with fewer observations occurring as we move further from this center. The beauty of the normal distribution lies in its symmetry and the way it characterizes many real-world datasets, making it a cornerstone in both statistics and machine learning.

2.2 The Players

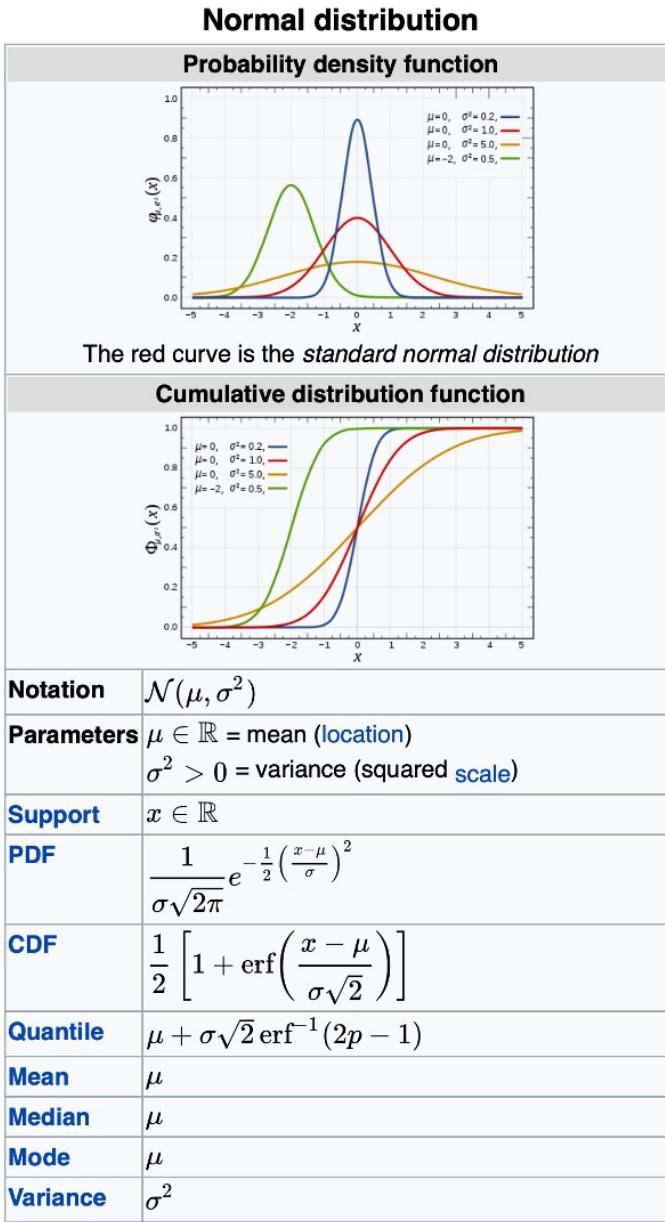


Figure 2.1: Summary of Normal Distribution with PDF and CDF



Carl Friedrich Gauss
discovered the normal
distribution in 1809 as a way to
rationalize the **method of least
squares**.

Figure 2.2: Carl Friedrich Gauss



Pierre-Simon Laplace
proved the **central limit
theorem** in 1810, consolidating
the importance of the normal
distribution in statistics.

Figure 2.3: Pierre-Simon Laplace

1. **Normal Distribution:** [Normal Distribution on Wikipedia]
2. **Carl Friedrich Gauss:** [Carl Friedrich Gauss on Wikipedia]
3. **Pierre-Simon Laplace:** [Pierre-Simon Laplace on Wikipedia]

2.3 Historical Context

The normal distribution, often referred to as the Gaussian distribution, has a rich history rooted in the development of probability theory and statistics. The distribution is named after the German mathematician Carl Friedrich Gauss, who used it extensively in his work on astronomy and measurement errors in the early 19th century. Gauss formalized the idea that errors in measurements tend to follow a symmetric pattern around the true value, leading to the bell-shaped curve we now associate with the normal distribution.

However, the concept of the normal distribution predates Gauss and can be traced back to the work of Abraham de Moivre, an 18th-century French mathematician. De Moivre first derived the normal distribution as an approximation to the binomial distribution when the number of trials becomes very large. His work laid the foundation for what would later be formalized and popularized by Gauss.

The normal distribution became particularly significant due to the Central Limit Theorem, which states that the sum of many independent random variables, regardless of their original distribution, tends to follow a normal distribution. This theorem, proven by mathematicians such as Pierre-Simon Laplace, helped establish the normal distribution as a fundamental tool in statistics and the natural sciences. Today, it is widely used not only because of its mathematical properties but also because it naturally arises in numerous real-world situations, making it one of the most important and recognizable distributions in probability and statistics.

2.4 Relative Standard Deviation (RSD)

In the context of the normal distribution, the spread of data points is characterized by the standard deviation (σ). A useful concept derived from this is the **Relative Standard Deviation (RSD)**, which is expressed as a percentage of the mean:

$$\text{RSD} = \left(\frac{\sigma}{\mu} \right) \times 100$$

Typical RSD Values: The RSD provides insight into how variable the data is relative to its mean. A low RSD (e.g., below 10%) indicates that the data points are closely clustered around the mean, whereas a higher RSD suggests more dispersion. In many real-world datasets, RSDs ranging between 5% and 20% are common, depending on the type of measurement and its inherent variability. RSD helps give a standardized view of spread that is independent of the scale of the data, which is particularly useful when comparing the variability between datasets.

2.5 Understanding the Notation $f(x|\mu, \sigma^2)$

The notation $f(x|\mu, \sigma^2)$ represents the **probability density function** of a normal distribution given the parameters μ (mean) and σ^2 (variance). This notation indicates a **conditional relationship**, where the probability density function $f(x)$ is dependent on the parameters μ and σ^2 .

In other words, μ and σ^2 define the specific characteristics of the distribution (where it is centered and how spread out it is), while x represents the variable for which the probability is being calculated. This notation highlights that the distribution is characterized by these parameters, and the output $f(x)$ tells us how likely it is to observe a particular value of x under that distribution.

2.6 The Nature of the Exponential Function in Normal Distribution

A key component of the normal distribution equation is the **exponential function**:

$$\exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

This part of the function gives the normal distribution its famous **bell shape**. The negative squared term in the exponent ensures that values closer to the mean (μ) have higher probabilities, while values further away have exponentially decreasing probabilities. This shape is what makes the normal distribution symmetric, with a single peak at the mean.

The function $\exp(-x^2)$ falls off rapidly as $|x|$ increases, which results in a smooth, continuous decline from the peak at the mean. This property is crucial because it reflects how, in many natural phenomena, values tend to cluster around an average, with extreme values being much less common. This elegant behavior is what makes the normal distribution so special and why it is so frequently used in statistical modeling.

2.7 Non-existence of an Analytical Anti-Derivative

Interestingly, the **normal distribution function does not have an anti-derivative** that can be expressed in closed form. This means that the area under the curve (which represents the cumulative probability) cannot be solved using traditional analytic integration techniques. Instead, it is computed numerically or looked up using statistical tables.

The integral of the normal distribution is given by:

$$\int_{-\infty}^{\infty} f(x|\mu, \sigma^2) dx = 1$$

However, there is no elementary function that represents this integral. This is why the **error function (erf)** is introduced in mathematics to help approximate these values:

$$\int e^{-x^2} dx = \frac{\sqrt{\pi}}{2} \operatorname{erf}(x) + C$$

The lack of an analytical solution means that in practice, probabilities for a normal distribution are often calculated using **numerical methods** or **precomputed tables**. This characteristic of the normal distribution makes it an interesting function from a mathematical standpoint, as it combines simplicity of form with complexity in integration.

2.8 Central Limit Theorem

The Central Limit Theorem (CLT) is a fundamental concept in probability theory that explains why the normal distribution is so prevalent in nature and in machine learning. The CLT states that when independent random variables are added together, their properly normalized sum tends to form a normal distribution, regardless of the original distribution of the variables. This remarkable property applies as long as the number of variables is sufficiently large and they have finite variances.

Mathematically, if we have a set of random variables X_1, X_2, \dots, X_n , each with mean μ_i and variance σ_i^2 , the sum or average of these variables will tend towards a normal distribution as n becomes large. The formal expression is:

where \rightarrow denotes convergence in distribution, and Z represents the standard normal distribution.

The CLT helps explain why the normal distribution is observed so often in practice. Many complex processes can be thought of as the sum of many small, independent effects. Whether it's measurement errors, human heights, or even financial market fluctuations, these effects combine to form a distribution that is approximately normal.

2.8.1 Importance in Machine Learning

The Central Limit Theorem has significant implications in machine learning. It provides the foundation for many statistical techniques and justifies the assumption of normality in models. For example:

Modeling Errors: In regression analysis, the residuals (errors) are often assumed to be normally distributed. This assumption allows us to derive confidence intervals and perform hypothesis

testing.

Feature Engineering: When aggregating data, such as calculating the mean of multiple features or observations, the resulting values tend to be normally distributed, making it easier to apply techniques that assume normality.

Sampling Distributions: The CLT allows us to make inferences about population parameters from sample data. Many machine learning algorithms rely on sampling and estimation, and the CLT ensures that the distribution of the sample mean approximates normality, which simplifies analysis and interpretation.

The CLT thus serves as a bridge between randomness and order, allowing us to apply the powerful tools of the normal distribution even in cases where the underlying data might not be normally distributed. It is this ability to generalize and predict outcomes that makes the normal distribution so central to both statistics and machine learning.

2.9 Applications in Machine Learning

The normal distribution plays a vital role in numerous machine learning algorithms and concepts. Its presence is seen in everything from model assumptions to data transformations. Below are some of the key areas where the normal distribution is commonly applied:

1. Gaussian Naive Bayes

Gaussian Naive Bayes is a classification algorithm based on Bayes' theorem. It assumes that the features follow a normal distribution, which allows the model to calculate probabilities efficiently. The assumption of normality simplifies the calculations, enabling rapid classification even with high-dimensional data. This assumption works well for many real-world datasets, making Gaussian Naive Bayes a popular choice for problems like spam detection and document classification.

2. Linear Regression Error Terms

In linear regression, the error terms (residuals) are often assumed to be normally distributed. This assumption allows us to make statistical inferences about the parameters of the regression model, such as constructing confidence intervals and conducting hypothesis tests. If the residuals are approximately normal, we can apply powerful statistical tools to evaluate model fit and make predictions.

3. Weight Initialization in Neural Networks

When training neural networks, the weights are often initialized using a normal distribution. For example, in the Xavier initialization method, weights are drawn from a normal distribution with a mean of zero and a variance that depends on the number of input and output nodes.

2.10. SUGGESTED ADDITIONAL CONTENT FOR THE ~~CHAPTER~~ Learning Reviews

This helps ensure that the neurons start with a diverse range of values, which prevents issues such as all neurons producing the same output. Proper initialization is crucial for efficient training, and using the normal distribution helps keep the gradients within a reasonable range during backpropagation.

4. Generative Models

The normal distribution is also used in generative models, such as Gaussian Mixture Models (GMMs), which assume that the data is generated from a mixture of several Gaussian distributions. GMMs are widely used in clustering problems, where they help to model the underlying distribution of the data and assign probabilities to different clusters.

5. Feature Scaling and Data Transformation

In many machine learning algorithms, it is beneficial for features to follow a normal distribution. Methods such as StandardScaler in scikit-learn standardize features by removing the mean and scaling to unit variance, resulting in a distribution with a mean of 0 and a standard deviation of 1. This process is especially important for algorithms that are sensitive to the scale of the input features, such as support vector machines (SVMs) and gradient descent optimization.

Importance of the Normal Distribution in Machine Learning

The normal distribution is not only a convenient assumption but also a useful tool in machine learning. Its prevalence in real-world phenomena and its mathematical properties make it indispensable for building, evaluating, and optimizing models. Whether it's in simplifying calculations through Gaussian Naive Bayes, making statistical inferences in linear regression, initializing neural networks effectively, or clustering data in GMMs, the normal distribution is at the core of numerous machine learning practices. By understanding and leveraging the normal distribution, practitioners can better model uncertainties and improve the robustness of their machine learning solutions.

2.10 Suggested Additional Content for the Chapter

- 1. Historical Context:** Add a brief history of the normal distribution, perhaps mentioning Carl Friedrich Gauss, who contributed to its development, and why it is often called the Gaussian distribution.
- 2. Visualization:** Add a visualization of the normal distribution with varying means and standard deviations to help illustrate how changes in μ and σ affect the shape of the curve.

Equation 3

Z-Score

$$z = \frac{x - \mu}{\sigma}$$



Take a portion of the probability with the Z-Score

Key ML Equation 3: Z-Score

$$z = \frac{x - \mu}{\sigma} \quad (3.0.1)$$

- x The **random variable**, representing the value within the distribution that is being standardized.
- μ The **mean** of the distribution, which is the average value and the point around which data tends to cluster.
- σ The **standard deviation** of the distribution, indicating the spread or dispersion of values around the mean. A larger σ suggests more variability.
- z The **Z-Score**, which represents how many standard deviations a particular x value is from the mean μ . It is a measure of relative position within the distribution.

Introduction: The Z-score represents the number of standard deviations a data point is from the mean.

Description: It is used to standardize data points within a dataset, making comparisons between different distributions possible.

Importance in ML: Z-scores are crucial in feature scaling and normalization, allowing different features to be compared and helping gradient-based algorithms converge faster by ensuring all features have a similar scale.

3.1 What is Behind the Equation

The Z-score, also known as the standard score, is a measure that describes the position of a value relative to the mean of a dataset, in units of the standard deviation. The Z-score calculation is particularly useful in the context of the normal distribution, allowing us to determine how far a particular value x lies from the mean μ when expressed in terms of the distribution's standard deviation σ . The Z-score is expressed by the formula:

$$z = \frac{x - \mu}{\sigma}$$

This metric is instrumental in transforming individual data points into a universal scale, where positive Z-scores indicate values above the mean, and negative Z-scores indicate values below the mean. The Z-score effectively normalizes different datasets, making comparisons straightforward.

3.2 Applications in Machine Learning

In machine learning, Z-scores are employed for several purposes, ranging from outlier detection to feature scaling.

- **Outlier Detection:** Z-scores can help identify outliers, as values with extremely high or low Z-scores typically indicate data points that lie far from the distribution's average behavior. These outliers might signify anomalies or valuable insights that need closer examination.
- **Feature Scaling:** Z-score normalization, also known as standardization, is a common feature scaling method. By transforming features using Z-scores, they are rescaled to have a mean of 0 and a standard deviation of 1. This type of normalization is especially valuable in algorithms like logistic regression, k-means clustering, and principal component analysis (PCA), where features with different ranges can negatively impact model performance.
- **Standard Normal Table Applications:** The Z-score is also used in combination with the standard normal distribution table to compute probabilities and p-values for hypothesis testing. In machine learning, this is often relevant in model evaluation and statistical testing.

3.3 Typical Range for Z-Scores

The typical range for Z-scores in a standard normal distribution is between -3 and 3 . Values beyond this range are considered rare and are often associated with outliers.

- $|Z| > 3$: Typically considered outliers. Values beyond three standard deviations are quite unusual in a normal distribution.
- $|Z| < 1$: Most values are likely within this range and close to the mean.
- $1 \leq |Z| < 2$: These values are still common, although they are further away from the mean.
- $|Z| > 2$: Values in this range start to become unusual, indicating either naturally rare observations or data quality issues.

The calculation of Z-score can be a valuable indicator of how far a data point deviates from what is expected, helping differentiate between normal variance and potential outliers.

3.4 Mathematical Notation and Concept

The equation for the Z-score is represented as follows:

$$z = \frac{x - \mu}{\sigma}$$

- x : Represents the actual observed value within the dataset.
- μ : The mean of the dataset, which represents the central tendency around which the dataset is distributed.
- σ : The standard deviation, which gives insight into the degree of variability within the data points. It tells us how spread out the values are around the mean.

The Z-score essentially tells us how many standard deviations away from the mean a given value x is, and whether it is to the left (negative z) or right (positive z) of the mean.

3.5 Mathematical Insights: Area Under the Normal Curve

One of the key properties of the Z-score is how it allows for the calculation of probabilities from the normal distribution. By converting a value to its Z-score, we can use the standard normal distribution (which has a mean of 0 and a standard deviation of 1) to determine the probability that a value is less than or greater than a given point.

For example, calculating the probability $P(Z \leq z)$ involves integrating the probability density

function of the normal distribution up to the given Z-score value:

$$P(Z \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

Since this integral has no closed-form solution, it is typically evaluated using numerical methods, standard normal distribution tables, or software. This integral defines the cumulative distribution function (CDF) of the normal distribution, which is essential in assessing cumulative probabilities for standard normal variables.

In practice, machine learning models often use pre-calculated tables or libraries to compute these probabilities for performance metrics, hypothesis testing, or evaluating the significance of model parameters.

3.6 Example: Z-Score in R

The Z-score calculation can be easily implemented in R to standardize a dataset or to identify outliers. Below is a simple example to illustrate how we can compute Z-scores for a given dataset and visualize the data distribution using a histogram:

```
# Load required libraries
library(ggplot2)

# Define the parameters for the normal distribution
mean_value <- 0
std_dev <- 1
z_score <- 1.25

# Define the cumulative distribution function (CDF) to calculate probability
p_value <- pnorm(z_score, mean = mean_value, sd = std_dev)

# Generate data for plotting the normal distribution curve
x_values <- seq(-4, 4, length.out = 1000)
y_values <- dnorm(x_values, mean = mean_value, sd = std_dev)

data <- data.frame(x = x_values, y = y_values)

# Plot the normal distribution curve with shaded area under the curve
plot <- ggplot(data, aes(x = x, y = y)) +
  geom_line(color = "black", size = 1) +
  geom_area(data = subset(data, x <= z_score), aes(x = x, y = y), fill = "steelblue", alpha = 0.5) +
  geom_vline(xintercept = z_score, color = "red", linetype = "dashed", size = 1) +
  ggtitle("Normal Distribution with Z-Score of 1.25") +
  xlab("Z") +
  ylab("Density") +
  theme_minimal()

# Print the plot
print(plot)
```

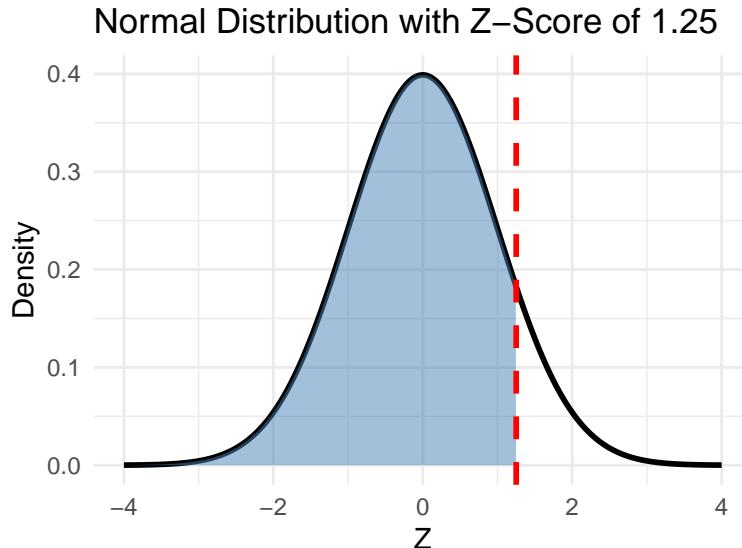


Figure 3.1: Default caption for all figures

```
# Print out the results
cat("Z-Score:", z_score, "\n")

## Z-Score: 1.25
cat("Probability (P-value) for Z <=", z_score, ":", round(p_value * 100, 2), "%\n")

## Probability (P-value) for Z <= 1.25 : 89.44 %
```

In this code block:

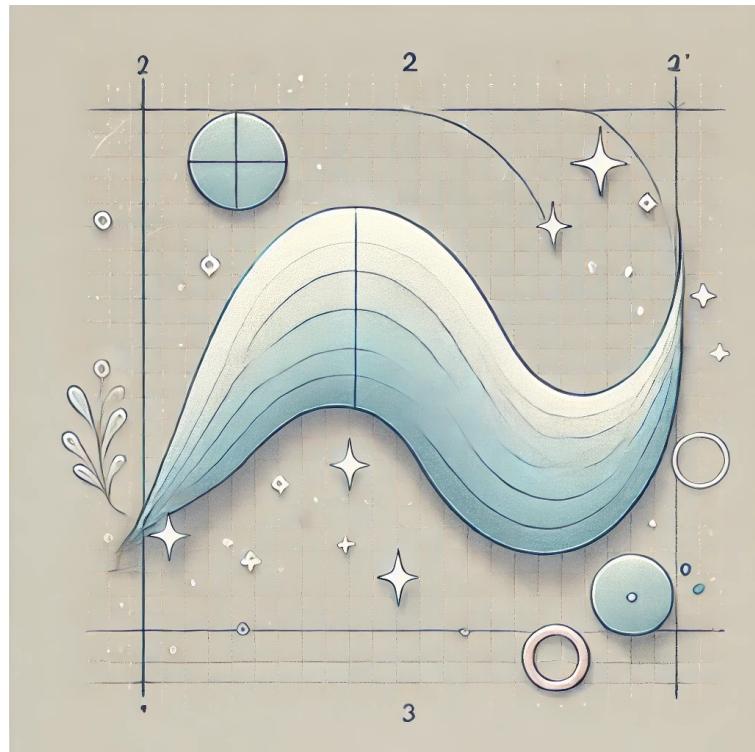
- * We first generate a dataset of 100 normally distributed random numbers with a mean of 50 and a standard deviation of 10.
- * We calculate the Z-scores for each value by subtracting the mean and dividing by the standard deviation of the dataset.
- * Finally, we plot the distribution of the Z-scores using a histogram to visualize their spread.

This example showcases how the Z-score helps normalize data to a common scale, making it easier to compare different values and identify outliers in the dataset.

Equation 4

Sigmoid Function

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



A smooth on or off transition

Key ML Equation 4: Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

- z The **input variable**, which is the linear combination of features and weights in a model. It represents the score calculated before applying the sigmoid transformation and can take any real value.
- $\sigma(z)$ The **sigmoid output**, which is the transformed value of z , mapped between 0 and 1. This output represents the probability that the input belongs to a particular class (e.g., 1 for positive, 0 for negative in binary classification).
- e The mathematical constant, approximately equal to 2.718, which forms the base of the natural logarithm. The exponential term e^{-z} causes the sigmoid function to asymptotically approach 0 for large negative z and 1 for large positive z .

Introduction: The sigmoid function is an activation function that outputs values between 0 and 1.

Description: It maps any input value into a range between 0 and 1, making it suitable for probability-related tasks.

Importance in ML: The sigmoid function is widely used in logistic regression and as an activation function in neural networks. It helps in modeling binary classification problems and introduces non-linearity into neural models.

4.1 What is Behind the Equation

The sigmoid function, denoted as $\sigma(x)$, is an activation function widely used in machine learning and data science. It maps any real-valued number to a value between 0 and 1, making it ideal for scenarios where probability-related interpretations are essential, especially binary classification problems.

This equation transforms an input x into an output within the range of 0 to 1, effectively normalizing the input space.

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x} = 1 - \sigma(-x) \quad (4.1.1)$$

This useful identity allows for simplifications in various mathematical expressions and facilitates gradient-based optimization in machine learning models. Note that the function is differentiable, which is crucial for backpropagation in neural networks.

In many fields, particularly in artificial neural networks, the term sigmoid function usually refers to the logistic function, a specific type of S-shaped, or “sigmoid,” curve. While other functions like the Gompertz curve or the ogee curve may appear similar because they share the same general S-shape, these are distinct mathematical functions with unique properties and applications.

A sigmoid function is characterized by its shape: it begins at a low value, then increases rapidly in the middle, and finally levels off at a high value. The logistic function, which is a common example, maps any real number (from negative to positive infinity) to an output between 0 and 1. This makes it especially useful when working with probabilities or values that need to be constrained within this range.

Other variations of sigmoid functions, like the hyperbolic tangent function (\tanh \tanh), also share the S-shaped curve but differ in their output ranges. The hyperbolic tangent function, for example, produces outputs between -1 and 1, which may be more suitable in cases where negative values are meaningful (such as for representing balanced or symmetric outcomes).

4.2 General Overview of Applications of Sigmoid Functions

- Artificial Neural Networks: Sigmoid functions are widely used as activation functions, which means they determine whether a neuron in a neural network should be “activated” (i.e., pass on its signal). In particular, the logistic function is popular because it maps values to a range between 0 and 1, simulating the all-or-nothing activation seen in

biological neurons. This behavior allows neural networks to make decisions in a gradual, continuous manner.

- Cumulative Distribution Functions in Statistics: The sigmoid function is also useful in statistics, where it can represent the cumulative distribution function (CDF) of a probability distribution. In this context, the function's smooth, continuous increase models the accumulation of probability from one extreme to another, representing how likely it is to observe values up to a certain point.
- Probabilistic Interpretation: When the output of the logistic sigmoid is interpreted as a probability (from 0 to 1), it enables the function to model binary classification problems, where there are only two possible outcomes (such as yes/no or true/false). By setting a threshold (e.g., 0.5), predictions can be made based on the probability that an input belongs to a certain class.
- Invertibility: The logistic sigmoid function is invertible, meaning it has an inverse function called the logit function. The logit function transforms probabilities (between 0 and 1) back into raw values from all real numbers, which can be useful in logistic regression and other statistical models that seek to model relationships in probability terms.

In summary, sigmoid functions, particularly the logistic function, are versatile tools in both machine learning and statistics due to their S-shape, bounded output range, and smooth transition, which makes them ideal for modeling gradual changes in various contexts.

4.3 Application to Machine Learning: Sigmoid Function in Binary Classification

In binary classification tasks, the sigmoid function is particularly valuable as it can be used to estimate the probability that a given input belongs to a certain class. Logistic regression, for example, uses the sigmoid function to model the probability of a binary outcome, with predictions being based on whether the probability exceeds a specified threshold (e.g., 0.5). By mapping outputs to a range between 0 and 1, the sigmoid function enables interpretable probability scores and smooth gradients for optimization.

4.4 How Sigmoid Works in Logistic Regression

In logistic regression, the sigmoid function serves as the “activation” that transforms the output of a linear model into a probability.

1. **Linear Combination of Features:** First, we calculate a linear combination of the input features. For example, in a logistic regression model with two features x_1 and x_2 , the model calculates:

$$z = w_1x_1 + w_2x_2 + b \quad (4.4.1)$$

where w_1 and w_2 are weights learned from the data, and b is the bias term. This value z represents a “score” that is not yet bounded between 0 and 1.

2. **Applying the Sigmoid Transformation:** Next, we pass this linear combination z through the sigmoid function, which transforms it into a probability:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (4.4.2)$$

This output $\sigma(z)$ now represents the estimated probability that the input belongs to the positive class (e.g., “spam” or “yes”).

3. **Thresholding for Classification:** Once we have the probability, we can decide on a threshold to classify the input. For example, if we set a threshold of 0.5, we would classify the input as belonging to the positive class if $\sigma(z) \geq 0.5$, and to the negative class otherwise. This threshold is often set based on the problem’s specific requirements and the balance of classes.

4.5 Example Walkthrough

Imagine you’re building a model to detect spam emails based on certain features, like the presence of certain words or the frequency of punctuation marks. Here’s how the process might work step-by-step:

1. **Input Features:** Let’s say we have two features:
 - x_1 : The frequency of the word “free.”
 - x_2 : The number of exclamation marks in the email.
2. **Model Weights:** After training on a dataset, the model might assign the following weights:
 - $w_1 = 1.2$ (suggesting “free” is an indicator of spam)

- $w_2 = 0.8$ (suggesting a high number of exclamation marks also correlates with spam)

Let's assume the bias term b is -1.5.

3. **Calculate z :** For an email where $x_1 = 2$ and $x_2 = 3$ (meaning “free” appears twice and there are three exclamation marks), we calculate:

$$z = (1.2 \times 2) + (0.8 \times 3) - 1.5 = 2.4 + 2.4 - 1.5 = 3.3 \quad (4.5.1)$$

4. **Apply Sigmoid to Get Probability:** Now, we pass z through the sigmoid function:

$$\sigma(z) = \frac{1}{1 + e^{-3.3}} \approx 0.964 \quad (4.5.2)$$

This probability, 0.964, indicates a high likelihood that the email is spam.

5. **Classification Decision:** If our threshold is 0.5, we would classify this email as “spam” since $0.964 > 0.5$. This probability-based approach allows us to adjust the threshold based on how conservative or lenient we want our classification to be.

4.6 Why Sigmoid is Useful for Optimization

The sigmoid function also has smooth gradients, meaning that even small changes in z lead to gradual changes in the output probability $\sigma(z)$. This smoothness is critical for optimization because it allows gradient-based algorithms, such as gradient descent, to make small, continuous updates to model parameters.

To demonstrate the smooth gradient of the sigmoid function mathematically, we can derive its derivative, which shows how the output changes with respect to changes in the input z . This derivative will show that the sigmoid function has a continuous and smooth gradient, making it suitable for gradient-based optimization methods.

4.6.1 Smooth Gradient of the Sigmoid Function

To understand why the sigmoid function has smooth gradients, let's derive its first derivative. We can calculate the derivative of $\sigma(z)$ with respect to z , which will tell us how sensitive the output $\sigma(z)$ is to small changes in z .

4.7 Deriving the Derivative of the Sigmoid Function

- Start with the sigmoid function:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (4.7.1)$$

- To find $\frac{d\sigma}{dz}$, we apply the chain rule. Let's rewrite the function in a form that will make it easier to differentiate:

$$\sigma(z) = (1 + e^{-z})^{-1} \quad (4.7.2)$$

- Now, differentiate with respect to z :

$$\frac{d\sigma}{dz} = -(1 + e^{-z})^{-2} \cdot (-e^{-z}) \quad (4.7.3)$$

- Simplify the expression:

$$\frac{d\sigma}{dz} = \frac{e^{-z}}{(1 + e^{-z})^2} \quad (4.7.4)$$

- Notice that we can rewrite e^{-z} in terms of $\sigma(z)$ because:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \Rightarrow e^{-z} = \frac{1 - \sigma(z)}{\sigma(z)} \quad (4.7.5)$$

- Substitute this back to express $\frac{d\sigma}{dz}$ in terms of $\sigma(z)$:

$$\frac{d\sigma}{dz} = \sigma(z)(1 - \sigma(z)) \quad (4.7.6)$$

4.7.1 Interpretation of the Derivative

The derivative of the sigmoid function, $\frac{d\sigma}{dz} = \sigma(z)(1 - \sigma(z))$, shows us that:

- The gradient is always positive and smooth for all values of z .
- This derivative is largest around $z = 0$, where the sigmoid curve is steepest, and it becomes smaller as z moves towards positive or negative infinity, where the curve flattens out.

This smooth gradient is crucial for gradient descent, as it allows the algorithm to make gradual updates to the weights. Because the derivative never suddenly changes, the model can adjust parameters smoothly without sudden jumps, aiding in stable and effective learning.

Thus, the sigmoid's smooth gradient makes it well-suited for optimization in machine learning.

4.8 Typical Ranges or Values for the Sigmoid Function

The sigmoid function is bounded by 0 and 1. At $x = 0$, $\sigma(x)$ returns 0.5. As $x \rightarrow \infty$, $\sigma(x) \rightarrow 1$; and as $x \rightarrow -\infty$, $\sigma(x) \rightarrow 0$. The sigmoid function's characteristic "S"-shaped curve makes it particularly suited for gradual transitions between classes in classification tasks.

4.9 Important Mathematical Identities

Some useful identities associated with the sigmoid function include:

1. Derivative of the sigmoid:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (4.9.1)$$

2. Relationship with the hyperbolic tangent:

$$\sigma(x) = \frac{1 + \tanh(\frac{x}{2})}{2} \quad (4.9.2)$$

These identities assist in understanding the behavior of the sigmoid function in neural networks and other machine learning algorithms.

4.10 Comparative Functions

Several functions exhibit similar “S”-shaped or smooth transition behaviors. Here are some of them:

$$\begin{aligned} -\operatorname{erf}\left(\frac{\sqrt{\pi}}{2}x\right) &= -\frac{x}{\sqrt{1+x^2}} \\ -\tanh(x) &= -\frac{2}{\pi} \arctan\left(\frac{\pi}{2}x\right) \\ -\frac{2}{\pi} \operatorname{gd}\left(\frac{\pi}{2}x\right) &= -\frac{x}{1+|x|} \end{aligned}$$

In addition to the sigmoid function, there are several other functions that exhibit smooth, “S”-shaped curves and are used for various applications in machine learning and mathematics. Below are six of these comparative functions, along with a brief description of their unique properties.

4.10.1 Error Function:

The error function, denoted as $\operatorname{erf}(x)$, is often used in probability and statistics, particularly in relation to the normal distribution. It approximates the integral of the Gaussian function, and its shape closely resembles that of the sigmoid. The error function has applications in fields such as signal processing and heat diffusion.

$$-\operatorname{erf}\left(\frac{\sqrt{\pi}}{2}x\right) \tag{4.10.1}$$

4.10.2 Reciprocal Square Root:

This function smooths inputs similarly to the sigmoid but has a different asymptotic behavior. Unlike the sigmoid, which approaches 1 and 0, this function approaches -1 and 1 as x tends to ∞ and $-\infty$, respectively. It is sometimes used in neural networks as an alternative activation function.

$$-\frac{x}{\sqrt{1+x^2}} \tag{4.10.2}$$

4.10.3 Hyperbolic Tangent:

The hyperbolic tangent function, $\tanh(x)$, is a popular alternative to the sigmoid in neural networks because its range is from -1 to 1, allowing for more balanced gradient flows during backpropagation. It provides output symmetry, which can sometimes lead to faster training.

$$-\tanh(x) \quad (4.10.3)$$

4.10.4 Arctangent:

This function maps inputs to a range between -1 and 1 similarly to the hyperbolic tangent but has a slightly different curve. The arctangent function is sometimes used in machine learning models when smaller gradients or more gradual transitions are preferred.

$$\frac{2}{\pi} \arctan\left(\frac{\pi}{2}x\right) \quad (4.10.4)$$

4.10.5 Gudermannian Function:

The Gudermannian function, $\text{gd}(x)$, provides a link between circular and hyperbolic functions without involving complex numbers. This function has applications in geometry and physics, particularly in mapping problems.

$$\frac{2}{\pi} \text{gd}\left(\frac{\pi}{2}x\right) \quad (4.10.5)$$

4.10.6 Inverse Absolute:

This function smoothly scales inputs, but its range only extends from -1 to 1. Its asymptotic properties are particularly useful when a smoother transition around zero is preferred.

$$\frac{x}{1 + |x|} \quad (4.10.6)$$

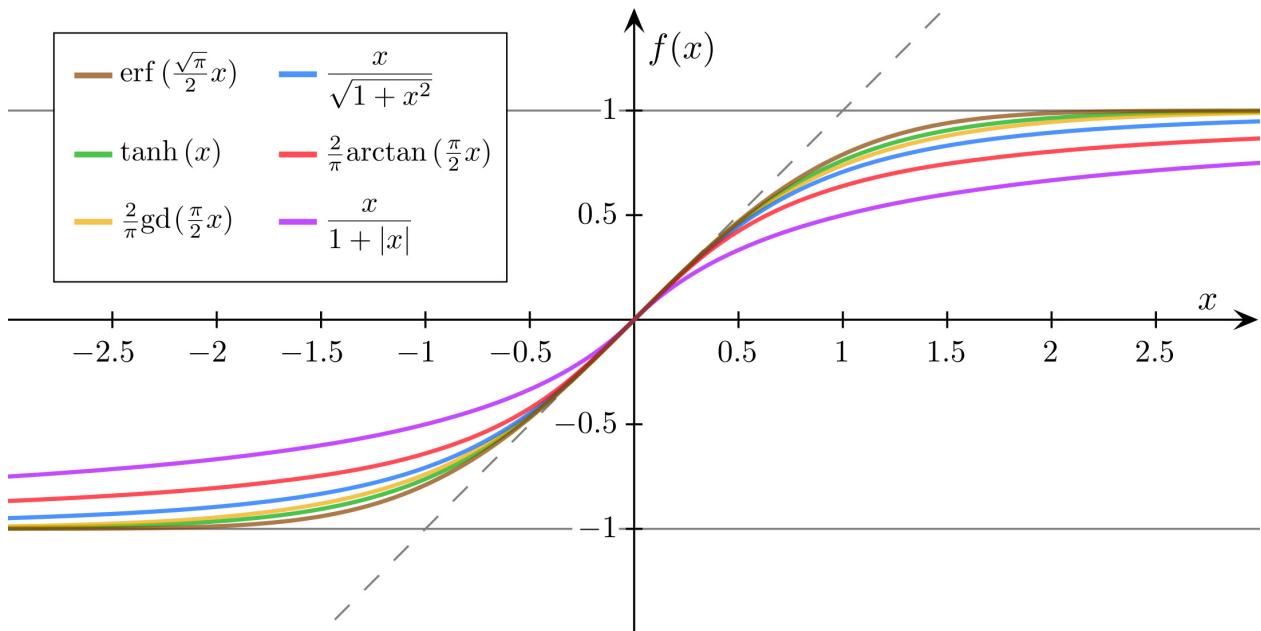


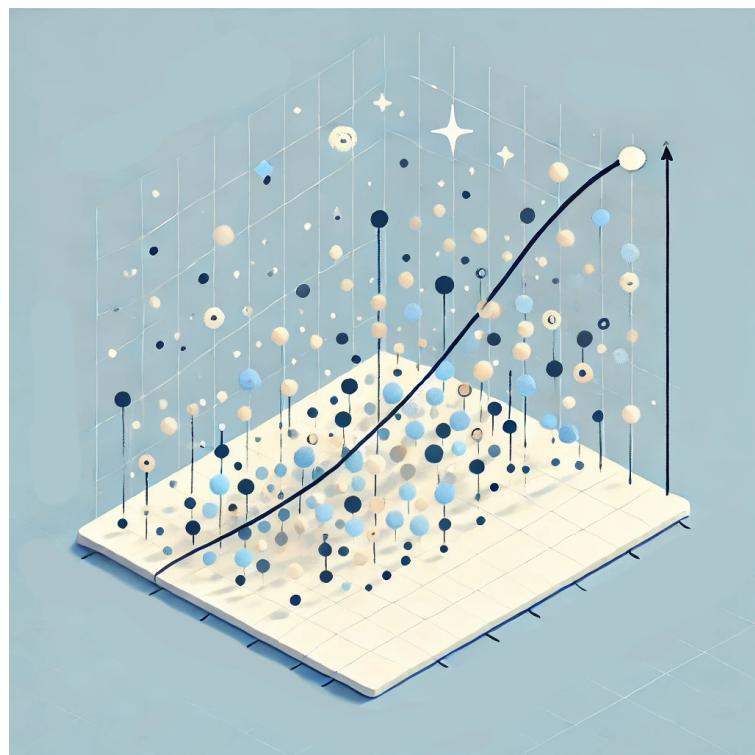
Figure 4.1: Summary of Alternate Sigmoid Like Curves

Note: All functions are normalized in such a way that their slope at the origin is 1

Equation 5

Correlation

$$\text{Correlation} = \frac{\text{Cov}(X,Y)}{\text{Std}(X) \cdot \text{Std}(Y)}$$



x_i and y_i are linearly related

Key ML Equation 5: Correlation

$$\text{Correlation} = \frac{\text{Cov}(X, Y)}{\text{Std}(X) \cdot \text{Std}(Y)} \quad (5.0.1)$$

(5.0.2)

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (5.0.3)$$

Cov(X, Y) The **covariance** between vectors X and Y , which measures the degree to which the two variables (or vectors) vary together.

Std(V) The **standard deviation of vector V** , which measures the spread of the elements of V around each mean.

\bar{x} The **mean of variable X** , calculated as $\frac{1}{n} \sum x_i$, where n is the number of data points.

\bar{y} The **mean of variable Y** , calculated as $\frac{1}{n} \sum y_i$, where n is the number of data points.

x_i The **individual data point** in the variable X , where i indexes each observation in the dataset.

y_i The **individual data point** in the variable Y , where i indexes each observation in the dataset.

Introduction: Correlation measures the strength and direction of the linear relationship between two variables.

Description: It ranges from -1 to 1, where values close to 1 or -1 indicate a strong relationship, and values close to 0 indicate no relationship.

Importance in ML: Correlation analysis helps in feature selection by identifying which features are related to the target variable. Strongly correlated features can be used for better predictions, whereas uncorrelated features can be removed to simplify models.

5.1 What is Behind the Equation

Correlation is a statistical measure that describes the strength and direction of a relationship between two variables. The most common measure of correlation is Pearson's correlation coefficient, denoted as r , which ranges from -1 to 1. A positive value indicates a direct relationship, while a negative value indicates an inverse relationship. A value of 0 implies no linear relationship.

The equation for Pearson's correlation r coefficient 5.0.3 can be compared to the normalized version 5.1.1:

$$r_{\mu} = \frac{\sum (x_i y_i)}{\sqrt{\sum x_i^2 \sum y_i^2}} \quad (5.1.1)$$

Where: - x_i and y_i are individual data points of variables x and y , - \bar{x} and \bar{y} are the means of variables x and y , respectively.

- **Variables:** The two variables, x and y , whose relationship we are measuring. These could be anything from economic indicators to user behaviors.
- **Data Points:** Each pair of values (x_i, y_i) represents an observation.
- **Means:** The mean values of x and y are used to standardize the data by subtracting them from each data point, which helps in normalizing the relationship.
- x_i and y_i are individual data points for the two variables X and Y .
- \bar{x} and \bar{y} are the means (averages) of the respective variables.
- **Numerator:** The numerator is the **covariance** between X and Y , which measures how the two variables vary together.
- **Denominator:** The denominator is the product of the **standard deviations** of X and Y , which normalizes the covariance and ensures that the correlation coefficient lies between -1 and 1.

The Pearson correlation measures the strength and direction of the **linear relationship** between two variables. This is the most commonly used form of correlation.

5.1.0.1 Use case:

This equation is used when you want to understand the **linear relationship** between two continuous variables and is the formula used by the `cor()` function in R with the default method (`method = "pearson"`).

5.1.0.2 Difference in Interpretation:

- The pearson equation 5.0.3 measures **how deviations from the mean of each variable are related**, i.e., it normalizes by the standard deviations of each variable.
- The second equation 5.1.1 compares the raw products of the variables, but it does **not center** the data by subtracting the means. This makes the second equation more sensitive to the scales and magnitudes of the variables, which is why it's often less preferred for calculating correlation in its standard sense.

5.1.1 Key Differences:

1. Normalization:

- The first equation involves centering the data by subtracting the means, which is important to account for the relationship of the variables in terms of deviations from their respective averages.
- The second equation does not involve centering the data. It uses the raw values of x_i and y_i , making it sensitive to the overall scale and magnitude of the data.

2. Covariance vs Raw Product:

- The first equation is essentially computing a normalized **covariance** (divided by the product of the standard deviations of X and Y).
- The second equation uses the **raw product** of the variables, which lacks the standardization and may lead to different numerical results.

3. Interpretation:

- The first equation (Pearson's correlation) measures the strength of a **linear relationship**, normalized by the variance in both variables.
- The second equation is more general but is not commonly used as a standard for measuring correlation. It may appear in specific contexts, such as **cosine similarity** when the data is already standardized or when normalization is not necessary.

5.1.2 Final point summary Pearson vs normed at

- **The first equation** 5.0.3 is the **Pearson correlation coefficient**, the standard method for measuring linear relationships, normalized by the standard deviations of both variables.
- **The second equation** 5.1.1 is a raw form that compares the products of the variables and is **not typically used** for Pearson's correlation, as it lacks normalization. It could be seen as a raw form of a **cosine similarity** (for vector comparison) or some other

variant of correlation, but it's not generally used for typical correlation measurements.

5.2 Historical Context

The concept of correlation was first developed by Sir Francis Galton in the 19th century. Galton's work laid the foundation for the field of regression analysis and correlational statistics. The Pearson correlation, developed by Karl Pearson in 1896, is the most widely used form of correlation in modern statistics.

5.3 Understanding the Notation

The formula for Pearson's correlation coefficient involves the concept of covariance, which measures how much two variables change together. The denominator in the formula standardizes the covariance by dividing by the product of the standard deviations of x and y .

- **Covariance:** Measures the degree to which two variables move together.
- **Standard Deviation:** Measures the spread of a variable around its mean.

5.4 Applications in Machine Learning

In machine learning, correlation is used to:

- Identify relationships between input features and target variables,
- Check for multicollinearity, where highly correlated features might skew results,
- Perform feature selection, as correlated features may provide redundant information.

For instance, if two features are highly correlated, one may be removed from the model to simplify it without losing predictive power.

5.5 Correlation in R

To compute the correlation coefficient in R, you can use the built-in `cor()` function:

```
# Example R code to compute Pearson's correlation coefficient
x <- c(1, 2, 3, 4, 5)
y <- c(5, 4, 3, 2, 1)
correlation <- cor(x, y)
# print{correlation}
cat(x)

## 1 2 3 4 5
cat(y)

## 5 4 3 2 1
cat(correlation)

## -1
```

5.6 Scatter Plots with Linear Overlays

In this extension of the correlation chapter, we will demonstrate how to visualize the relationship between two variables using scatter plots and add a linear regression line to the plots. We will use two different sets of scatter data for illustration.

5.6.1 Example 1: Scatter Data with Linear Overlay, $N_{pdf}(\mu, \sigma)$

$$N_{pdf}(\mu, \sigma) = N_{pdf}(\mu = 1, \sigma = 1)$$

```
# Example 1: Create scatter data
set.seed(42) # For reproducibility

# Generate example data
x1 <- rnorm(100, mean = 5, sd = 2) # 100 random data points from a normal distribution
y1 <- 3 + 2 * x1 + rnorm(100, mean = 0, sd = 1) # Linear relationship with some noise

# Plotting the scatter plot with a linear regression line
plot(x1, y1, main = "Scatter Plot with Linear Overlay Ex1",
      xlab = "X1", ylab = "Y1", pch = 19, col = "blue")
abline(lm(y1 ~ x1), col = "red", lwd = 2) # Add linear regression line
```

Scatter Plot with Linear Overlay Ex1

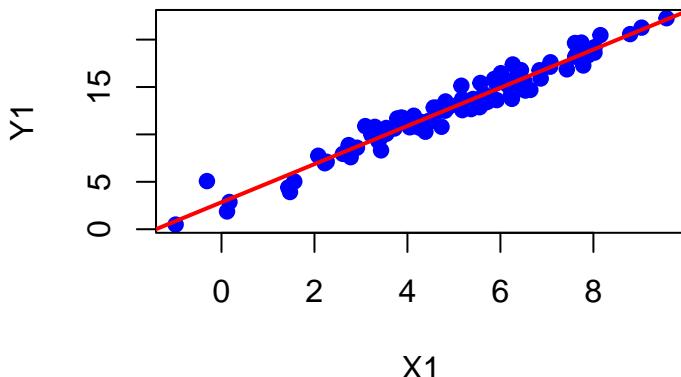


Figure 5.1: Default caption for all figures

```
# Calculate Pearson's correlation coefficient
r1 <- cor(x1, y1)
cat("Pearson's Correlation for Example 1: ", r1, "\n")

## Pearson's Correlation for Example 1: 0.9775592
```

5.6.2 Example 1: Scatter Data with Linear Overlay and some noise

```
# Example 2: Create another set of scatter data
set.seed(24) # For reproducibility

# Generate another set of example data
x2 <- rnorm(100, mean = 10, sd = 5) # 100 random data points from a different normal distribution
y2 <- 1 + 0.5 * x2 + rnorm(100, mean = 0, sd = 2) # Another linear relationship with noise
```

```
# Plotting the scatter plot with a linear regression line
plot(x2, y2, main = "Scatter Plot with Linear Overlay Ex2",
      xlab = "X2", ylab = "Y2", pch = 19, col = "green")
abline(lm(y2 ~ x2), col = "purple", lwd = 2) # Add linear regression line
```

Scatter Plot with Linear Overlay Ex2

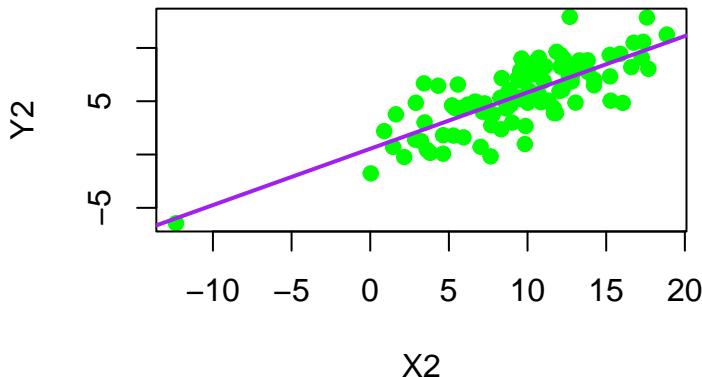


Figure 5.2: Default caption for all figures

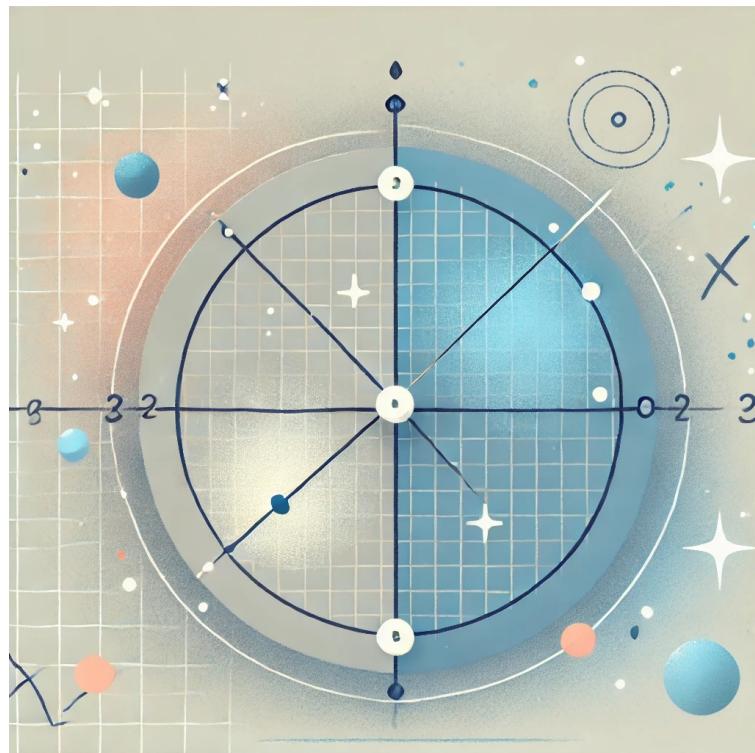
```
# Calculate Pearson's correlation coefficient
r2 <- cor(x2, y2)
cat("Pearson's Correlation for Example 2: ", r2, "\n")

## Pearson's Correlation for Example 2:  0.783414
```


Equation 6

Cosine Similarity

$$S_C(A, B) := \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$



You're comparative angle and measure

Key ML Equation 6: Cosine Similarity

$$\text{cosine similarity} = S_C(A, B) := \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \cdot \sqrt{\sum_{i=1}^n b_i^2}} \quad (6.0.1)$$

- A The **first vector** A , which represents the features (such as word frequencies or embedding values) of the first document or data point.
- B The **second vector** B , which represents the features of the second document or data point.
- $A \cdot B$ The **dot product** of vectors A and B , which computes the sum of the products of corresponding components of the two vectors. This measures how aligned the vectors are in the same direction.
- $\|A\|$ The **norm of vector** A , calculated as $\|A\| = \sqrt{\sum A_i^2}$, which measures the magnitude (length) of vector A .
- $\|B\|$ The **norm of vector** B , calculated as $\|B\| = \sqrt{\sum B_i^2}$, which measures the magnitude (length) of vector B .
- similarity** The **cosine similarity** between vectors A and B , which quantifies the cosine of the angle between the vectors, effectively measuring the directional alignment or similarity between them.

Introduction: Cosine similarity measures the cosine of the angle between two vectors in a multi-dimensional space.

Description: It is commonly used to determine how similar two documents are, regardless of their size, by comparing their word vectors.

Importance in ML: Cosine similarity is crucial in text mining and information retrieval applications. It is often used in clustering and classification algorithms for textual data, where measuring the similarity between vectors is needed.

6.1 What is Behind the Equation

6.1.1 Cosine Similarity Technically speaking

Cosine similarity 6.0.1 measures the cosine of the angle between two vectors \mathbf{A} and \mathbf{B} in a multi-dimensional space. The equation for cosine similarity is given by:

Where: - $\mathbf{A} \cdot \mathbf{B}$ is the dot product of vectors \mathbf{A} and \mathbf{B} , - $\|\mathbf{A}\|$ and $\|\mathbf{B}\|$ are the magnitudes (or norms) of vectors \mathbf{A} and \mathbf{B} .

This formula results in a value between -1 and 1, where: - 1 indicates the vectors are identical (or point in the same direction), - 0 indicates the vectors are orthogonal (i.e., no similarity), - -1 indicates the vectors are diametrically opposed (i.e., completely dissimilar).

Cosine similarity is a measure used to compute the **cosine of the angle** between two non-zero vectors in an inner product space. It is commonly used to measure the similarity between documents in text mining, recommendation systems, and machine learning tasks. The basic idea is to calculate how similar two vectors are, regardless of their magnitudes, by comparing their directions.

Where: - $\mathbf{A} \cdot \mathbf{B}$ is the **dot product** of vectors \mathbf{A} and \mathbf{B} , - $\|\mathbf{A}\|$ and $\|\mathbf{B}\|$ are the **norms** (or magnitudes) of the vectors \mathbf{A} and \mathbf{B} , respectively.

Cosine similarity ranges from -1 to 1: - 1 indicates that the vectors are identical (i.e., they point in the same direction), - 0 indicates orthogonality (no similarity), - -1 indicates that the vectors are diametrically opposed (pointing in opposite directions).

- **Dot Product:** The dot product $A \cdot B$, $A \cdot B = \sum_{i=1}^n a_i b_i$, is calculated as the sum of the products of corresponding elements in the vectors. It measures how much two vectors point in the same direction. Mathematically, the dot product is represented as: Where a_i and b_i represent the individual components of vectors A and B , respectively. This formula sums up the element-wise product of the two vectors.
- **Magnitude:** The magnitude of a vector A is calculated as $\|A\| = \sqrt{\sum_{i=1}^n a_i^2}$, where a_i represents the individual components of vector A . The magnitude measures the length of the vector.
- **Cosine of the Angle:** The cosine of the angle between the vectors is the core of cosine similarity. If the vectors are pointing in the same direction, the cosine value will be 1, indicating high similarity.

6.1.2 Cosine Similarity Mathematically speaking

Mathematically, the cosine similarity between two vectors \mathbf{A} and \mathbf{B} is given by:

$$\text{similarity}(\mathbf{A}, \mathbf{B}) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

Where: - $\mathbf{A} \cdot \mathbf{B}$ is the dot product of vectors \mathbf{A} and \mathbf{B} , - $\|\mathbf{A}\|$ and $\|\mathbf{B}\|$ are the norms (magnitudes) of the vectors.

Now, considering linear separability, for two classes of data to be linearly separable, the vectors corresponding to those classes should be **distinct** in terms of their direction. This can be formalized as follows:

- If the vectors for the two classes are aligned (i.e., pointing in the same direction), their cosine similarity will be close to **1**, indicating that the classes are **highly similar** and not linearly separable.
- If the vectors are **orthogonal**, their cosine similarity will be **0**, implying that they are **completely different** and hence easily separable by a linear classifier.

To represent this mathematically, if we have two vectors \mathbf{A}_1 and \mathbf{A}_2 representing the two classes, the cosine similarity between them can be used as follows:

$$\mathbf{A}_1 \cdot \mathbf{A}_2 = \|\mathbf{A}_1\| \|\mathbf{A}_2\| \cos(\theta)$$

Where: - θ is the angle between the two vectors.

If the vectors are **orthogonal** (i.e., $\theta = 90^\circ$), then:

$$\cos(90^\circ) = 0 \Rightarrow \mathbf{A}_1 \cdot \mathbf{A}_2 = 0$$

This means that the cosine similarity is 0, and the two classes are perfectly separable using a linear classifier.

For the vectors to be **linearly separable**, a good condition is that the cosine similarity between the class vectors should be **low** or **zero**. If the cosine similarity is high (close to 1), it means that the vectors are pointing in similar directions, and the classes might not be easily separable by a linear boundary.

In summary: In summary:

- **Cosine similarity close to 1:** Vectors are pointing in the same direction, classes are **not separable**.
- **Cosine similarity close to 0:** Vectors are orthogonal, classes are **easily separable**.

This is crucial for algorithms like **SVMs** and **Logistic Regression**, which rely on the geometric properties of the data for finding the decision boundary.

6.2 Historical Context

Cosine similarity has become a cornerstone of text mining and Natural Language Processing (NLP) due to its ability to compare documents and their word distributions. Its origins and application milestones are tied closely to developments in both computational linguistics and information retrieval systems. Below is an expanded historical overview of cosine similarity, focusing on its key milestones in the evolution of machine learning, text mining, and NLP.

6.2.1 The Historical Development of Cosine Similarity

- Early Foundations: Vector Space Models

Cosine similarity owes much of its origin to the Vector Space Model (VSM), a mathematical framework developed in the 1950s and 1960s for representing text documents as vectors in a multi-dimensional space.

1950s-1960s: Early work in information retrieval (IR) focused on using Boolean models and term frequency approaches to determine document relevance. However, these models lacked the ability to handle the nuances of natural language. The Vector Space Model (VSM), first popularized by Gerard Salton and colleagues at Cornell University in the late 1960s, brought about a shift in how text was represented. In VSM, a document was transformed into a vector, where each dimension represented the frequency of a term or a word in that document. This allowed for documents to be compared based on their word distributions, thus setting the stage for using similarity measures.

- The Birth of Cosine Similarity in Information Retrieval (1970s-1980s)

The concept of cosine similarity, specifically as a measure of the angular distance between two document vectors, emerged as a way to quantify the similarity between documents without being affected by their lengths.

1970s: Salton and McGill in their seminal book *Introduction to Modern Information Retrieval* (1983) formalized the use of cosine similarity as part of the Vector Space Model. The key insight was that the cosine of the angle between two document vectors provided a measure of how similar the documents were, regardless of their lengths or word counts. This normalized measure became the foundation for retrieval models in IR systems. 1980s-1990s: Cosine similarity began to dominate the field of information retrieval, where the goal was to rank documents in response to queries. The Vector Space Model, with cosine similarity as its central metric, was widely used in database search engines and document retrieval systems. During this period, IR research emphasized the importance of efficiently comparing documents in large datasets using similarity measures like cosine similarity. The Rise of Natural Language Processing (NLP) and Text Mining (1990s-2000s)

With the rise of computational linguistics and NLP, cosine similarity found itself at the heart of various machine learning and text analysis applications. The key development in this period was the shift from purely mathematical and statistical approaches to semantics-driven models in language processing.

- Cosine Similarity into the 1990s:

As large datasets and corpora of text became available (such as the Reuters-21578 dataset for news classification), cosine similarity was used to measure the similarity between documents and query-document pairs in information retrieval. This time also marked the rise of TF-IDF (Term Frequency-Inverse Document Frequency), a weighting scheme that was commonly combined with cosine similarity to enhance document comparison by weighting terms based on their importance. Late 1990s-2000s: Latent Semantic Analysis (LSA), introduced by Deerwester et al. in 1990, expanded the role of cosine similarity. LSA used singular value decomposition (SVD) to reduce the dimensionality of term-document matrices, with cosine similarity serving as the central metric to identify semantic similarity between documents. This was particularly useful for handling synonymy and polysemy (words having multiple meanings).

The development of word embeddings (early models like Latent Dirichlet Allocation (LDA), and later Word2Vec, GloVe, and FastText) built upon the concepts introduced by LSA, enabling more nuanced similarity measures for word and document comparisons. Cosine Similarity in Modern Machine Learning (2010s-present)

With the advent of deep learning and large-scale natural language models like BERT and GPT, cosine similarity continues to be widely used in modern NLP tasks, but with significant advancements.

- Cosine similarity in the birth of the ML age 2010s:

Cosine similarity remains one of the most widely used metrics in document clustering, topic modeling, and recommendation systems. The advent of word embeddings like Word2Vec (by Tomas Mikolov et al. at Google) revolutionized cosine similarity, as it enabled semantic word vector representations rather than just term frequency counts.

Word2Vec and GloVe create dense vector representations of words, and cosine similarity became a powerful way to compare these dense vectors and measure the similarity of word meanings. This was a significant leap from counting word occurrences to capturing the contextual relationships between words, documents, or even entire corpora.

Deep Learning Integration:

In modern applications, transformer-based models such as BERT, GPT, and T5 leverage pre-trained embeddings that can be compared using cosine similarity for tasks like semantic

textual similarity and document ranking. These models capture fine-grained meanings of words and sentences, and cosine similarity is often used to compare the output of these models for tasks such as document retrieval, question answering, and text classification.

- Cosine Similarity in Information Retrieval and Search Engines (2010s-present)

The field of search engines has seen a continual refinement of cosine similarity, especially in the era of semantic search:

Query-Document Matching: Traditional search engines used cosine similarity to rank documents by relevance to a given query. With the advent of semantic search, search engines now use more sophisticated embeddings (e.g., from BERT or T5) to map both queries and documents into a shared semantic space, where cosine similarity is used to measure their relevance, not just based on keyword overlap but on the underlying meaning of the content.

Reinforcement Learning and Ranking: Modern search engines like Google, Bing, and Yahoo still use cosine similarity in ranking documents, although they have become more sophisticated with reinforcement learning-based approaches that adjust ranking according to user interactions. The query understanding has evolved from simple keyword matching to contextual semantic understanding powered by neural networks.

- Milestones in the Application of Cosine Similarity in ML

Early Applications in Information Retrieval (1970s-1980s): Cosine similarity was first used in information retrieval to compare documents based on term frequency vectors. **Introduction of LSA (1990s):** LSA expanded the role of cosine similarity, allowing for the comparison of documents in a semantic space.

Word Embeddings (2000s-present): Cosine similarity became integral in word embeddings (e.g., Word2Vec, GloVe), which represent words in continuous vector spaces, enabling the comparison of words and documents at a semantic level. **Transformer Models (2018-present):** Cosine similarity continues to be used in modern semantic search engines, document clustering, and semantic textual similarity tasks, powered by state-of-the-art transformer-based language models like BERT and GPT.

6.3 Cosine Similarity - The basics and the usage

Cosine similarity is a fascinating measure, and there are a few interesting mathematical identities, properties, and quirks related to it. Here are some of them:

6.3.1 Cosine Similarity and the Angle Between Vectors:

Cosine similarity essentially measures the cosine of the angle between two vectors in multidimensional space. This brings us to some interesting geometric properties:

- Zero Similarity: If the cosine similarity is 0 , it implies that the angle between the vectors is 90° (i.e., the vectors are orthogonal). In text mining, this would mean two documents have no terms in common.
- Negative Similarity: A negative cosine similarity indicates that the vectors are pointing in opposite directions (i.e., the angle is greater than 90° but less than 180°).
- Unit Vectors: If the vectors A and B are normalized to unit vectors (i.e., $\|A\| = \|B\| = 1$), then cosine similarity becomes equivalent to the dot product $A \cdot B$. This property makes cosine similarity a normalized version of the dot product, and it's particularly useful in highdimensional spaces like text data, where magnitude can vary widely across documents.

6.3.2 Cosine Similarity and the Inner Product:

Cosine similarity is directly related to the inner product (or dot product) of vectors.

- When you normalize the vectors (making them unit vectors), cosine similarity becomes the inner product itself. This is why cosine similarity is sometimes referred to as a normalized form of the dot product.
- In vector spaces, the inner product captures both the magnitude and direction of the vectors, and cosine similarity isolates only the directional similarity by normalizing out the magnitudes.

3. **Connection with the Euclidean Distance:

- Interestingly, cosine similarity is connected to Euclidean distance through a well-known identity:

$$\text{Cosine Similarity} = 1 - \frac{\|A - B\|^2}{\|A\|^2 + \|B\|^2} = \frac{A \cdot B}{\|A\|\|B\|}$$

This identity shows that the cosine similarity can be interpreted in terms of the squared Euclidean distance between the vectors. High cosine similarity means low Euclidean distance, and vice versa. This provides a duality between these two commonly used measures of similarity in machine learning.

6.4 Cosine Similarity and Linear Separability

Let's begin by recalling the formulation of linear separability from the matrix formulation and connecting it to cosine similarity:

1. **Linear Separability Condition:** For a set of data points \mathbf{x}_i with labels y_i , the condition for linear separability is:

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) > 0, \quad \forall i = 1, 2, \dots, n$$

This condition means that all data points are correctly classified by the hyperplane, which is determined by \mathbf{w} and b .

2. **Cosine Similarity and Separation:** The cosine similarity between two vectors \mathbf{x}_i and \mathbf{w} (where \mathbf{w} is the normal vector to the hyperplane) is given by:

$$\text{similarity}(\mathbf{x}_i, \mathbf{w}) = \frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{w}\| \|\mathbf{x}_i\|}$$

- When cosine similarity is high (close to 1), the vector \mathbf{x}_i is pointing in a similar direction to \mathbf{w} , and the data point \mathbf{x}_i is not well separated from other points in the feature space.
- When cosine similarity is low (close to 0), the vectors are orthogonal, indicating that the points from different classes are well-separated, and the decision boundary (hyperplane) is more likely to be effective at separating the two classes.
- **Linear Separability and Margin Maximization:** In Support Vector Machines (SVM), the goal is to find the optimal hyperplane that maximizes the margin between the two classes. The margin is the perpendicular distance from the closest data points (the support vectors) to the hyperplane. The SVM optimization problem is formulated as:

$$\text{Maximize } \frac{2}{\|\mathbf{w}\|}$$

Subject to the constraint:

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad \forall i$$

This ensures that the data points are classified correctly with the maximum margin, which is geometrically related to the angle between the weight vector \mathbf{w} and the data points \mathbf{x}_i .

6.4.1 How Cosine Similarity Relates to the Margin in SVM:

The margin in SVM is inversely proportional to the magnitude of the weight vector \mathbf{w} . A larger margin corresponds to a lower value of $\|\mathbf{w}\|$. Since cosine similarity measures the

alignment between the vectors, it directly relates to the angle between the weight vector \mathbf{w} and the data vectors \mathbf{x}_i :

- A large cosine similarity means the vectors are pointing in the same direction, making it harder for the classifier to separate the classes with a hyperplane.
- A low cosine similarity (especially close to 0) means the vectors are more orthogonal, making the data points more separable.

6.4.2 Full Mathematical Representation Connecting Linear Separability and Cosine Similarity:

To summarize the relationship between cosine similarity and linear separability, we have:

1. Cosine similarity between the weight vector \mathbf{w} and each data point \mathbf{x}_i :

$$\text{similarity}(\mathbf{x}_i, \mathbf{w}) = \frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{w}\| \|\mathbf{x}_i\|}$$

2. For linear separability, the data points must satisfy the condition:

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) > 0, \quad \forall i = 1, 2, \dots, n$$

3. The margin between the classes is maximized when the cosine similarity between the weight vector \mathbf{w} and the data points \mathbf{x}_i is low (i.e., the vectors are orthogonal):

$$\text{Maximize margin: } \frac{2}{\|\mathbf{w}\|}$$

This leads to a large margin when cosine similarity is small, ensuring that the data points are easily separable by the hyperplane.

Thus, the cosine similarity between the weight vector \mathbf{w} and the data points \mathbf{x}_i plays a key role in determining the separability of the data. When cosine similarity is low, it indicates that the data is well separated and linearly separable.

6.4.3 Cosine Similarity with Sparse Vectors:

- One interesting property of cosine similarity is its robustness in high-dimensional, sparse data (such as text data). In sparse high-dimensional spaces, many of the components of the vectors are zero, and the cosine similarity reduces to comparing the non-zero entries of the vectors. This makes it very useful in scenarios like text analysis or recommendation systems, where documents or items are represented by high-dimensional sparse vectors, and many entries in these vectors are zero.

6.4.4 Cosine Similarity as a Measure of Clustering Quality:

- Cosine similarity is widely used in clustering algorithms, such as k-means and k-medoids, especially in text clustering. It helps measure how well different data points (such as documents) are grouped together. Interestingly, when cosine similarity is used in clustering:
- Cosine similarity doesn't rely on absolute magnitude: Unlike Euclidean distance, it doesn't care about the vector magnitudes, so documents of different lengths or sizes can still be compared meaningfully.
- Cluster Centers in Text Mining: When performing clustering with text data, the center of a cluster (the centroid) might not be a point in the vector space, but it's computed as the mean of cosine similarities with respect to all other points in the cluster.

6.4.5 Cosine Similarity and Vector Scaling:

- Scaling the vectors (by multiplying them by a constant) does not affect the cosine similarity. This is because cosine similarity is a function of the angle between vectors, not their magnitude. Hence, if two vectors point in the same direction but have different magnitudes, their cosine similarity will be 1. This is a useful property in NLP, where scaling the vector components doesn't change their relative orientation.

6.4.6 Cosine Similarity as a Measure of Clustering Quality:

- Cosine similarity is widely used in clustering algorithms, such as k-means and k-medoids, especially in text clustering. It helps measure how well different data points (such as documents) are grouped together. Interestingly, when cosine similarity is used in clustering:
- Cosine similarity doesn't rely on absolute magnitude: Unlike Euclidean distance, it doesn't care about the vector magnitudes, so documents of different lengths or sizes can still be compared meaningfully.
- Cluster Centers in Text Mining: When performing clustering with text data, the center

of a cluster (the centroid) might not be a point in the vector space, but it's computed as the mean of cosine similarities with respect to all other points in the cluster.

6.4.7 Cosine Similarity and Vector Scaling:

- Scaling the vectors (by multiplying them by a constant) does not affect the cosine similarity. This is because cosine similarity is a function of the angle between vectors, not their magnitude. Hence, if two vectors point in the same direction but have different magnitudes, their cosine similarity will be 1. This is a useful property in NLP, where scaling the vector components doesn't change their relative orientation.

6.4.8 Cosine Similarity as a Measure of Word Similarity in NLP:

- In NLP, word embeddings like Word2Vec and GloVe represent words as vectors. Cosine similarity is often used to measure semantic similarity between words based on their vector representations. Words with similar meanings (like “king” and “queen”) will have high cosine similarity, even though they are not identical in form.

6.4.9 Cosine Similarity as an Angular Measure:

- The cosine similarity is essentially the cosine of the angle between the two vectors. If the two vectors are exactly the same, the cosine similarity is 1 (indicating a 0° angle). If they are perpendicular, the cosine similarity is 0 (indicating a 90° angle). If they are opposite, the cosine similarity is -1 (indicating a 180° angle).
- This geometric interpretation is useful in visualizing vector spaces. For instance, when comparing documents, vectors pointing in the same direction are highly similar, and vectors pointing in opposite directions are highly dissimilar.

6.4.10 Cosine Similarity and Textual Data:

- In many information retrieval systems (like search engines), documents are represented by TF-IDF vectors (Term Frequency-Inverse Document Frequency). Cosine similarity is used to measure the relevance of a document to a query by comparing the query vector to the document vectors. Interestingly, because of the high-dimensional and sparse nature of TF-IDF vectors, cosine similarity becomes one of the most efficient ways to compare documents, avoiding issues that might arise with traditional distance measures like Euclidean distance.

6.4.11 Cosine Similarity and Regularization:

- Cosine similarity also plays a role in some machine learning models, especially those used for regularization. For instance, in sparse representations (such as Lasso regression or feature selection in text), cosine similarity can be used to measure how closely related a feature is to another, helping in the regularization process by identifying similar features or words.

6.5 Relation to Spherical Geometry:

Cosine similarity is deeply connected to spherical geometry. In fact, it is the spherical distance (or the angular distance) between two points on the unit sphere. The closer two vectors are in the direction they point, the smaller the angle between them, and therefore, the higher the cosine similarity. - This connection becomes evident when vectors are normalized (as unit vectors) and placed on the surface of a unit sphere. In this setting, the cosine of the angle between the vectors directly gives the cosine of the spherical angle between them, which is central in many algorithms that perform clustering and classification based on angles, such as kmeans clustering on normalized data.

6.6 Cosine Similarity and Kernel Methods:

Cosine similarity can also be seen as a kernel function in the context of machine learning. In the context of support vector machines (SVMs) and other kernel-based learning methods, cosine similarity can serve as a linear kernel for high-dimensional vector spaces. - The linear kernel in SVM is essentially equivalent to calculating the cosine similarity between data points when the data is centered and normalized. Thus, it provides a way to map the data into a high-dimensional space without explicitly calculating the dot product in that space.

6.7 Cosine Similarity and Information Retrieval:

In information retrieval (IR), cosine similarity is used to compare the vector representations of documents. The documents are typically represented in terms of their TF-IDF (Term Frequency-Inverse Document Frequency) vectors, and cosine similarity is used to measure how similar a query is to a document. - In this context, cosine similarity measures the angular distance between the query vector and document vectors. Despite differences in the document lengths, the cosine similarity focuses on the relative frequency of words, making it especially useful for content based recommendation systems and document clustering.

6.8 Cosine Similarity and the Norms of Vectors:

Cosine similarity is inherently related to vector norms. In particular, Euclidean distance (L₂ norm) is closely tied to cosine similarity, but cosine similarity itself is not a norm in the traditional sense. Let's break this down:

- Euclidean Norm (L₂ norm): The L₂ norm (also called the Euclidean norm) of a vector $\mathbf{A} = [a_1, a_2, \dots, a_n]$ is defined as:

$$\|\mathbf{A}\|_2 = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$$

It measures the “length” or magnitude of the vector.

- Cosine Similarity: Cosine similarity between two vectors \mathbf{A} and \mathbf{B} is given by:

$$\text{similarity}(\mathbf{A}, \mathbf{B}) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2}$$

Notice that cosine similarity involves normalizing the vectors by their L₂ norms. In this case, cosine similarity is related to the cosine of the angle between the vectors, which is a directional measure rather than a magnitude measure.

- 2. Cosine Similarity vs. Norms:
- Cosine similarity itself is not a norm, but it is closely related to the concept of a normalized inner product.
- A norm is a function that satisfies the following properties:

 - Non-negativity: $\|\mathbf{A}\| \geq 0$ with equality if and only if $\mathbf{A} = 0$,
 - Homogeneity: $\|\alpha\mathbf{A}\| = |\alpha| \|\mathbf{A}\|$ for any scalar α ,
 - Triangle inequality: $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$.

Cosine similarity does not satisfy the triangle inequality or the other properties needed to qualify as a norm. Therefore, cosine similarity is not a norm. Instead, it's a measure of similarity between vectors that is normalized to account for the lengths of the vectors, so that only the angle between them influences the result.

6.9 Can Cosine Similarity vs What can be used as a Norm?

Cosine similarity itself cannot be used as a norm, because, as mentioned, it doesn't satisfy the necessary properties of a norm. However, cosine distance, which is defined as:

$$\text{distance}(\mathbf{A}, \mathbf{B}) = 1 - \text{similarity}(\mathbf{A}, \mathbf{B}) = 1 - \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2}$$

can serve as a distance measure in a space where similarity is measured by cosine. While cosine similarity is not a norm, cosine distance satisfies the triangle inequality and other

properties of a metric, making it a valid distance measure.

Thus, cosine distance is a valid mathematical structure to measure dissimilarity between vectors, and although it isn't a norm itself, it's related to vector distances in a space defined by cosine similarity. 4. Connections with Other Norms:

While cosine similarity is related to the L2 (Euclidean) norm, it also appears in the context of other norms. Let's explore a few: - L1 Norm (Manhattan Norm): The L1 norm of a vector $\mathbf{A} = [a_1, a_2, \dots, a_n]$ is defined as:

$$\|\mathbf{A}\|_1 = |a_1| + |a_2| + \dots + |a_n|$$

- The cosine similarity measure is independent of the magnitude of the vectors but still considers their direction. In some applications, we could use cosine similarity in conjunction with the L1 norm to measure the similarity between high-dimensional sparse vectors (common in text analysis) where non-zero entries are few and far between.
- L ∞ Norm (Maximum Norm): The L ∞ norm (also called the Chebyshev norm) is given by:

$$\|\mathbf{A}\|_\infty = \max(|a_1|, |a_2|, \dots, |a_n|)$$

- Although the L ∞ norm measures the largest element in the vector, it doesn't directly impact cosine similarity, but it's worth noting that the interaction between cosine similarity and other norms can become relevant in scenarios such as image retrieval or clustering, where the components of the vectors may represent pixel intensities or other quantities.

6.10 Cosine Distance:

6.10.1 Cosine distance is defined as:

$$\text{Cosine Distance } (\mathbf{A}, \mathbf{B}) = 1 - \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

This is the angular distance between two vectors in a high-dimensional space, and it ranges from 0 (when the vectors are identical) to 2 (when the vectors are diametrically opposed). The closer the cosine distance is to 0 , the more similar the vectors are.

6.10.2 Can Cosine Similarity be Used as a Norm? (Cosine Distance as a Metric)

While cosine similarity cannot be a norm due to the reasons above, the cosine distance, which is defined as:

$$\text{Cosine Distance}(\mathbf{A}, \mathbf{B}) = 1 - \text{similarity}(\mathbf{A}, \mathbf{B}) = 1 - \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2}$$

can be used as a distance metric. Cosine distance satisfies the necessary properties of a distance measure: 1. Non-negativity: Cosine Distance $(\mathbf{A}, \mathbf{B}) \geq 0$ with equality if and only if $\mathbf{A} = \mathbf{B}$. 2. Symmetry: Cosine Distance $(\mathbf{A}, \mathbf{B}) = \text{Cosine Distance}(\mathbf{B}, \mathbf{A})$. 3. Identity of indiscernibles: Cosine Distance $(\mathbf{A}, \mathbf{B}) = 0$ if and only if $\mathbf{A} = \mathbf{B}$. 4. Triangle Inequality: This property holds for cosine distance, meaning it satisfies the conditions to be a valid distance metric.

However, cosine distance is a metric (not a norm), and while it is related to the angle between vectors, it is not directly related to the magnitude of the vectors.

Using Cosine Distance as a Norm (Modified Version) While cosine similarity cannot directly be used as a norm, there are ways to modify or adapt concepts like cosine distance for scenarios where norms are needed. However, this requires a modification to the usual properties of cosine similarity.

Transforming Cosine Similarity into a Norm-like Measure: 1. Unit Vectors (Normalization):

In many applications, vectors are normalized to unit length, meaning $\|\mathbf{A}\|_2 = 1$ and $\|\mathbf{B}\|_2 = 1$. In this case, cosine similarity becomes a direct measure of how aligned or opposite the vectors are. The modified “norm” is then based on the angular distance between the vectors. However, this is still angular distance, not a true norm, because the triangle inequality may not hold in all cases.

2. Scaling with Cosine Distance:

To make cosine distance behave more like a norm, you would need to scale the vectors in such a way that cosine distance follows the norm properties. For example, if we scale the vectors before calculating the cosine distance and then apply a normalization strategy, we could create an approximation of a norm-like structure. However, this would still require further adjustments to make it truly behave as a norm in all cases.

6.10.3 Why Cosine Similarity is Not a Norm

Cosine similarity itself is not a norm due to the following reasons: 1. Non-negativity: A norm must satisfy the condition that it is non-negative and equal to 0 only when the vector is the zero vector (i.e., $\|\mathbf{A}\| = 0$ if and only if $\mathbf{A} = \mathbf{0}$). - Cosine similarity ranges from -1 to 1, and it can be negative when the vectors are diametrically opposed, which violates the non-negativity property required for norms. Specifically, cosine similarity measures the cosine of the angle between vectors, and negative values represent an angle greater than 90° (opposite directions), which is conceptually different from the magnitude measure required for norms. 2. Homogeneity: A norm must also satisfy the property of homogeneity, where scaling a vector by a scalar should scale the norm by the absolute value of that scalar:

$$\|\alpha\mathbf{A}\| = |\alpha|\|\mathbf{A}\|$$

- Cosine similarity doesn't satisfy this property because it is based on the angle between vectors, and scaling the vectors doesn't change the angle. However, norms require the scaling of the vector to scale the resulting value. This means cosine similarity doesn't respond to scaling in the way norms are supposed to.
3. Triangle Inequality: A norm must satisfy the triangle inequality, which states:

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$$

- Cosine similarity does not satisfy this property. Cosine similarity is only concerned with the angle between vectors and does not scale linearly with the magnitude of vectors, meaning it doesn't adhere to the triangle inequality required by norms.

Because of these issues, cosine similarity is not a norm in the strict mathematical sense.

6.10.4 Modified Norm in Specialized Contexts:

in some specialized contexts, especially those involving unit vectors on the unit sphere, the cosine distance can be informally referred to as an angular norm or spherical norm. Here's why: - Unit Vectors: When vectors are normalized to have a length of 1 (i.e., unit vectors), the cosine similarity becomes a direct measure of how aligned or opposite the vectors are. The distance between these vectors is angular because it reflects the angle between them. - Spherical Geometry: In spherical geometry, the distance between points on the surface of a sphere is often computed using the cosine of the angle between them. in this context, cosine distance (or cosine similarity) could be informally referred to as an angular distance or spherical norm. These are not true norms in the mathematical sense, but the terminology is used to describe the nature of the distance being calculated on the surface of a sphere.

6.11 Triangle inequality for cosine similarity

The ordinary triangle inequality for angles (i.e., arc lengths on a unit hypersphere) gives us that

$$|\angle AC - \angle CB| \leq \angle AB \leq \angle AC + \angle CB.$$

Because the cosine function decreases as an angle in $[0, \pi]$ radians increases, the sense of these inequalities is reversed when we take the cosine of each value:

$$\cos(\angle AC - \angle CB) \geq \cos(\angle AB) \geq \cos(\angle AC + \angle CB).$$

Using the cosine addition and subtraction formulas, these two inequalities can be written in terms of the original cosines,

$$\begin{aligned} \cos(A, C) \cdot \cos(C, B) + \sqrt{(1 - \cos(A, C)^2) \cdot (1 - \cos(C, B)^2)} &\geq \cos(A, B) \\ \cos(A, B) &\geq \cos(A, C) \cdot \cos(C, B) - \sqrt{(1 - \cos(A, C)^2) \cdot (1 - \cos(C, B)^2)} \end{aligned}$$

This form of the triangle inequality can be used to bound the minimum and maximum similarity of two objects A and B if the similarities to a reference object C is already known. This is used for example in metric data indexing, but has also been used to accelerate spherical k -means clustering the same way the Euclidean triangle inequality has been used to accelerate regular k -means.

6.11.1 Soft cosine measure

A soft cosine or (“soft” similarity) between two vectors considers similarities between pairs of features. The traditional cosine similarity considers the vector space model (VSM) features as independent or completely different, while the soft cosine measure proposes considering the similarity of features in VSM, which help generalize the concept of cosine (and soft cosine) as well as the idea of (soft) similarity.

For example, in the field of natural language processing (NLP) the similarity among features is quite intuitive. Features such as words, n -grams, or syntactic n -grams can be quite similar, though formally they are considered as different features in the VSM. For example, words “play” and “game” are different words and thus mapped to different points in VSM; yet they are semantically related. In case of n -grams or syntactic n -grams, Levenshtein distance can be applied (in fact, Levenshtein distance can be applied to words as well).

For calculating soft cosine, the matrix \mathbf{s} is used to indicate similarity between features. It can

be calculated through Levenshtein distance, WordNet similarity, or other similarity measures. Then we just multiply by this matrix. Given two N -dimension vectors a and b , the soft cosine similarity is calculated as follows:

$$\text{soft_cosine}(a, b) = \frac{\sum_{i,j}^N s_{ij}a_i b_j}{\sqrt{\sum_{i,j}^N s_{ij}a_i a_j} \sqrt{\sum_{i,j}^N s_{ij}b_i b_j}},$$

where s_{ij} = similarity (feature i , feature j). If there is no similarity between features ($s_{ii} = 1, s_{ij} = 0$ for $i \neq j$), the given equation is equivalent to the conventional cosine similarity formula.

The time complexity of this measure is quadratic, which makes it applicable to real-world tasks. Note that the complexity can be reduced to subquadratic. An efficient implementation of such soft cosine similarity is included in the Gensim open source library.

6.12 Cosine Similarity Applications in Machine Learning

Cosine similarity plays a crucial role in various machine learning tasks, particularly in text mining and information retrieval:

- **Text Classification:** It can be used to compare a document's word vector with a predefined class vector, allowing for the classification of documents based on similarity.
- **Clustering:** In clustering, cosine similarity is used to measure the similarity between data points (documents, images, etc.), helping algorithms group similar items together.
- **Recommendation Systems:** In collaborative filtering for recommendation systems, cosine similarity is used to measure the similarity between users or items based on their preferences or behaviors.

6.13 Cosine Similarity in R

Here's how to calculate the cosine similarity between two vectors using R:

```
# Example vectors
A <- c(1, 3, 4, 5)
B <- c(2, 4, 6, 8)

# Function to calculate cosine similarity
cosine_similarity <- function(A, B) {
  return(sum(A * B) / (sqrt(sum(A^2)) * sqrt(sum(B^2))))
}

# Calculate and print cosine similarity
similarity <- cosine_similarity(A, B)
cat("Cosine Similarity between A and B: ", similarity, "\n")

## Cosine Similarity between A and B:  0.9970545
```

6.14 Cosine Similarity in R, 3D-plot

```
# Load necessary libraries
library(ggplot2)

# Example 3D vectors
A <- c(1, 2, 3) # Vector A
B <- c(4, 5, 6) # Vector B

# Function to calculate cosine similarity
cosine_similarity <- function(A, B) {
  return(sum(A * B) / (sqrt(sum(A^2)) * sqrt(sum(B^2))))
}

# Calculate cosine similarity
similarity <- cosine_similarity(A, B)
cat("Cosine Similarity between A and B: ", similarity, "\n")

## Cosine Similarity between A and B: 0.9746318

# Prepare data for ggplot
vector_data <- data.frame(
  x = c(0, A[1], 0, B[1]), # X-coordinates
  y = c(0, A[2], 0, B[2]), # Y-coordinates
  z = c(0, A[3], 0, B[3]), # Z-coordinates
  group = c("A", "A", "B", "B")
)

# Convert data to 2D projection for plotting in ggplot
vector_data_2d <- data.frame(
  x = vector_data$x / (vector_data$z + 3), # Simple projection
  y = vector_data$y / (vector_data$z + 3), # Simple projection
  group = vector_data$group
)

# Create the 2D plot using ggplot
ggplot(vector_data_2d, aes(x = x, y = y, color = group)) +
  geom_segment(aes(xend = 0, yend = 0), size = 1.5) +
  geom_point(size = 4) +
  scale_color_manual(values = c("blue", "red")) +
  theme_minimal() +
  labs(title = "3D Cosine Similarity between Vectors A and B",
       x = "X", y = "Y") +
  theme(legend.position = "none") +
  coord_fixed(ratio = 1) # Fix aspect ratio for proper projection

# Save the plot as a PDF
ggsave("cosine_similarity_2d_plot.pdf")
```

3D Cosine Similarity between Vec

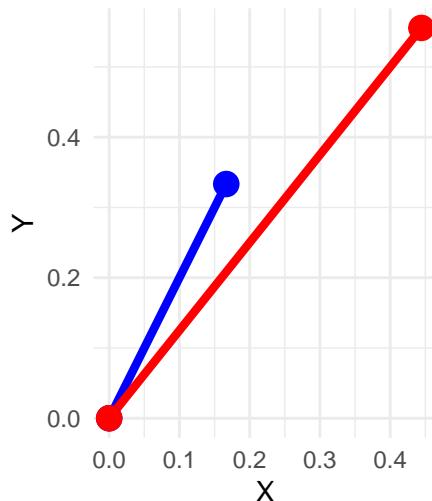
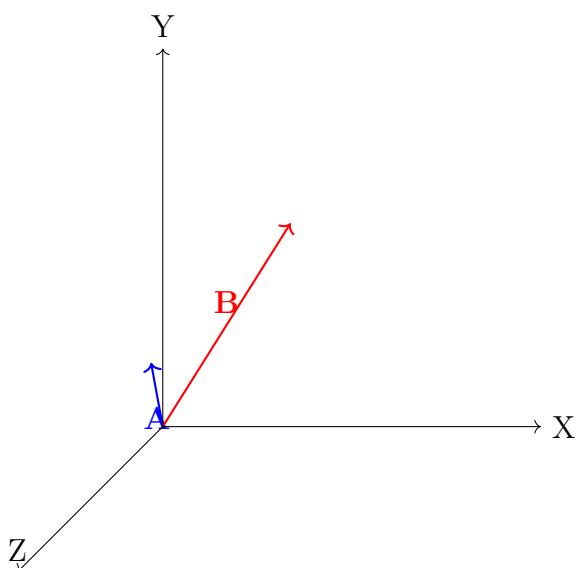


Figure 6.1: Default caption for all figures

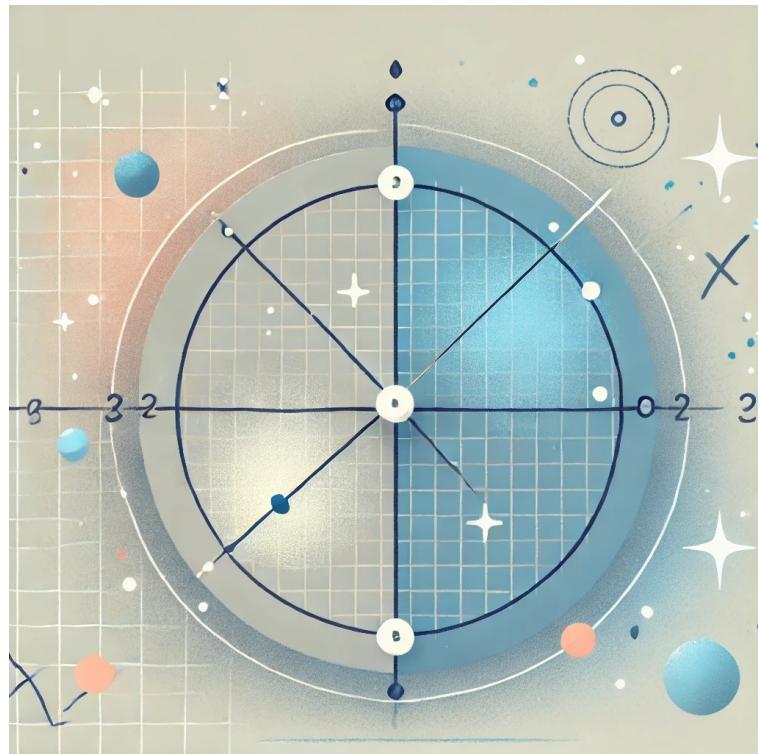
6.14.1 Tikz Output - cosine similarity



Equation 7

Naive Bayes

$$P(y|x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(x_1, \dots, x_n)}$$



BAYES somthing Whimsical - not yet

Key ML Equation 7: Naive Bayes

$$P(y|x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(x_1, \dots, x_n)} \quad (7.0.1)$$

x_i The **feature** x_i , representing the i -th feature of the data point (such as a specific word frequency in text classification, or a specific measurement in other domains).

y The **class label** y , which represents the class or category that the data point belongs to.

$P(y)$ The **prior probability** of class y , which represents the initial belief about the probability of class y occurring in the dataset.

$P(x_i|y)$ The **likelihood** of feature x_i given class y , which represents the probability of observing the feature x_i if the data point belongs to class y .

$P(x_1, \dots, x_n)$ The **evidence** or **normalizing constant**, which ensures that the total probability sums to 1. It is typically computed as $P(x_1, \dots, x_n) = \sum_y P(y) \prod_{i=1}^n P(x_i|y)$, but for classification, it serves as a normalization term.

Introduction: Naive Bayes is a probabilistic classifier based on Bayes' theorem with strong independence assumptions.

Description: It assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature, which simplifies computation.

Importance in ML: Naive Bayes is used for classification tasks, especially in text classification and spam detection. Its simplicity and efficiency make it a popular choice for high-dimensional datasets.

7.1 Naive Baye's - overview

Naive Bayes is a fundamental classification technique in machine learning that utilizes Bayes' theorem to make predictions based on feature independence. It is particularly valued for its simplicity and efficiency, making it suitable for large datasets and real-time applications.

Overview of Naive Bayes

Naive Bayes classifiers are a family of linear probabilistic classifiers that assume the features are **conditionally independent** given the target class. This assumption simplifies the computation of probabilities, allowing for quick and efficient classification. Despite its “naive” assumption, this method has proven effective in various applications, especially in natural language processing (NLP) tasks like spam detection and sentiment analysis.

Key Features::

Probabilistic Model: Naive Bayes uses conditional probability to predict class membership based on the likelihood of features occurring within each class.

Independence Assumption: It assumes that the presence of a particular feature does not affect the presence of any other feature, which simplifies calculations significantly.

Fast and Scalable: The algorithm is highly scalable, requiring only a small amount of training data to estimate parameters effectively, making it suitable for large datasets¹³⁴. Applications

Naive Bayes is widely used in various domains:

Text Classification: Such as spam filtering and sentiment analysis. Medical Diagnosis: Classifying diseases based on symptoms. Recommendation Systems: Suggesting products based on user preferences.

In order to discuss the Naive Bayes' equation and its application we first introduce the Bayes' equation.

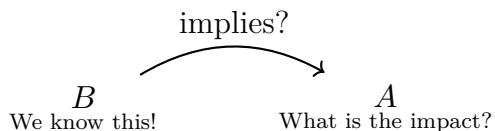
7.2 Introduction to Bayesian Theory

7.2.1 Introduction the traditional formulation of Bayes' theory

Bayes' Theorem is a powerful mathematical tool used in probability theory to calculate conditional probabilities. It is named after Reverend Thomas Bayes and articulates a fundamental relationship between the conditional probabilities of two events.

The basic concept revolves around updating our initial beliefs or estimates (prior probabilities) in light of new evidence or information (likelihood) to arrive at a revised belief (posterior probability). This process is particularly useful in various fields such as statistics, machine learning, and data analysis, where making predictions based on incomplete information is common.

The implication is simple in the sense we want to understand how “the occurrence of B” affects “the occurrence of A”



7.2.2 Bayes Formula

Definition 7.1: Traditional Bayes'

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad (7.2.1)$$

- A The **event of interest**, which represents the hypothesis or condition for which the probability is being computed.
- B The **evidence or observed event**, which represents the known condition or observation influencing the computation of $P(A | B)$.
- $P(A)$ The **prior probability** of event A , representing the initial belief about the likelihood of A occurring, before considering the evidence B .
- $P(B)$ The **marginal probability** of B , representing the total likelihood of observing B across all possible scenarios.
- $P(B | A)$ The **likelihood** of observing evidence B given that A is true. It quantifies how probable the evidence is under the hypothesis.
- $P(A | B)$ The **posterior probability** of A given B , representing the updated belief about the likelihood of A after considering the evidence B .

7.2.3 Explanation of the player in Bayes' formula

Prior Probability $P(A)$:

This is your initial belief about the probability of an event before any additional information is considered. For instance, if you are trying to determine the probability that a randomly chosen person is left-handed (A), your prior might be the overall percentage of left-handed people in the population.

Likelihood $P(B|A)$:

This reflects how probable the evidence (B) is, assuming the hypothesis (A) is true. For example, if a diagnostic test is known to detect a disease 95% of the time when the disease is present, then $P(B|A) = 0.95$.

Marginal Probability $P(B)$:

This is the total probability of observing the evidence (B) under all possible circumstances. It's often calculated as the sum of probabilities for all possible values of A : $P(B) = P(B|A) \cdot P(A) + P(B|\neg A) \cdot P(\neg A)$, where $\neg A$ represents the complement of A .

Posterior Probability $P(A|B)$:

This is what you're solving for. It's your updated belief about the probability of A after taking into account the new evidence B .

7.2.4 Rational Considerations of the meaning of the Bayes' formula

Taking a closer look at Bayes' formula:

$$P(A|B) = \frac{\text{Posterior}}{\text{Likelihood}} \cdot \frac{\text{Prior}}{\text{Evidence}}$$

$$P(A|B) = P(B|A) \cdot \frac{P(A)}{P(B)}$$

The key to tying this all together is to understand the following key expression:

$$\text{Posterior} = \text{Likelihood} \cdot \frac{\text{Prior}}{\text{Evidence}}$$

7.3 Fundamental Bayes' Probability Space

We begin by showing the probability space via a Venn diagram visualization.

7.3.1 Probability Space Venn Diagram for traditional Bayesian

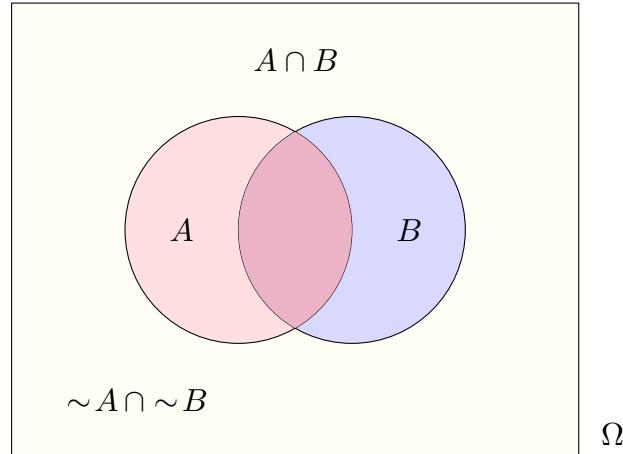


Figure 7.1: Venn diagram representing the intersection and complements of sets A and B .

The Venn diagram depicts a Bayes' probability space, with sets A and B as distinct events within the universal event space Ω . The area $A \cap B$ illustrates the joint occurrence of events A and B , while $\sim A \cap \sim B$ denotes the region where neither event takes place. This visualization aids in understanding Bayesian inference, particularly how evidence B updates the probability of event A according to Bayes' Theorem.

Identities derived by inspection:

$$\begin{aligned} P(A) &= P(A \cap B) + P(A \cap \sim B) \\ P(\sim A) &= P(\sim A \cap B) + P(\sim A \cap \sim B) \\ P(B) &= P(A \cap B) + P(\sim A \cap B) \\ P(\sim B) &= P(A \cap \sim B) + P(\sim A \cap \sim B) \end{aligned}$$

Fundamental Axioms:

$$\begin{aligned} \{P(A), P(B)\} &\in [0, 1] \triangleq [0, \Omega] \\ P(A \cap B) &= P(B \cap A) \\ P(\overline{A \cap B}) &= \Omega + P(A \cap B) - P(A) - P(B) \end{aligned}$$

7.4 Kolmogorov Axioms of Probability

The Kolmogorov axioms {non-negativity, normalization, and additivity} provide a way to derive a general expression for conditional probability, forming the basis of Bayes' theorem. We want a way to efficiently calculate Bayesian probability updates. This is a derivation grounded in the axiomatic structure of probability.

7.4.1 The Kolmogorov Axioms

1. **Non-negativity:** $P(E) \geq 0$ for any event E .
2. **Normalization:** $P(\Omega) = 1$, where Ω is the sample space.
3. **Additivity:** If A and B are disjoint events ($A \cap B = \emptyset$), then:

$$P(A \cup B) = P(A) + P(B).$$

7.4.2 Definition of Conditional Probability

Conditional probability is defined as:

$$P(A | B) = \frac{P(A \cap B)}{P(B)} \quad \text{if } P(B) > 0.$$

This definition can be derived systematically using the axioms of probability and set theory.

7.4.3 Axiomatic Derivation of $P(A | B)$

Step 1: Define the restricted sample space

When conditioning on event B , we are effectively restricting the sample space to B . In this restricted sample space: - The total probability of B becomes 1, i.e., $P_B(\Omega) = 1$, where P_B denotes the new probability measure conditioned on B .

Step 2: Normalize probabilities within B

The probability of any event A , given B , should reflect the proportion of A within B . This leads to the proportionality:

$$P(A | B) \propto P(A \cap B).$$

To normalize this proportion, divide by the total probability of B :

$$P(A | B) = \frac{P(A \cap B)}{P(B)}.$$

Step 3: Justify using Kolmogorov's axioms

This definition satisfies the axioms of probability:

1. **Non-negativity:** Since $P(A \cap B) \geq 0$ and $P(B) > 0$, $P(A | B) \geq 0$.
2. **Normalization:** The total probability in the restricted space sums to 1:

$$P(B | B) = \frac{P(B \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1.$$

3. **Additivity:** For disjoint events A_1, A_2 , where $A_1 \cap A_2 = \emptyset$:

$$P((A_1 \cup A_2) | B) = \frac{P((A_1 \cup A_2) \cap B)}{P(B)}.$$

By the additivity axiom:

$$P((A_1 \cup A_2) \cap B) = P((A_1 \cap B)) + P((A_2 \cap B)),$$

so:

$$P((A_1 \cup A_2) | B) = \frac{P(A_1 \cap B)}{P(B)} + \frac{P(A_2 \cap B)}{P(B)} = P(A_1 | B) + P(A_2 | B).$$

7.4.4 Axiomatic Derivation of $P(B | A)$

By symmetry of intersection ($A \cap B = B \cap A$), the same reasoning applies:

$$P(B | A) = \frac{P(B \cap A)}{P(A)} \quad \text{if } P(A) > 0.$$

7.4.5 Summary of the Axiomatic Derivation

1. Start from the definition of probability within a restricted sample space (conditioning).
2. Use normalization within the restricted space (B) to derive $P(A | B)$.
3. Verify that the result satisfies Kolmogorov's axioms (non-negativity, normalization, additivity).
4. Symmetry of intersection ensures the analogous formula for $P(B | A)$.

Thus:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \quad P(B | A) = \frac{P(A \cap B)}{P(A)}.$$

These results naturally follow from the probability axioms and set theory.

7.4.6 The Role of Commutativity in Bayesian Probability Calculations

The salient question arising at this juncture is:

How do we relate the Venn diagram in Figure: 7.1 to Bayes' theorem?

And, how do we calculate using theorem?

In set theory:

$$A \cap B = B \cap A, \quad \text{and} \quad A \cup B = B \cup A.$$

These properties reflect the invariance of the resulting set under the order of operation. Similarly, in a metric space, where probability is conceptualized as a measure P , the measure of the intersection ($A \cap B$) is invariant under the order of A and B . This invariance is the cornerstone that allows Bayes' relation to hold up under all valid probability distributions. This property allows us to express the joint probability of two events A and B in multiple equivalent ways:

$$P(A \cap B) = P(B \cap A).$$

This symmetry is not arbitrary; it arises directly from the axioms of set theory. Specifically, the intersection operation in sets is commutative, meaning the order of the operands does not affect the result. In a probabilistic context, this translates to the joint occurrence of A and B being independent of which event is referenced first.

From a metric perspective, probability can be thought of as a measure defined on a ***sigma-algebra of subsets**. The intersection of two subsets corresponds to the overlap of their measurable areas. The measure of this overlap does not depend on the order in which the subsets are combined.

This commutativity directly enables the derivation of Bayes' theorem. The theorem hinges on the equivalence of two different factorizations of the joint probability $P(A \cap B)$:

$$P(A \cap B) = P(A)P(B | A) = P(B)P(A | B)$$

Without the commutative property, such equivalence could not be guaranteed. This would undermine the consistency of conditional probability calculations and, by extension, the entire framework of Bayesian analysis.

7.4.7 Computational Implications for Bayesian Analysis

The commutative property simplifies Bayesian calculations by providing a consistent framework to calculate conditional probabilities. Specifically, we can use:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \quad \text{and} \quad P(B | A) = \frac{P(A \cap B)}{P(A)}.$$

These forms are essential in computational Bayesian analysis because they allow direct computation of conditional probabilities using readily available components: the joint probability and the marginal probabilities.

By relying on commutativity, we avoid ambiguity and reduce the complexity of Bayesian inference to manageable and consistent operations. This makes Bayesian analysis not only theoretically sound but also practically implementable in diverse applications such as machine learning, statistical inference, and decision-making under uncertainty. The commutativity of intersection streamlines Bayesian probability calculations, reducing them to the ratio of the events' intersection to the total probability of A or B . This pivotal concept underpins Bayesian analysis.

$$P(B)P(A | B) = P(A \cap B) = P(B \cap A) = P(A)P(B | A)$$

This consideration streamlines the computation of Bayesian probabilities. By conceptualizing it as the ratio of the intersection of events A and B to the total probability of A or B , one can readily calculate the conditional probabilities inherent in Bayes' theorem.

Thus, we arrive at essential expressions for Bayesian inference computation:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

7.5 Probability Space Graph over Ω

Interestingly Enough - these relations yield a probability symmetry originating from Ω .

Interestingly enough, the relationships visualized in this graph yield a probability symmetry originating from Ω . This graph is meaningful because it provides a structured way to represent the relationships between various conditional probabilities in the probability space Ω . By mapping these relationships, we gain insights into the fundamental symmetries and connections that define conditional probability.

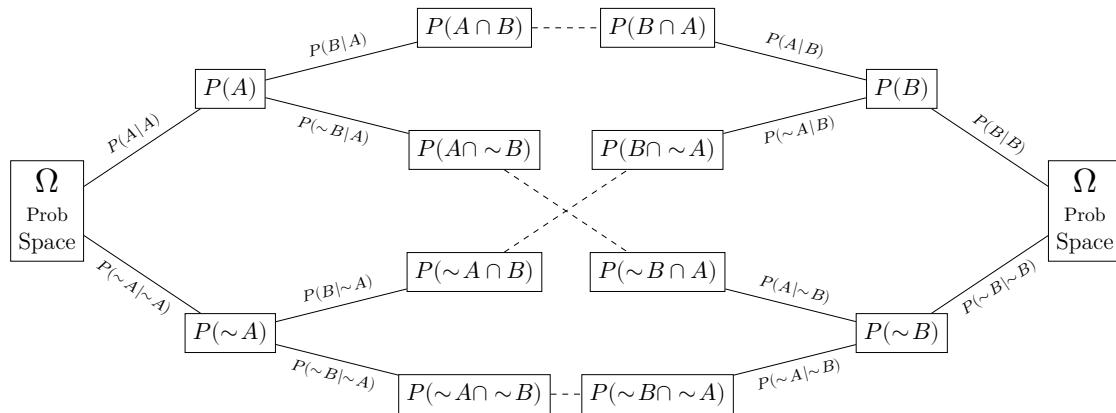


Figure 7.2: Probability Graph Ω bridged by commutative intersection $P(A \cap B) = P(B \cap A)$

7.5.1 Why is this graph meaningful?

This type of visualization allows us to:

- Understand how conditional probabilities link different events.
- Highlight the inherent symmetry in relationships such as $P(A \cap B) = P(B \cap A)$, emphasizing commutative properties.
- Clarify how joint probabilities, complements, and conditional probabilities interact within the space Ω .

Each **edge** in the graph represents a **conditional probability** that connects adjacent nodes.

For example:

- An edge from Ω to $P(A)$ trivially represents $P(A | A)$, the probability of A given A .
- Similarly, edges like $P(B | A)$ connect events $P(A)$ and $P(A \cap B)$, reinforcing the dependency between them.

Through graph re-ordering we can better highlight inherent symmetries around Ω . We can then better understand how these conditional probabilities form a consistent framework that governs the relationships between different events.

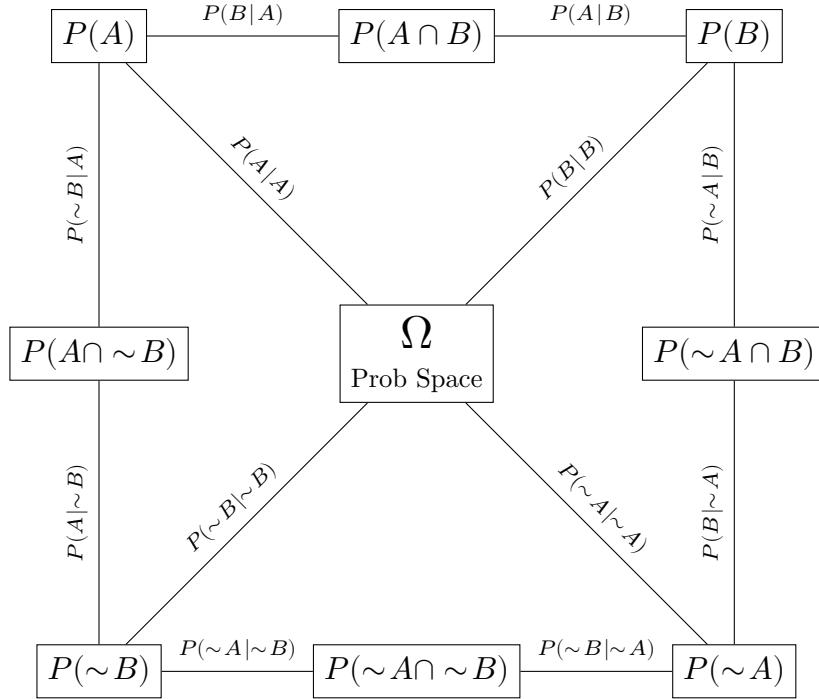


Figure 7.3: Probability Graph displaying inherent symmetry around Ω

Probability theory provides a robust framework for understanding and quantifying uncertainty. At the heart of this framework lies the concept of a probability space, represented here as Ω . The mappings within this space, as visualized in the following graphs, reveal deep insights into how probabilities interact, connect, and relate to each other.

This visualization captures the relationships between various events and their intersections, along with the conditional probabilities that bridge these relationships. Each node in the graph represents a specific probability or set of probabilities, while the edges linking them denote conditional probabilities, which express the likelihood of one event given the occurrence of another. By structuring these connections graphically, we can achieve a more intuitive grasp of the mathematical underpinnings of probability theory.

Why This Graph is Meaningful ?

Symmetry and Commutative Properties:

The graph highlights a fundamental property of probability: the commutative nature of intersections. For instance, $P(A \cap B) = P(B \cap A)$. This symmetry is clearly visualized through the commutative links, represented as dashed lines.

Conditional Dependencies:

Each edge encodes a conditional probability, such as $P(A | B)$ or $P(B | A)$. These conditional probabilities provide critical insights into how one event influences another. The graph captures this influence in a structured and visually accessible way.

Partitioning of Probability Space:

The diagram divides Ω into mutually exclusive and exhaustive components, such as $P(A \cap B)$, $P(A \cap \sim B)$, $P(\sim A \cap B)$, and $P(\sim A \cap \sim B)$. This partitioning reveals how probabilities distribute across the space and helps in understanding joint and marginal probabilities.

Inherent Symmetry:

The visualization emphasizes the inherent symmetry in probability relationships. For example, the mappings around Ω are visually balanced, reinforcing the idea that Ω serves as the universal set encompassing all events.

Practical Interpretations:

Such mappings are particularly useful in Bayesian inference, where updating beliefs based on observed data requires careful computation of conditional probabilities. This graph provides a foundation for such calculations by making the relationships explicit.

Insights Gained from the Visualization

Connectivity:

The direct edges between nodes illustrate how probabilities propagate through the system. For example, $P(A)$ connects to $P(A \cap B)$ via the conditional probability $P(B | A)$. This connectivity mirrors real-world decision-making processes where one piece of information influences another.

Comparative Understanding:

By juxtaposing related probabilities, such as $P(A \cap B)$ and $P(B \cap A)$, the graph enables a side-by-side comparison that reinforces mathematical equivalence.

Structure of Relationships:

The graph imposes a structure on the probability space, making abstract relationships more tangible. This structured view is invaluable for teaching, analysis, and practical applications such as decision analysis and machine learning.

How to Use This Graph

Study Conditional Probabilities:

Trace the edges to understand the relationship between nodes. For example, $P(A \cap B)$ is linked to $P(A)$ through $P(B | A)$, showing the dependence of the joint probability on the

marginal and conditional probabilities.

Analyze Bayesian Inference:

Use the graph to identify how evidence updates beliefs. For instance, given $P(B | A)$ and $P(A)$, one can compute $P(A \cap B)$, which forms the basis of Bayesian updates.

Visualize Complex Spaces:

For advanced topics, such as joint distributions or higher-dimensional spaces, extend the graphing technique to include additional events and their intersections.

In summary, this probability space graph provides an essential tool for exploring and understanding the relationships between events in a probabilistic framework. By visualizing Ω and its components, we gain not only a mathematical appreciation but also a conceptual understanding of the intricate dance of probabilities.

7.6 Maringal and Bayes' probability space sums

In Bayesian analysis, the interrelation of event probabilities and their complements constructs a full probability model. The following table clarifies this by mapping the joint probabilities of events A and b along with their complements in the probability space Ω . It exemplifies the key tenet that all outcome probabilities sum to unity, aligning with Ω . The table shows the network for calculating marginal probabilities and grasping Bayesian inference's core, which will be demonstrated in subsequent examples. Note: While we can normally employ $\Omega = 1$, we can also commonly allow some value like $\Omega = 750$ and then the numbers sum to this value.

	$\{A\}$	$\{\sim A\}$	Marginal
$\{B\}$	$P(A \cap B)$	$+ P(\sim A \cap B)$	$= P(B)$
$\{\sim B\}$	$+ P(A \cap \sim B)$	$+ P(\sim A \cap \sim B)$	$= P(\sim B)$
Marginal	\parallel	\parallel	\parallel
	$P(A)$	$+ P(\sim A)$	$= \Omega$

Figure 7.4: Marginals and sums network of Traditional Bayes' Space

7.7 bayes' by example

7.7.1 Classical Bayes' word problem example (“Red Ace Example”)

Imagine you have a deck of cards, and you want to know the probability that a card is an ace given that it's a red card. Here, the event A is “the card is an ace,” and B is “the card is red.”

- The prior probability $P(A)$ is $4/52$ since there are 4 aces in a 52-card deck.
- The likelihood $P(B | A)$ is $1/2$ because half of the aces are red.
- The marginal probability $P(B)$ is $26/52$ since there are 26 red cards in the deck.

Plugging these values into Bayes' Theorem:

So, the probability that a card is an ace given that it's red is $1/13$.

Bayes' Theorem thus provides a robust framework for updating beliefs and making informed decisions based on new data. It's a cornerstone of statistical inference and finds applications in various fields, from spam filtering in emails to disease diagnosis in medicine

7.7.2 Considering the “Classic Example” and why we can do better

To candidly address the challenges inherent in understanding probability spaces and the diverse applications of Bayes' Theorem, it is essential to delve beyond the surface of typical examples. By doing so, we aim to unravel the deeper intricacies of Bayes' formula. Our objective here is not just to reiterate the theorem but to derive a series of identities and explore the concept of probability metric spaces.

By fostering a profound comprehension from foundational principles and underlying identities, we intend to cultivate a structured approach to employing Bayes' Theorem in various scenarios. This is particularly vital in contexts where the traditional formulation of Bayes' may not be the most appropriate or insightful application. Through this exploration, we aspire to formulate more nuanced and, ostensibly, more compelling scenarios that demonstrate the versatility and power of Bayes' Theorem.

In this journey, we will dissect the theorem, examining each component and its role in shaping the overall framework. We'll explore how prior beliefs are updated with new evidence, and how this process is influenced by the underlying probability structures. By doing so, we aim to provide readers with a toolkit for thinking critically about probability and making more informed decisions in the face of uncertainty.

Ultimately, this exploration is not just about understanding a mathematical formula; it's about embracing a new way of thinking. It's about shifting our perspective from seeing Bayes' Theorem as a mere tool to appreciating it as a fundamental principle that underpins a vast array of applications—from scientific research to everyday decision-making.

7.7.3 A visual understanding of Bayes

7.7.4 Deriving the Rules of Probability

To derive the rules of probability, consider the general case involving two random variables (X) and (Y).

7.7.5 Example Context

In a medical example:

- (X) could represent the **presence or absence of cancer**
- (Y) could represent the **outcome of a diagnostic test**

These variables vary between individuals in generally unknown ways, so they are referred to as **random variables** (or **stochastic variables**).

7.7.6 Variable Definitions

We assume:

- (X) can take values (x_i), where $i = 1, \dots, L$
- (Y) can take values (y_j), where $j = 1, \dots, M$

7.7.7 Observational Setup

Suppose there are N total trials in which we sample both (X) and (Y). Define:

- (n_{ij}): the **number of trials** where ($X = x_i$) and ($Y = y_j$)
- (c_i): the **number of trials** where ($X = x_i$) (irrespective of (Y))
- (r_j): the **number of trials** where ($Y = y_j$) (irrespective of (X))

7.7.8 Joint Probability

The **joint probability** $p(X = x_i, Y = y_j)$ is the probability that:

- (X) takes the value (x_i),
- and (Y) takes the value (y_j).

It is calculated as the fraction of trials (n_{ij}), where ($X = x_i$) and ($Y = y_j$), divided by the total number of trials (N):

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}.$$

This represents the proportion of observations falling into the corresponding cell (i, j) in a contingency table.

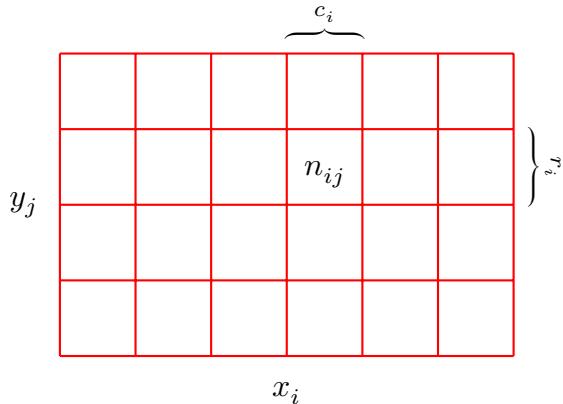


Figure 7.5: Red Grid for Bayes Calc (edit later)

7.8 Sum and Product Rules of Probability

To derive the **sum and product rules of probability**, consider the following:

- We have two **random variables**:
- X , which takes values $\{x_i\}$ where $i = 1, \dots, L$ – Y , which takes values $\{y_j\}$ where $j = 1, \dots, M$
- For this illustration:
 - $L = 5$
 - $M = 3$

7.8.1 Definitions:

- **Total instances:** Let N denote the total number of instances of these variables.
- **Joint instances:** The number of instances where $X = x_i$ and $Y = y_j$ is n_{ij} . This corresponds to the **number of instances in the (i, j) cell of the array**.
- **Column totals:**
 - The total number of instances in column i , corresponding to $X = x_i$, is c_i .
- **Row totals:**
 - The total number of instances in row j , corresponding to $Y = y_j$, is r_j .

7.8.2 Summary of Relationships:

$$\sum_{j=1}^M n_{ij} = c_i \quad (\text{Sum over rows for each column } i)$$

$$\sum_{i=1}^L n_{ij} = r_j \quad (\text{Sum over columns for each row } j)$$

7.8.3 Observations:

- The joint frequency n_{ij} determines how many instances satisfy $X = x_i$ and $Y = y_j$.
- Marginal counts c_i and r_j are computed as the **sum of joint frequencies** over their respective rows or columns.

7.9 Connecting Bayes' to Naive Bayes'

In the context of extending your original Bayesian framework to Naive Bayes, the key modification comes from how the events A and B are treated when there is an assumption of conditional independence between features in a classification context.

In Naive Bayes, the assumption is that the features (or events) are conditionally independent given the class. This assumption drastically simplifies the computation of the posterior probability, making it computationally efficient for tasks like classification. Here's how the core concepts from traditional Bayesian inference are adapted in Naive Bayes:

Naive Bayes: Conditional Independence Assumption In traditional Bayesian analysis, you might compute the probability of a set of features given a class using a general conditional probability like:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Where A and B represent some events. However, in the context of Naive Bayes, if A_1, A_2, \dots, A_n are features and B is the class, the assumption of conditional independence leads to:

$$P(B | A_1, A_2, \dots, A_n) = \frac{P(B) \prod_{i=1}^n P(A_i | B)}{P(A_1, A_2, \dots, A_n)}$$

This assumption simplifies the likelihood term $P(A_1, A_2, \dots, A_n | B)$ as the product of individual conditional probabilities:

$$P(A_1, A_2, \dots, A_n | B) = \prod_{i=1}^n P(A_i | B)$$

This is what makes Naive Bayes “naive” - it assumes that each feature (or piece of evidence) contributes independently to the probability of the class.

7.10 Raw Bayes Text Here First

7.11 Modifications in Venn Diagram Representation

1. Traditional Bayesian Space: The Venn diagram you provided earlier represents events A and B within the universal space Ω . Each event may overlap or be disjoint, and their interactions are captured through their joint, marginal, and conditional probabilities.
2. Naive Bayes Representation: The key modification is that, in the Naive Bayes model, we assume that each feature A_1, A_2, \dots, A_n is conditionally independent given the class B . This simplification means the probability space is represented differently, potentially with separate Venn diagrams for each feature given B .

In the case of Naive Bayes, you can think of the events A_1, A_2, \dots, A_n as having their own Venn diagrams, where each feature A_i contributes independently to the classification decision. The overlap between these features and B is still shown in the diagram, but the assumption of conditional independence implies that these overlaps do not interact as they would in a more general Bayesian framework.

Modified Probability Space for Naive Bayes For Naive Bayes, we are interested in the posterior probability of the class B given the features A_1, A_2, \dots, A_n :

$$P(B | A_1, A_2, \dots, A_n) = \frac{P(B) \prod_{i=1}^n P(A_i | B)}{P(A_1, A_2, \dots, A_n)}$$

This equation breaks down into: 1. Prior: $P(B)$ - the prior probability of class B . 2. Likelihood: $\prod_{i=1}^n P(A_i | B)$ - the product of the conditional probabilities for each feature given the class. 3. Normalization: $P(A_1, A_2, \dots, A_n)$ - the likelihood of the observed features across all classes, used for normalization.

7.12 Modifications to Your Original Bayesian Framework

The Venn diagram for Naive Bayes will primarily have distinct regions where each feature A_1, A_2, \dots, A_n is considered conditionally independent of each other, but dependent on the class B . These dependencies are shown as separate slices or parts of the overall space, but they do not overlap in the way traditional Bayesian probability might depict interactions between events. 1. The Venn diagram of the original Bayesian framework might have complex intersections between events A and B , as they can influence each other. 2. In the Naive Bayes model, the overlap between events is reduced due to the assumption of conditional independence. Each feature is considered independently given the class, simplifying the model.

Thus, while the structure of the probability space in both models follows a similar form, the key difference is how the conditional independence assumption in Naive Bayes simplifies the joint probability computation. Instead of modeling the full interaction between features, Naive Bayes reduces it to the product of individual feature probabilities conditioned on the class.

7.13 Naive Bayes' as Filter

Yes, Naive Bayes can be thought of as a filter in certain contexts, especially in classification tasks. The idea behind this is that Naive Bayes filters out irrelevant features (or at least gives them less weight) due to its assumption of conditional independence.

Key Aspects of Naive Bayes as a Filter:

1. Feature Relevance: Naive Bayes classifies by calculating the probability of each class given the features. However, because of the conditional independence assumption, it treats each feature independently and does not model interactions between features. As a result, irrelevant features that do not help discriminate between classes will not have a significant impact on the classification. In this sense, Naive Bayes inherently “filters out” noisy features that don’t provide discriminative power.
2. Probabilistic Thresholding: Naive Bayes computes the posterior probability of each class using Bayes’ theorem, and the class with the highest probability is chosen. Features that are more relevant to the class have a stronger influence on the posterior probability. Features that are less relevant have a smaller influence because their conditional probability given the class is likely to be uniform or less informative. This allows Naive Bayes to “filter” out noisy features by giving them a lower weight in the final decision.
3. Handling Redundant Features: Because Naive Bayes treats each feature as independent, it does not take into account any redundant or correlated features (which may be a source of noise in other classification algorithms). In cases where features are highly correlated, Naive Bayes essentially treats them separately, and the model does not suffer as much from overfitting due to feature redundancy. This simplification can act as a form of regularization, which also helps in “filtering out” the unnecessary complexity that might arise from correlated features.

Naive Bayes in Practice - Text Classification: In tasks like spam detection, Naive Bayes acts as a filter by giving higher probabilities to features (e.g., certain words) that are more indicative of spam or ham (non-spam). The algorithm does not consider the relationship between words but simply filters through each word's likelihood in the context of spam and ham categories. - Feature Selection: In situations where you have a large number of features, Naive Bayes naturally reduces the influence of irrelevant features. While it's not a feature selection method per se (like some others that explicitly remove unimportant features), its simplification (treating features independently) effectively reduces the impact of less important features.

Limitations of Naive Bayes as a Filter: While Naive Bayes does perform a form of feature weighting implicitly, it's important to note that it's not a filtering algorithm in the strict sense, like some other models (e.g., Decision Trees, Random Forests, or feature selection algorithms) that explicitly discard features. Naive Bayes doesn't eliminate features or explicitly ignore them; instead, it simply assigns lower probabilities to features that don't significantly contribute to class distinctions.

In Summary: Naive Bayes can act as a filter in that it implicitly “filters out” irrelevant or less relevant features by not considering the dependencies between them and by giving less weight to features that don't distinguish well between classes. However, it doesn't perform explicit feature selection or removal like other models that focus on filtering features based on their importance.

7.14 Example to classical Bayes' vs. Naive Bayes'

In the scenario presented in the image, the problem involves a medical test for a condition (illness), with the concern being a false positive (i.e., when the test indicates illness, but the person is actually healthy). Let's break down how Naive Bayes could be used to address this problem, step by step.

Key Components of the Scenario:

Events:

- A : The person has the illness (cancer).
- $\sim A$: The person is illness-free.
- B : The test result is positive.
- $\sim B$: The test result is negative.

Given Data:

- $P(A | B)$: The probability that the person has the illness given a positive test result.
- $P(\sim A | B)$: The probability that the person is illness-free given a positive test result.
- The contingency table is used to provide the necessary data for calculating these probabilities.

From the image, we have the following values from the contingency table:

Test Result	Illness (A)	Iliness-Free ($\sim A$)
Positive	8	103
Negative	2	895

So:

- Total number of tests = $8 + 2 + 103 + 895 = 1008$.
- Number of true positives (illness and positive test, $A \cap B$) = 8.
- Number of false positives (illness-free and positive test, $\sim A \cap B$) = 103.
- Number of true negatives (illness-free and negative test, $\sim A \cap \sim B$) = 895.
- Number of false negatives (iliness and negative test, $A \cap \sim B$) = 2.

7.15 USing Bayes and the tie to physics

7.15.1 Physics Analogies and Bayesian Reasoning Models

A key concept in information theory, which yields its implied relationship to physics, is the average value (or expectation) of the information content of events (symbols arriving at the end of a communication line):

$$H(x) = - \sum_x P(x) \log_2 P(x)$$

This symbol $H(x)$ is often called information entropy, drawing an analogy to physics (specifically thermodynamics), where entropy is a measure of disorder, particularly the number of specific ways that a set of particles may be arranged in space. As Feynman (2000) illustrated, we can understand this by considering a hypothetical situation A : a volume containing a single molecule of gas with X equally likely locations. In this uniform distribution across X locations, the information needed to describe any possible location is given by:

$$H_A(x) = - \sum_x P(x) \log_2 P(x) = -x \left(\frac{1}{x} \right) \log_2 \left(\frac{1}{x} \right) = -\log_2 \left(\frac{1}{x} \right)$$

Thus, any message indicating the particle's location has this information content. While this scenario assumes uniform distribution, one can generalize this approach to compute the average information content across any distribution.

Now, consider performing work to compress the gas, halving the original volume (situation B). Here, half the terms in the sum are zeroed out, the remaining probabilities double, and the bits required to convey the particle's location are given by:

$$H_B(x) = -\log_2 \left(\frac{1}{2x} \right) = -\log_2 \left(\frac{1}{x} \right) - 1 = H_A(x) - 1$$

In this way, the volume compression effectively reduces the information entropy (expected information content) of messages indicating the particle's location by one bit.

7.15.2 Extending the Analogy to Biologics Lot-to-Lot Variance

In biologics, particularly in the production of drug substances, there is an analogous concept where we can think of “information” as the variability or uncertainty in quality attributes across different production lots. Consider the production of several biologics lots, where we measure attributes such as purity, potency, and concentration. Each lot can be thought of

as a unique “event” with an associated probability distribution, akin to particle locations in thermodynamics.

For small set sizes, such as when $n = 3$ to 5 , we face higher uncertainty and variance in these quality attributes across lots. Information entropy, in this case, represents the uncertainty in predicting the quality attributes for future lots based on initial measurements. The entropy $H(x)$ of these quality attributes might indicate how “disordered” or variable our lot-to-lot outcomes are under initial conditions.

Imagine now performing a series of adjustments after each production run. By incorporating learnings from each new lot (e.g., 1 through 5), we might make adjustments to the process parameters to reduce variability. This is analogous to “compressing” the space in which our attributes vary, akin to reducing the gas volume. For example, if we find that a specific pH adjustment improves consistency, we may decrease the uncertainty (information entropy) about the quality attributes of subsequent lots. The expected entropy of future lots 6, 7, 8, etc., becomes:

$$H_{\text{new}}(x) = H_{\text{initial}}(x) - \Delta H$$

where ΔH represents the reduction in entropy due to process optimization and learnings from previous lots.

Thus, the iterative nature of Bayesian reasoning in this context reflects how biologics manufacturing often relies on sequential data collection and process optimization. Just as reducing a gas’s volume constrains the possible locations for particles, process adjustments reduce the variability in quality attributes, thereby lowering the information entropy and enhancing the predictability of future lot outcomes.

In practice, this iterative approach of learning and adjustment can reduce the “noise” in biologic production and improve consistency, making the analogy between thermodynamic entropy and biologics lot-to-lot variance a powerful conceptual tool.

7.16 Using Naive Bayes

Naive Bayes assumes conditional independence between features (in this case, the presence of illness and the test result). However, since we are just using the test result as evidence for illness or not, the Naive Bayes classifier simplifies to using Bayes' theorem directly.

The posterior probability of having the illness given a positive test result is given by Bayes' theorem:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Where: - $P(A | B)$: Probability of having the illness given a positive test result. - $P(B | A)$: Probability of testing positive given that the person has the illness (True Positive Rate). - $P(A)$: Prior probability of having the illness (based on the population's illness rate). - $P(B)$: Probability of testing positive, which is the total probability of a positive test (True Positives + False Positives).

Step-by-Step Calculation: 1. Prior Probability $P(A)$: This is the probability of having the illness in the population:

$$P(A) = \frac{\text{Number of illness cases}}{\text{Total number of tests}} = \frac{8 + 2}{1008} = \frac{10}{1008} \approx 0.0099$$

2. Likelihood $P(B | A)$: This is the probability of testing positive given the person has the illness:

$$P(B | A) = \frac{\text{True Positives}}{\text{Total illness cases}} = \frac{8}{10} = 0.8$$

$$P(B | A) = \frac{\text{True Positives}}{\text{Total illness cases}} = \frac{8}{10} = 0.8$$

3. Evidence $P(B)$: This is the total probability of a positive test result:

$$P(B) = \frac{\text{True Positives} + \text{False Positives}}{\text{Total number of tests}} = \frac{8 + 103}{1008} = \frac{111}{1008} \approx 0.1102$$

4. Posterior Probability $P(A | B)$: Now, applying Bayes' theorem:

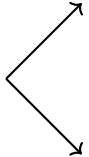
$$P(A | B) = \frac{P(B | A)P(A)}{P(B)} = \frac{(0.8)(0.0099)}{0.1102} \approx 0.0717$$

Thus, the probability that a person has the illness given a positive test result is approximately 7.17%.

Conclusion: Even though the test is positive, the low prior probability of the illness (0.99%) and the high false positive rate (103 out of 895 illness-free people test positive) result in a relatively low posterior probability for the person actually having the illness. This demonstrates how Naive Bayes uses the data and Bayes' theorem to provide a more accurate estimate of the illness probability, given the test results.

7.17 after area

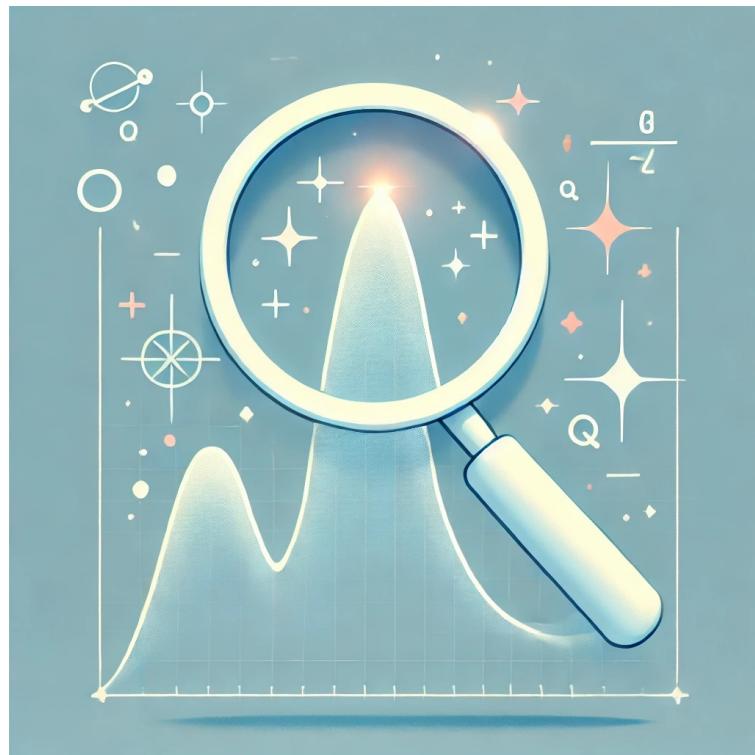
7.18 Bayes start building



Equation 8

Maximum Likelihood

$$\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^n P(x_i | \theta)$$



Parametric based Maximum

Key ML Equation 8: Maximum Likelihood Estimation

$$\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^n P(x_i | \theta) \quad (8.0.1)$$

x_i	The observed data point or the i -th element of the dataset, which the likelihood function evaluates based on the parameter θ . More specifically belongs to a set of data namely $x_i \in \{x_1, x_2, \dots, x_n\} = X$.
θ	The parameter(s) to be estimated using MLE. These can represent different characteristics of the distribution (e.g., mean, variance). For the illustration here we are only considering one value for θ . More generally belongs to a set of data namely $\theta_i \in \{\theta_1, \theta_2, \dots, \theta_m\} = \Theta$.
$\hat{\theta}$	The MLE estimate of the parameter θ , representing the value that maximizes the likelihood function.
$P(x_i \theta)$	The likelihood function evaluated for the data point x_i given the parameter θ . It represents the probability of observing x_i under the assumed parameter.
\prod	The product operator , which combines the likelihoods of all observed data points $x_i \in \{x_1, x_2, \dots, x_n\} = X$ to compute the overall likelihood of the dataset.
$\arg \max_{\theta}$	The argument of the maximum , which says to return the value of θ that maximizes the likelihood function.

Introduction: MLE is a method used to estimate the parameters of a statistical model.

Description: It finds the parameter values Θ that maximize the likelihood of making the observations given the model $P(x_i | \theta)$, and a specific data set X .

Importance in ML: MLE is used to train many machine learning models, including logistic regression and Gaussian mixture models. It provides a way to fit models to data and make statistical inferences.

8.1 What is Behind the Equation

The Maximum Likelihood Estimation equation, while compact in its representation, embodies several fundamental statistical concepts and operational components that work together to find optimal parameter estimates. Let's break down its key elements and underlying principles:

Core Components of the MLE

1. Probability Model Structure

- The equation centers on $P(x_i|\theta)$, which represents our assumed probabilistic model.
- This model defines how we believe our data is probabilistically represented.
- The vertical bar “|” indicates conditional probability: “probability of data x_i given parameter θ ”.

2. Independence Assumption

- We assume, that the product \prod implies that each probability $P(x_i|\theta)$ for a given x_1 observation is independent.
- This allows the benefit to multiply individual probabilities of each of the n observations: $P(X|\theta) = P(x_1|\theta) \times P(x_2|\theta) \times \dots \times P(x_n|\theta)$.

3. Optimization Framework

- The operator $\arg \max_{\theta}$ indicates we're searching for the best parameter value.
- This search spans the entire parameter space for θ .
- We seek the global maximum of the likelihood function with respect to θ .

Mathematical Properties

1. Monotonic Transformation

- Often converted to log-likelihood for computational efficiency leveraging the identity property of logs via $\log(\prod P(x_i|\theta)) = \sum \log P(x_i|\theta)$.
- This transforms products into sums while preserving the maximum value at some θ

2. Derivatives and Critical Points

- Finding maximum often involves setting derivatives to zero $\frac{\partial}{\partial \theta} \log L(\theta) = 0$.
- Second derivatives confirm maximum vs. minimum via the Hessian.

Statistical Implications

1. Sample Size Effects

- Larger n typically leads to more accurate estimates.
- Product terms count grows with each additional observation.

- Asymptotic properties become more relevant with large n.

2. Parameter Space

- θ can be single-valued or multi-dimensional
- Parameter constraints affect the optimization
- Solution must lie within valid parameter space

Practical Considerations

1. Computational Aspects

- Numerical underflow can occur with large products
- Log transformation helps prevent computational issues
- Optimization algorithms may need constraints

2. Model Selection Impact

- Choice of probability model affects likelihood structure
- Different models lead to different estimation procedures
- Model assumptions must be validated. This equation thus represents a sophisticated interplay between probability theory, optimization methods, and statistical inference, providing a systematic framework for parameter estimation in statistical modeling and machine learning applications.

8.2 Historical Context

The method of Maximum Likelihood Estimation was first formally introduced by **Ronald A. Fisher** in his 1922 paper “*On the Mathematical Foundations of Theoretical Statistics*”. Fisher’s groundbreaking work laid the foundation for modern statistical theory, and MLE quickly became one of the most widely adopted estimation methods. However, the ideas behind MLE can be traced back even further to **Carl Friedrich Gauss** and **Pierre-Simon Laplace**, who used related concepts in the context of least squares fitting and probability theory.

The modern form of MLE was developed by **Ronald A. Fisher** in his groundbreaking 1921 paper, “*On the Mathematical Foundations of Theoretical Statistics*”. Fisher introduced MLE as a method to estimate parameters of a statistical model by maximizing the likelihood of observing the data given those parameters. Fisher’s work formalized MLE and established its theoretical underpinnings, including the idea that MLE estimators are asymptotically unbiased, efficient, and consistent under regular conditions. Fisher’s development of MLE was revolutionary because it provided a clear, consistent framework for estimating model parameters that was both flexible and powerful. Over the years, MLE has found widespread

use in many statistical applications, ranging from simple regression models to complex models in machine learning and bioinformatics.

Impact and Adoption

- Fisher's ideas were revolutionary because they provided a unified framework for parameter estimation that was both robust and flexible, applicable to many different types of statistical models.
- MLE quickly became the cornerstone of statistical inference, replacing older methods like method of moments or least squares in many applications.

Further Developments

- In the decades following Fisher's work, **Bayesian methods** and **generalized likelihood approaches** expanded on MLE, but it remained a dominant method in frequentist statistics.
- Modern developments have extended MLE to complex models and high-dimensional data, incorporating computational techniques like the **expectation-maximization (EM) algorithm** for handling missing data and latent variables.

MLE was first formally developed by Ronald Fisher in the 1920s, building on earlier work by Gauss and Laplace. It became a foundational tool in statistics due to its general applicability and strong theoretical properties, such as consistency and efficiency. Over the years, it has remained a key method in both classical and modern statistical modeling.

8.3 MLE General Introduction

Maximum Likelihood Estimation (MLE) is one of the most widely used and fundamental methods in statistical inference. It provides a systematic approach for estimating the parameters of a probabilistic model based on observed data. At its core, MLE seeks to find the parameter values that make the observed data most likely, given a specific model. This method plays a pivotal role in various fields, including regression analysis, bioinformatics, machine learning, and econometrics, where it is frequently used to **fit models to data** and make predictions.

What is Maximum Likelihood Estimation?

In simple terms, **Maximum Likelihood Estimation** aims to determine the parameter values of a statistical model that would make the observed data as probable as possible. Suppose we have a statistical model parameterized by a set of parameters θ . Given some observed data $X = (x_1, x_2, \dots, x_n)$, the likelihood function $\mathcal{L}(\theta)$ expresses the probability of observing this data as a function of the model parameters θ . The goal of MLE is to find the parameter values $\hat{\theta}$ that maximize this likelihood function:

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta | X)$$

This can be thought of as searching for the set of parameters under which the model would most likely generate the data we have observed. The value of $\hat{\theta}$ that achieves this maximization is called the **Maximum Likelihood Estimate**.

The Likelihood Function

The likelihood function is central to MLE. For a given set of parameters θ , it calculates the probability of observing the data set X . For independent and identically distributed (i.i.d.) random variables, the likelihood is the product of the individual probabilities (or probability density functions). For example, if the data $X = (X_1, X_2, \dots, X_n)$ follows a probability distribution $f(X_i|\theta)$, the likelihood function is:

$$\mathcal{L}(\theta|X) = \prod_{i=1}^n f(X_i|\theta)$$

In practice, it's often easier to work with the **log-likelihood function**, which is the natural logarithm of the likelihood function. The logarithmic transformation turns products into sums, simplifying the maximization process:

$$\ell(\theta|X) = \log \mathcal{L}(\theta|X) = \sum_{i=1}^n \log f(X_i|\theta)$$

Maximizing the log-likelihood function is mathematically equivalent to maximizing the likelihood itself, but it is computationally simpler and more stable.

8.3.1 Likelihood Function

Given an assumed probability distribution (e.g., normal distribution, exponential distribution, etc.), the likelihood function represents the probability of observing the given data under that distribution. For instance, if X_1, X_2, \dots, X_n are observed data points from a distribution with parameter θ , the likelihood function $\mathcal{L}(\theta)$ is expressed as:

$$\mathcal{L}(\theta|X) = P(\{x_1, x_2, \dots, x_n\} | \theta)$$

Where P denotes the probability (or probability density) of observing the data given the parameters.

The Left-Hand Side: $\theta | X$: The notation $\theta | X$ on the left-hand side indicates the conditional dependence of the parameters on the observed data. We are interested in the distribution of the parameters θ given the data X , and the likelihood function $L(\theta | X)$ is designed to express this. It's clear that the likelihood is defined in terms of parameters, with the data acting as the conditioning factor, further emphasizing that the model is being adjusted to fit the data.

Likelihood Function Structure: The likelihood function $L(\theta | X)$ represents the probability of observing the data $X = (x_1, x_2, \dots, x_n)$, given a particular set of parameters θ . Notice the conditioning notation $P(x_i | \theta)$ on the right-hand side, where x_i represents the observed data and θ represents the unknown parameters that govern the model. The likelihood function $L(\theta | X)$ then depends on the parameters θ that we wish to estimate.

The Concept of ‘Fitting the Model’: This equation reinforces the idea of “fitting models to data” because we are searching for the parameter values $\hat{\theta}$ that maximize the likelihood of the data under the model. We do not alter the observed data x_i ; rather, we adjust the model’s parameters θ so that the model explains the data as best as possible. Maximizing the likelihood function allows us to choose the parameters that make the data as probable as possible under the assumed model.

8.3.2 The Log-Likelihood

Log-Likelihood for Convenience: The log-likelihood function, which is the logarithm of the likelihood, turns the product of probabilities into a sum of log-probabilities. This transformation simplifies the maximization process and stabilizes numerical calculations, especially when dealing with small probabilities that could otherwise lead to computational underflow. The loglikelihood still retains the same relationship: the model's parameters θ are adjusted, not the data x_i .

The log-likelihood function is typically easier to work with. By taking the natural logarithm of the likelihood function, we convert the product of probabilities into a sum, which is computationally more stable:

$$\ell(\theta|X) = \log \left(\prod_{i=1}^n f(X_i|\theta) \right) = \sum_{i=1}^n \log f(X_i|\theta)$$

This transformation simplifies differentiation and allows us to apply optimization techniques more efficiently.

Useful identities

We remind the reader of some useful log identities that are useful in the manipulation of these sorts of expressions:

$$\begin{aligned}\log(a \times b) &= \log(a) + \log(b) \\ \log(p^a) &= a \log(p) \\ \log(p^{a+b}) &= a \log(p) + b \log(p) \\ \log\left(\prod_{i=1}^n p^{x_i}\right) &= \sum_{i=1}^n x_i \log(p) \\ \prod_{i=1}^n p^{x_i} &= p^{\sum x_i}\end{aligned}$$

8.3.3 Maximum Likelihood Estimation

The maximum likelihood estimator $\hat{\theta}$ is the value of θ that maximizes the log-likelihood function. Mathematically, this is equivalent to solving the following optimization problem:

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta|X)$$

In cases where the likelihood function is differentiable, we can use calculus to find the maximum. Specifically, we take the derivative of the log-likelihood with respect to θ and set

it equal to zero:

$$\frac{\partial}{\partial \theta} \ell(\theta | X) = 0$$

This yields a set of equations that can be solved to find the optimal parameter estimates.

8.4 Connection to Other Estimation Methods

Maximum Likelihood Estimation is closely related to other methods of statistical estimation, including **least squares estimation**. In fact, for many common models, MLE is equivalent to the least squares method. For example, in the case of linear regression where the errors are assumed to follow a normal distribution with constant variance, the MLE for the regression parameters turns out to be the same as the least squares solution.

From a **Bayesian** perspective, MLE can be seen as a special case of **maximum a posteriori (MAP) estimation** with a uniform prior. In frequentist inference, MLE is a specific instance of an **extremum estimator**, where the objective function is the likelihood function.

8.5 Theoretical Properties of MLE

One of the reasons MLE is so widely used is that it has desirable **asymptotic properties** under certain regularity conditions:

- **Consistency:** As the sample size grows, the MLE $\hat{\theta}$ converges to the true parameter value θ . In other words, MLE produces increasingly accurate estimates as more data becomes available.
- **Asymptotic Normality:** Under regular conditions, the distribution of the MLE becomes approximately normal as the sample size increases. This means that the estimator becomes approximately unbiased and its variance shrinks as the sample size grows.
- **Efficiency:** MLE is often asymptotically efficient, meaning it achieves the lowest possible variance among unbiased estimators. In fact, it attains the **Cramer-Rao lower bound**, which sets a lower bound on the variance of an unbiased estimator.

However, MLE is not without its challenges. In some cases, the likelihood function may not have a closed-form solution, requiring numerical optimization methods such as **gradient descent** or **expectation-maximization (EM)**. Additionally, MLE estimators can be **sensitive to model misspecification**, and in small sample sizes, they may exhibit **bias** or **overfitting**.

Maximum Likelihood Estimation is a powerful and widely used statistical method for parameter estimation. By maximizing the likelihood function, we can obtain parameter estimates that are consistent, efficient, and asymptotically normal. This method is applicable in a wide range of statistical models and provides a solid foundation for both frequentist and Bayesian approaches to statistical inference.

In the following sections, we will delve deeper into MLE, explore examples of its application, and discuss how to implement MLE in various contexts. We will also examine the limitations of MLE and compare it to alternative estimation methods.

8.6 The Role of the Fisher Information

The **Fisher Information** is a measure of the amount of information that an observable random variable carries about an unknown parameter. It plays a significant role in MLE, especially in assessing the variance of the estimator.

8.6.1 Fisher Information

$$I(\theta) = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \log \mathcal{L}(\theta) \right]$$

The Fisher Information is used to compute the **Cramer-Rao lower bound**, which sets a lower limit on the variance of an unbiased estimator.

8.7 The Role of the Cramer-Rao Lower Bound

The **Cramer-Rao lower bound (CRLB)** plays an important role in statistical estimation theory, particularly in determining the efficiency of an estimator. It provides a theoretical limit on the variance of any unbiased estimator of a parameter, thereby offering a benchmark against which the performance of different estimators can be evaluated.

Definition of the Cramer-Rao Lower Bound

The **Cramer-Rao lower bound** establishes that, for any unbiased estimator $\hat{\theta}$ of a parameter θ , the variance of $\hat{\theta}$ cannot be smaller than the inverse of the **Fisher Information** $I(\theta)$:

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)}$$

Where: - $I(\theta)$ is the **Fisher Information** of the parameter θ , which quantifies the amount of information the observed data provides about the parameter.

8.8 2. Intuition Behind the Cramer-Rao Lower Bound

The CRLB gives an indication of how well a parameter can be estimated based on the available data:

- The **Fisher Information** $I(\theta)$ essentially measures the sensitivity of the likelihood function to changes in the parameter θ . A high value of $I(\theta)$ implies that the likelihood function changes steeply with respect to θ , indicating that the parameter is easier to estimate.
- Conversely, a lower value of $I(\theta)$ means the likelihood is relatively flat, suggesting that the parameter is harder to estimate accurately.

The **CRLB** provides the lowest possible variance for an unbiased estimator, which means that no other unbiased estimator can achieve a variance lower than this bound. Thus, it serves as a benchmark for evaluating the efficiency of an estimator.

8.9 3. Efficiency of Estimators

An estimator that **attains the CRLB** is called an **efficient estimator**. Such an estimator is said to be optimal because it achieves the minimum possible variance for unbiased estimators, making it the best possible estimator in terms of accuracy among the class of unbiased estimators.

For example, in the case of estimating the mean of a normal distribution with known variance, the sample mean is an efficient estimator because its variance matches the CRLB.

8.10 4. Role in Statistical Inference

- **Theoretical Benchmark:** The CRLB is often used as a **theoretical benchmark** to determine how well an estimator is performing. If an estimator has a variance close to the CRLB, it is considered efficient and optimal. If the variance is far from the CRLB, there is potential room for improvement.
- **Unbiased Estimators:** The CRLB is only applicable to **unbiased estimators**. An estimator is unbiased if its expected value equals the true parameter value, meaning that it correctly estimates the parameter “on average”. The CRLB gives us the lower variance bound for such unbiased estimators, but it does not provide information about biased estimators.
- **Asymptotic Efficiency:** Many maximum likelihood estimators (MLEs) are known to be **asymptotically efficient**, meaning that as the sample size grows to infinity, the variance of the MLE approaches the CRLB. This property is one of the reasons why MLE is such a popular and powerful method for parameter estimation.

8.11 5. Practical Implications

- **Model Evaluation:** In practical applications, the CRLB helps in evaluating the effectiveness of different models and estimation techniques. By comparing the variance of an estimator to the CRLB, statisticians can determine whether their chosen estimator is nearly optimal or whether alternative methods might yield more accurate estimates.
- **Trade-off Between Bias and Variance:** Although the CRLB is useful for unbiased

estimators, many real-world problems involve estimators that are biased but have a lower overall **mean squared error (MSE)**. The **bias-variance trade-off** often leads to the choice of biased estimators that offer lower MSE compared to strictly unbiased ones.

8.12 6. Examples

- **Normal Distribution Example:** For a normal distribution with mean μ and variance σ^2 , if σ^2 is known, the CRLB for estimating μ is:

$$\text{Var}(\hat{\mu}) \geq \frac{\sigma^2}{n}$$

where n is the sample size. In this case, the sample mean $\hat{\mu}$ is an efficient estimator, as it achieves this bound.

- **Poisson Distribution Example:** For a Poisson-distributed random variable with parameter λ , the CRLB for estimating λ is:

$$\text{Var}(\hat{\lambda}) \geq \frac{\lambda}{n}$$

where n is the number of observations. Here, the sample mean of the Poisson-distributed data is an efficient estimator for λ .

The **Cramer-Rao lower bound** plays a critical role in statistical estimation by providing a **lower limit on the variance** of any unbiased estimator of a parameter. It serves as a benchmark for determining the **efficiency** of an estimator, with efficient estimators achieving the CRLB and thus having the minimum possible variance among unbiased estimators. The CRLB also highlights the relationship between Fisher Information and the precision of parameter estimates, with larger Fisher Information indicating a tighter bound and more accurate estimation.

Understanding the CRLB helps in assessing estimator quality, choosing appropriate estimation techniques, and determining whether an estimator is optimal in the sense of having the smallest possible variance. It is a fundamental concept in statistical theory that aids in ensuring the precision and reliability of parameter estimates in practice.

8.13 Estimating Parameters with MLE

In this section, we will demonstrate how to apply MLE to estimate parameters for a simple model. We'll use both a theoretical example and a data-driven case study.

8.13.1 Example: Estimating the Mean of a Normal Distribution

Given a dataset of independent and identically distributed (i.i.d.) observations X_1, X_2, \dots, X_n from a normal distribution with unknown mean μ and variance σ^2 , the MLE for μ is the sample mean.

8.13.1.1 MLE for the Mean of a Normal Distribution

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$$

This estimator is obtained by maximizing the log-likelihood for the normal distribution.

8.14 Properties of MLE

MLE estimators have desirable properties under certain regularity conditions:

- **Consistency:** As the sample size increases, the MLE converges to the true parameter value.
- **Asymptotic Normality:** The distribution of the MLE approaches a normal distribution as the sample size grows large.
- **Efficiency:** MLE is often asymptotically efficient, meaning it achieves the lowest possible variance among unbiased estimators.

8.15 Challenges and Limitations of MLE

While MLE is a powerful method, it comes with its own set of challenges:

- **Complexity of Maximization:** In some cases, the likelihood function may not have a closed-form solution, requiring numerical optimization.
- **Non-Identifiability:** The parameters may not be identifiable if the likelihood function does not contain enough information to estimate the parameters uniquely.
- **Overfitting:** MLE can overfit the data if the model is too complex relative to the available data.

8.16 GPT Dump Conclusion

In this chapter, we introduced the principles of Maximum Likelihood Estimation, a widely used method for parameter estimation in statistical models. We discussed the likelihood function, the MLE process, its properties, and common challenges associated with MLE. In the next chapter, we will explore Bayesian Estimation and compare it with MLE.

8.17 Claude Dump

8.17.1 Maximum Likelihood Estimation

The Maximum Likelihood Estimation (MLE) stands as one of the most fundamental and powerful methods in statistical inference and machine learning. This chapter explores how MLE provides a systematic approach to parameter estimation by finding values that maximize the likelihood of observing the given data.

8.18 Historical Context and Development

The concept of maximum likelihood was first developed by R.A. Fisher in the 1920s, though similar ideas can be traced back to Carl Friedrich Gauss and Pierre-Simon Laplace. Fisher's revolutionary paper "On the Mathematical Foundations of Theoretical Statistics" (1922) formally introduced MLE, establishing it as a cornerstone of modern statistical theory.

8.19 Theoretical Foundation

8.19.1 The Likelihood Function

The likelihood function represents the probability of observing the data given a set of parameters. For a dataset $X = \{x_1, \dots, x_n\}$ and parameters θ , the likelihood function is:

$$L(\theta|X) = P(X|\theta) = \prod_{i=1}^n P(x_i|\theta)$$

8.19.2 Log-Likelihood

We often work with the log-likelihood for computational convenience:

$$\ell(\theta|X) = \log L(\theta|X) = \sum_{i=1}^n \log P(x_i|\theta)$$

8.20 The Maximum Likelihood Principle

Key ML Equation 9: Maximum Likelihood Estimator

The Maximum Likelihood Estimator $\hat{\theta}_{MLE}$ is defined as:

$$\arg \max_{\theta} \prod_{i=1}^n P(x_i|\theta)$$

$$\arg \max_{\theta} L(\theta | X) = \arg \max_{\theta} \sum_{i=1}^n \log P(x_i | \theta)$$

8.21 Properties of MLE

8.21.1 Consistency

Under regular conditions, MLE estimates converge to the true parameter values as sample size increases.

8.21.2 Asymptotic Normality

For large samples, the MLE is approximately normally distributed around the true parameter value.

8.21.3 Efficiency

MLE achieves the Cramér-Rao lower bound asymptotically, making it an efficient estimator.

8.22 Applications in Machine Learning

8.22.1 Gaussian Distribution Example

Example 8.1: MLE for Gaussian Distribution

For a normal distribution $N(\mu, \sigma^2)$, the MLEs are:

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

8.22.2 Linear Regression Connection

Lemma 8.1: MLE in Linear Regression

The MLE for linear regression parameters under Gaussian noise assumptions coincides with the least squares solution.

When analytical solutions aren't available, numerical methods like gradient descent or Newton-Raphson are used:

```
def mle_optimization(data, initial_params):
    def negative_log_likelihood(params):
        return -np.sum(log_probability(data, params))

    result = optimize.minimize(negative_log_likelihood,
                               initial_params,
                               method='BFGS')
    return result.x
```

8.23 MLE Section

8.24 MLE Section

8.25 MLE Section

8.26 MLE Section

8.27 MLE Section

8.28 Determining MLE: A Generalized Approach

Maximum Likelihood Estimation (MLE) is one of the most widely used methods in statistical estimation, offering a flexible and efficient way to estimate parameters of a statistical model based on observed data. At its core, MLE aims to find the set of parameters that maximize the likelihood of the observed data under a given model. In the context of this section, we will explore the general methodology for computing the MLE, without assuming any specific form for the underlying probability density function (PDF). Instead, this flowchart serves as a broad, multidimensional illustration of the steps required to derive the MLE, providing a systematic approach to parameter estimation applicable to a wide range of statistical models.

8.28.1 General Overview of the MLE Process

The flowchart presented here details the sequential steps involved in the computation of the MLE for a given dataset. It is a generalized workflow that does not assume a particular distribution or parameterization for the data. Instead, it provides a robust framework for identifying the parameters of an unknown model, ensuring that the resulting MLE is well-defined and mathematically sound. By focusing on the methodology of parameter estimation, the flowchart offers a clear visualization of the steps involved in determining the maximum likelihood estimates (MLE) that best describe the observed data.

1. Initialization and Likelihood Calculation: The first step in this process is the specification of the likelihood function $L(\Theta)$, where $\Theta = [\theta_1, \theta_2, \dots, \theta_m]$ is the vector of parameters to be estimated, and $X = [x_1, x_2, \dots, x_n]$ represents the observed data. For a given model, the likelihood function $L(\Theta)$ is defined as the probability of observing the data given a particular set of parameter values. In this generalized framework, the likelihood function is assumed to take the form $L(\Theta) = \prod_{i=1}^n P(x_i | \Theta)$, where $P(x_i | \Theta)$ represents the conditional probability of the data point x_i given the parameter vector Θ .
2. Log-Likelihood and Gradient Calculation: Since the likelihood function is often complicated and involves large products of probabilities, it is common practice to work with the loglikelihood function, $\ell(\Theta) = \ln L(\Theta)$. Taking the logarithm of the likelihood simplifies the expression and transforms products into sums, making it more tractable for optimization. The log-likelihood is given by $\ell(\Theta) = \sum_{i=1}^n \ln P(x_i | \Theta)$. The next step involves computing the gradient of the log-likelihood with respect to the parameters, denoted as $\nabla_{\Theta} \ell(\Theta)$. This gradient represents the first-order partial derivatives of the log-likelihood function with respect to each parameter θ_k , and it provides the direction in which the log-likelihood increases most rapidly.
3. Solving for the Local Maxima: The goal of MLE is to maximize the log-likelihood function. To find the parameter values that maximize $\ell(\Theta)$, we set the gradient equal to zero, i.e., $\nabla_{\Theta} \ell(\Theta) = 0$. This condition ensures that the log-likelihood function has a

local maximum at the estimated parameter values. Solving this system of equations typically requires numerical methods, especially in high-dimensional settings, where an analytical solution may not be feasible.

4. Hessian Matrix and Curvature Evaluation: After finding the critical points where the gradient is zero, we next examine the nature of these critical points. This is accomplished by evaluating the second-order partial derivatives of the log-likelihood function, which form the Hessian matrix $H(\Theta)$. The Hessian matrix is defined as $H(\Theta) = \frac{\partial^2 \ell(\Theta)}{\partial \Theta^2}$, and it encapsulates information about the curvature of the log-likelihood surface. A negative-definite Hessian matrix indicates that the critical point corresponds to a local maximum. The negative definiteness of the Hessian can be verified by examining its eigenvalues. If all eigenvalues of the Hessian are negative, the critical point is indeed a maximum, and the corresponding parameter values are the MLE.
5. Eigenvalue Check and Final Validation: The final step in the process is to verify whether the Hessian matrix is negative-definite, which confirms that the solution corresponds to a local maximum. This can be done by checking if all the eigenvalues of the Hessian matrix are negative. In the flowchart, this is represented as checking the condition $\mathbf{v}^\top H(\Theta)\mathbf{v} < 0$ for all non-zero vectors \mathbf{v} , or equivalently by checking whether the eigenvalues of H are negative (i.e., $H\mathbf{v} = \lambda\mathbf{v}$ with $\lambda_i < 0$ for all i). If the Hessian is negative-definite, the local maximum is confirmed, and the MLE for the parameters is validated.
6. Iteration and Convergence: In some cases, especially in high-dimensional or complex problems, the optimization may require iterative steps to converge to the global maximum. The algorithm may re-seed the initial parameter estimates Θ and attempt to find new local maxima until convergence is achieved. The algorithm halts once the MLE is determined, and the estimated parameters Θ are printed as the final result.

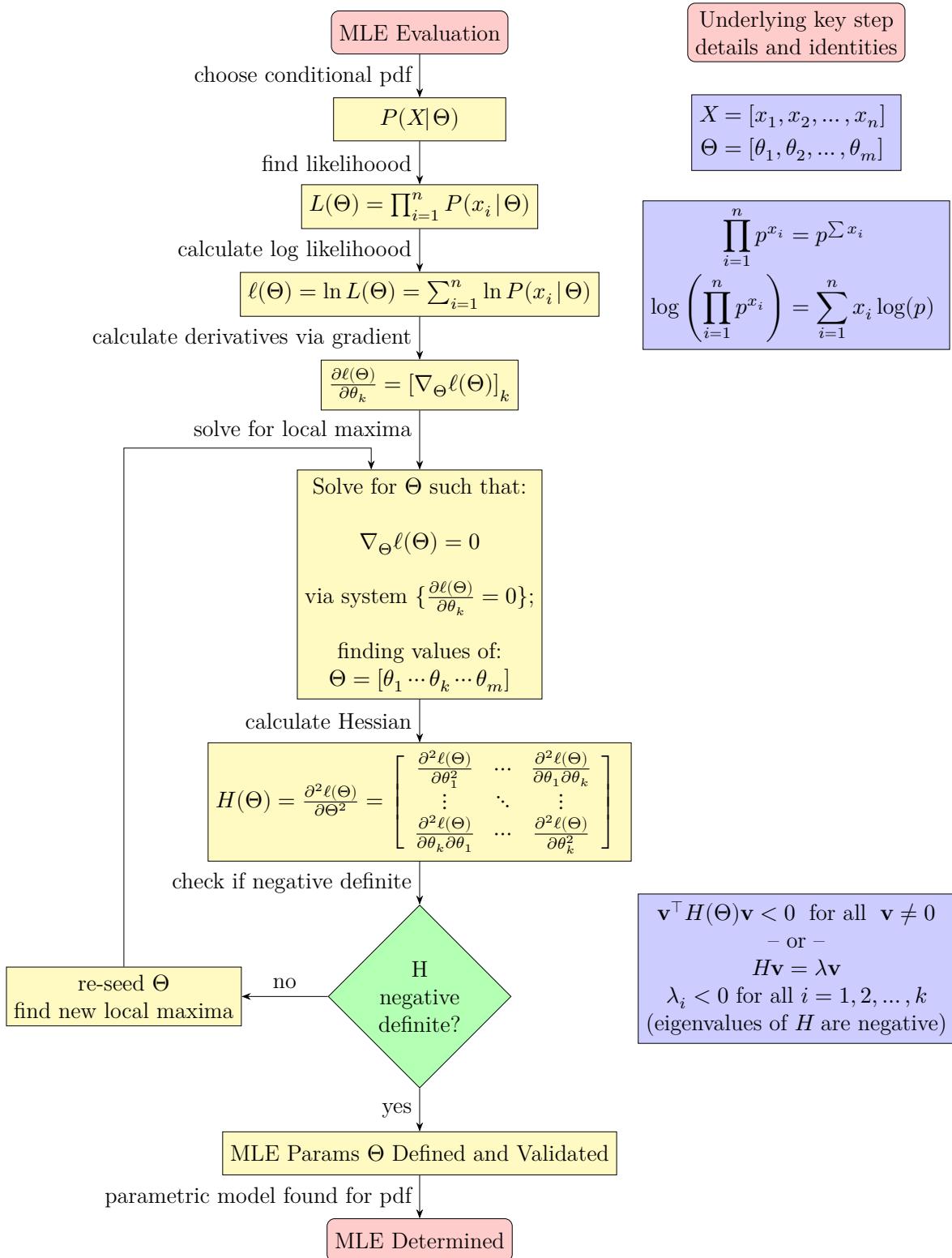


Figure 8.1: MLE Flow Chart

This generalized flowchart for determining the MLE highlights the flexibility and power of the maximum likelihood method. By not assuming a specific form for the underlying PDF, it provides a versatile framework that can be applied to a wide variety of statistical models. Whether the data follows a normal distribution, a Poisson distribution, or any other form, the steps outlined in this process can be adapted to estimate the parameters of the model. The primary goal of this workflow is to demonstrate the general steps involved in determining the MLE, from likelihood calculation and gradient computation to the evaluation of the Hessian matrix and verification of the local maximum.

Through these steps, one can confidently arrive at the parameter estimates that best describe the data according to the principle of maximum likelihood. The beauty of this approach lies in its generality, making it an invaluable tool for statisticians, data scientists, and anyone interested in understanding the underlying relationships in complex datasets. By following the methodical procedure outlined here, practitioners can ensure that the MLE provides a robust and accurate representation of the data, laying the foundation for further analysis and modeling.

8.29 STEP By STEP Description of MLE (worth it ?)

Equation 9

Linear Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + \epsilon$$



Line Fitting Straight Assumption

Key ML Equation 10: Linear Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + \epsilon \quad (9.0.1)$$

y The **dependent variable**, representing the outcome or the target we aim to predict.

x_i The **independent variable(s)**, representing the features or predictors for the i -th observation.

β_0 The **intercept**, representing the expected value of y when all predictors are zero.

β_i The **coefficient** for predictor x_i , representing the change in the expected value of y for a one-unit change in x_i , while holding all other variables constant.

ϵ The **error term**, accounting for the variability in y that cannot be explained by the linear combination of predictors.

Introduction: Linear regression is a simple and widely used method for predictive analysis.

Description: It models the relationship between a dependent variable and one or more independent variables by fitting a linear equation to observed data.

Importance in ML: Linear regression is fundamental for understanding relationships between variables and is used in predictive modeling to estimate outcomes.

9.1 Introduction

Linear regression is a fundamental statistical approach used to model the relationship between a dependent variable and one or more independent variables. The objective is to establish a linear relationship between inputs (predictors) and an output variable that allows for interpretation, prediction, and inference.

This chapter discusses Equation 9, which formulates the core concept of linear regression. The method is commonly applied across various fields such as finance, economics, medicine, and more.

9.2 Mathematical Formulation

The linear regression equation is represented by Equation 9. Here, the model is built by finding the coefficients $\beta_0, \beta_1, \dots, \beta_n$ that minimize the error term ϵ .

The general approach for estimation in linear regression is based on minimizing the residual sum of squares (RSS):

$$RSS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_n x_{in}))^2$$

The goal is to determine values for the coefficients β_i that minimize the RSS , resulting in the best fit line through the observed data.

9.3 Understanding the Components

- **Dependent Variable (y):** The outcome that we are trying to predict.
- **Independent Variables (x_i):** Predictors or features used in the model.
- **Coefficients (β_i):** Parameters that describe the influence of the predictors on the outcome.
- **Error Term (ϵ):** Represents the portion of y that cannot be explained by the linear model.

9.4 Estimation using Ordinary Least Squares (OLS)

Ordinary Least Squares (OLS) is the most common method used to estimate the parameters in linear regression. The idea is to minimize the difference between the observed values and the values predicted by the model. This difference is captured by the residual sum of squares, which is minimized to obtain the optimal values of the coefficients.

9.5 Implementation in R

Below, we demonstrate how to fit a linear regression model in R using a simple mock dataset. The dataset consists of a response variable y and two predictors x_1 and x_2 .

```
# Generating Mock Data
set.seed(123)
x1 <- rnorm(100, mean = 5, sd = 2)
x2 <- rnorm(100, mean = 10, sd = 3)
epsilon <- rnorm(100, mean = 0, sd = 1)
y <- 3 + 2*x1 - 0.5*x2 + epsilon

# Creating a data frame
data <- data.frame(y, x1, x2)

# Linear Regression Model
model <- lm(y ~ x1 + x2, data = data)

# Summary of the Model
summary(model)

##
## Call:
## lm(formula = y ~ x1 + x2, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -1.8730 -0.6607 -0.1245  0.6214  2.0798 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 3.38862   0.43979   7.705 1.13e-11 ***
## x1          1.93341   0.05243  36.873  < 2e-16 ***
## x2         -0.49206   0.03300 -14.912  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9513 on 97 degrees of freedom
## Multiple R-squared:  0.9442, Adjusted R-squared:  0.943 
## F-statistic: 820.2 on 2 and 97 DF,  p-value: < 2.2e-16
```

9.6 Limitations or Pitfalls of Linear Regression

While linear regression is a simple and often effective approach, it comes with certain limitations and potential pitfalls. Understanding these drawbacks is crucial when determining whether linear regression is the best choice for a given dataset. Here are some of the key limitations:

1. **Assumption of Linearity:** Linear regression assumes that there is a linear relationship between the independent and dependent variables. If the true relationship is non-linear, the linear model will not adequately capture the complexity of the data, leading to poor predictive performance.
2. **Outliers Can Have a Large Influence:** Linear regression is sensitive to outliers, which can disproportionately affect the fit of the model. A single extreme value in the data can significantly alter the coefficient estimates, leading to a biased or misleading interpretation of the model (more on this below).
3. **Multicollinearity:** When independent variables are highly correlated, the estimates of the regression coefficients can become unstable. This phenomenon, known as multicollinearity, makes it difficult to determine the independent effect of each predictor on the dependent variable.
4. **Non-Constant Variance (Homoscedasticity):** Linear regression assumes that the variance of the errors is constant across all levels of the independent variables. If this assumption is violated (heteroscedasticity), it can lead to inefficient estimates and affect hypothesis testing.
5. **Assumption of Normality in Errors:** Linear regression assumes that the error terms are normally distributed. If the errors deviate significantly from normality, this can impact confidence intervals and the reliability of hypothesis tests.
6. **Overfitting in High Dimensions:** When there are many predictors relative to the number of observations, linear regression can easily overfit the data, capturing noise instead of the underlying trend. Regularization methods such as Lasso or Ridge regression may be more appropriate in such cases.

9.7 Impact of Outliers on Goodness of Fit

Outliers are extreme values that do not follow the general trend of the rest of the data. In linear regression, outliers can have a major influence on the estimated model parameters, ultimately affecting the goodness of fit. Let's discuss how outliers impact the model:

- **Effect on Coefficients:** Since linear regression aims to minimize the residual sum of squares (RSS), any point that lies far away from the regression line will have a disproportionately large impact on the resulting line. The presence of outliers can pull the regression line toward themselves, distorting the relationship between the dependent and independent variables.
- **Leverage and Influence:** Outliers can be categorized into two types: those that have high leverage and those that have high influence. Leverage points are those that lie far from the mean of the predictors, and they can significantly affect the position of the regression line even if their response values align with the trend. Influential points are outliers that change the overall fit of the model, and removing these points can substantially alter the estimated coefficients.
- **Example of Outlier Influence:** The effect of an outlier can be illustrated using a simple plot:

```
# Generating Mock Data with an Outlier
set.seed(456)
x1 <- rnorm(100, mean = 5, sd = 2)
epsilon <- rnorm(100, mean = 0, sd = 1)
y <- 3 + 2*x1 + epsilon

# Adding an Outlier
x1 <- c(x1, 20)
y <- c(y, 100)

# Fitting a Linear Model
data_with_outlier <- data.frame(y, x1)
model_with_outlier <- lm(y ~ x1, data = data_with_outlier)

# Plotting
plot(x1, y, main = "Effect of Outlier on Linear Regression",
      xlab = "x1", ylab = "y", pch = 19)
abline(model_with_outlier, col = "red")
```

Effect of Outlier on Linear Regression

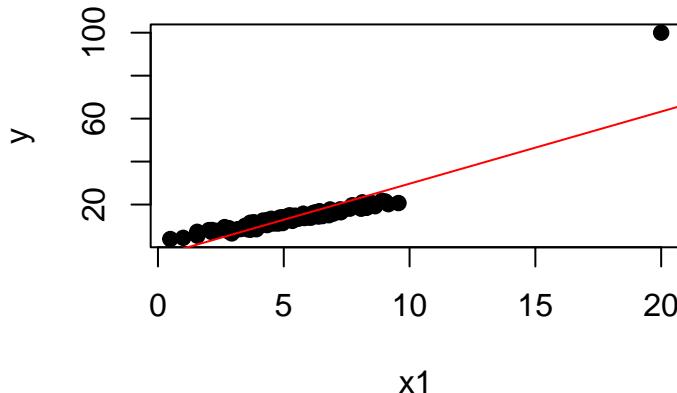


Figure 9.1: Default caption for all figures

In the plot above, the red regression line is influenced by the single outlier, deviating significantly from what it would be without that extreme point. In practice, techniques such as robust regression can be employed to minimize the impact of such outliers.

9.7.1 Residual Sum of Squares (RSS) and L2 Norm

9.7.1.1 Residual Sum of Squares (RSS)

The Residual Sum of Squares (RSS) is a key concept in linear regression, representing the total deviation of the observed values from the values predicted by the model. The objective of linear regression is to minimize the RSS, which is mathematically defined as:

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

9.7.1.2 The L2 Norm

The RSS minimization problem is related to the concept of the L2 norm. In the context of linear regression, the L2 norm represents the Euclidean distance between the actual observed values and the predicted values of the model. The L2 norm is defined as:

$$\|\mathbf{r}\|_2 = \left(\sum_{i=1}^n r_i^2 \right)^{1/2}$$

where $r_i = y_i - \hat{y}_i$ are the residuals. The RSS is simply the squared L 2 norm of the residual vector, meaning that the linear regression objective can be understood as minimizing the L2 distance between the predicted and actual values.

This connection to the L2 norm is one of the reasons why linear regression is so sensitive to

outliers. Because each residual is squared, larger residuals have a disproportionate impact, which leads to a model that might be unduly influenced by a few data points.

9.7.1.3 Comparison with L1 Norm

An alternative to the L2 norm is the L1 norm, which forms the basis for Lasso Regression. The L1 norm is defined as:

$$\|\mathbf{r}\|_1 = \sum_{i=1}^n |r_i|$$

Minimizing the L1 norm instead of the L2 norm results in a model that is less sensitive to outliers. This is because the L1 norm does not square the residuals, giving less weight to extreme errors. In fact, L1-based models often produce sparse solutions, making them particularly useful in high-dimensional settings where feature selection is important.

9.7.1.4 Regularization and the L2 Penalty

In practice, Ridge Regression (a form of regularized linear regression) adds an L2 penalty term to the loss function to prevent overfitting:

$$L(\beta) = RSS + \lambda \sum_{j=1}^n \beta_j^2$$

The λ term controls the amount of regularization, shrinking the coefficients toward zero and reducing model complexity. This can be particularly useful when dealing with multicollinearity or high-dimensional datasets.

Practical Implementation with Regularization

Below is an example in R where we use the `glmnet` package to fit a Ridge regression model:

```
# Load the glmnet package
library(glmnet)

# Generating Mock Data
set.seed(789)
x1 <- rnorm(100, mean = 5, sd = 2)
x2 <- rnorm(100, mean = 10, sd = 3)
epsilon <- rnorm(100, mean = 0, sd = 1)
y <- 3 + 2*x1 - 0.5*x2 + epsilon

# Creating a matrix for predictors (required by glmnet)
X <- as.matrix(data.frame(x1, x2))

# Fitting a Ridge Regression Model
ridge_model <- glmnet(X, y, alpha = 0) # alpha = 0 for Ridge
```

```
# Plotting the coefficient paths
plot(ridge_model, xvar = "lambda", label = TRUE)
title(main = "Coefficient Paths for Ridge Regression")
```

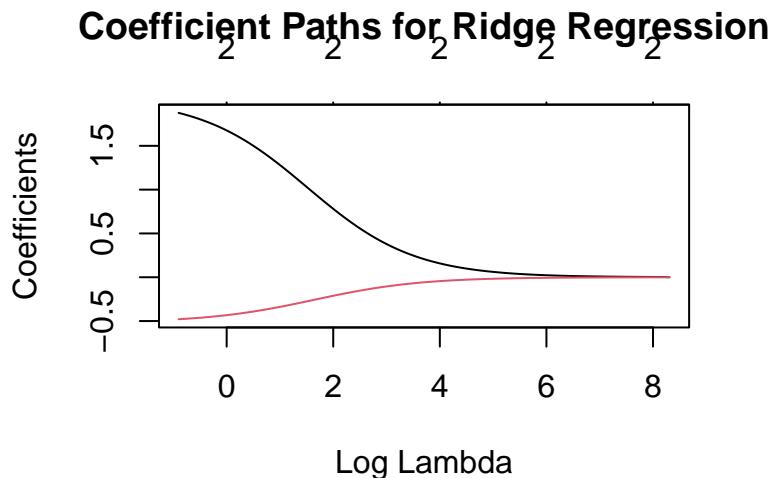


Figure 9.2: Default caption for all figures

The plot above shows how the coefficients shrink as the regularization parameter λ increases, helping to control overfitting while still capturing the main patterns in the data.

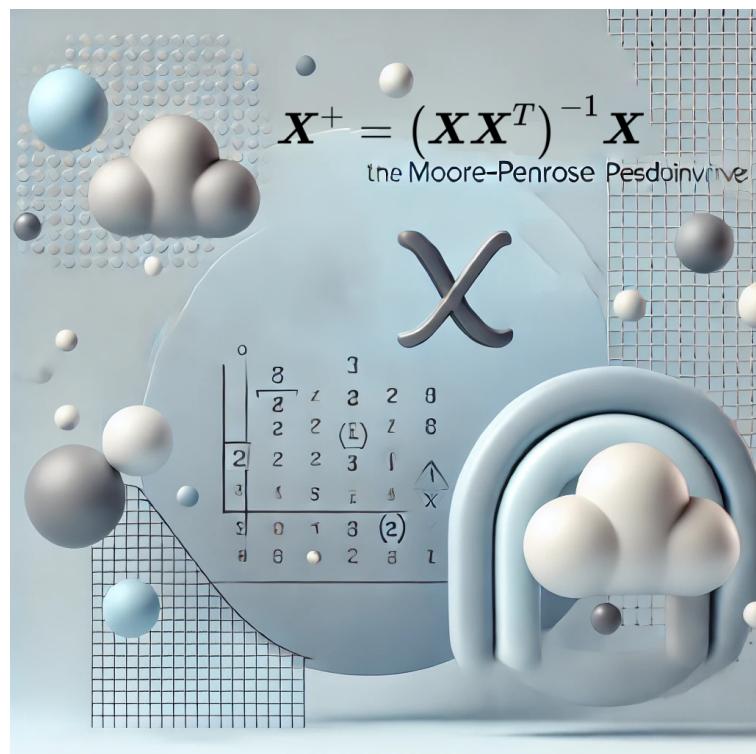
9.8 Equation Summary

In this extended discussion of linear regression, we highlighted some of the key limitations of the method, discussed the significant impact that outliers can have on the model, and introduced the concept of RSS and its relation to the L2 norm. We also explored how these concepts are extended in regularized linear regression, such as Ridge regression, to address issues like overfitting and multicollinearity. Linear regression remains a powerful and widely used tool, but it's important to be mindful of the underlying assumptions and potential pitfalls. By understanding the limitations and exploring regularization techniques, we can build more robust and effective predictive models.

Equation 10

Ordinary Least Squares (OLS)

$$\hat{\beta} = (X^T X)^{-1} X^T y$$



Multivariate Linear Model Fitting

Key ML Equation 11: Maximum Likelihood Estimation

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (10.0.1)$$

 $\hat{\beta}$

The **estimated parameter vector**, representing the values that maximize the likelihood function in the context of linear regression.

 X

The **design matrix**, a matrix of independent variables (features) where each row corresponds to an observation and each column to a feature in the dataset.

 y

The **vector of observed values**, representing the dependent variable or the target values corresponding to the observations in the dataset.

 X^T

The **transpose of the design matrix**, which switches the rows and columns of the matrix X .

 $(X^T X)^{-1}$

The **inverse of the matrix product** of X^T and X , a necessary step to solve for $\hat{\beta}$ using the normal equation.

Introduction: OLS is a method for estimating the unknown parameters in a linear regression model.

Description: It minimizes the sum of squared residuals between the observed and predicted values.

Importance in ML: OLS is a fundamental technique for regression analysis. It helps determine the best-fit line for the data, making it useful for predictive modeling.

10.1 Linear Model Intro

We have linear model - and assume the error, in this setting ϵ , is normally distributed. This is a very important consideration since it frames how the solution for the linear regression problem is obtained!

This is a VECTOR equation of the overall data set:

$$Y = A + BX + [\varepsilon]$$

Here we denote $[\varepsilon]$ as a VECTOR of normal error for the outcome vector Y, where

$$\varepsilon \sim N(0, \sigma^2) \text{ is i.i.d, Markov}$$

This is the linear model estimator for a sample over i scalar equation (notice the hat in \hat{y}_i):

$$\hat{y}_i = \alpha + \beta x_i$$

$$\varepsilon_i = \hat{y}_i - y_i$$

We would like to determine α and β . So, we are trying to minimize this equation, the residuals of error:

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

Now we consider this from a MLE perspective, but first consider the expected value of y_i :

What is dist of y_i ?, look at the estimated value of y_i :

$$\begin{aligned} E(y_i) &= E(\alpha + \beta x_i + \varepsilon_i) \\ &= \alpha + \beta x_i + E(\varepsilon_i) \\ &= \alpha + \beta x_i \end{aligned}$$

Now also consider the variance of y_i :

$$\begin{aligned} \text{Var}(y_i) &= \text{Var}(\alpha + \beta x_i + \varepsilon_i) \\ &= \text{Var}(\varepsilon_i) \\ &= \sigma^2 \end{aligned}$$

Note the pdf for y_i is all about the expected value and the variance of the data. So we can say this in a specific way ... take a look at this to better understand the definitions here.

$$y_i \sim N(E(y_i), \text{Var}(y_i)) \sim N(\alpha + \beta x_i, \sigma).$$

Here is the “burn your noodle idea here” ... These y_i are independent, but NOT identically distributed, since the ‘mean’ or expected value depends on x_i , as you should inspect $E(y_i)$ carefully and see $y_i \sim N(\alpha + \beta x_i, \sigma)$... A little tricky but important to realize.

Ok so lets write down the pdf for for y_i :

$$\Rightarrow f(y_i | \alpha, \beta, \sigma^2) = \frac{e^{-(y_i - \alpha - \beta x_i)^2 / 2\sigma^2}}{\sqrt{2\pi\sigma^2}}$$

This is a pdf for one single y_i , recalling from above that this has it’s own pdf, and each y_i is not identically distributed. So, we have to consider the likelihood.

The likelihood ...

$$\begin{aligned} \Rightarrow L(\alpha, \beta, \sigma^2 | y_1 \dots y_n) &= \prod_{i=1}^n f(y_i | \alpha, \beta, \sigma^2) = \prod_{i=1}^n \frac{e^{-(y_i - \alpha - \beta x_i)^2 / 2\sigma^2}}{\sqrt{2\pi\sigma^2}} \\ L(\alpha, \beta, \sigma^2 | y_1 \dots y_n) &= \frac{\prod_{i=1}^n \exp(-(y_i - \alpha - \beta x_i)^2 / 2\sigma^2)}{(2\pi\sigma^2)^{n/2}} \\ &= \frac{\exp(\sum_{i=1}^n \frac{-(y_i - \alpha - \beta x_i)^2}{2\sigma^2})}{(2\pi\sigma^2)^{n/2}} \end{aligned}$$

We designate the log L with a small l , and you get the log form of the likelihood (please recall the tricks for manipulating products and sums of logs):

$$l(\alpha, \beta, \sigma^2 | Y_1, \dots, Y_n) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 / 2\sigma^2 - \frac{n}{2} \log(2\pi\sigma^2)$$

So we are back to the MLE ideas: With the log-likelihood, we can maximized with the first derivative. “Nothing new to see here ... move along folks’’ lol. As such, we will determine the model parameters α and β using the MLE approach from before on our linear model log-likelihood expression.

So lets talk about what this looks like for β

$$y_i = \alpha + \beta x_i + \varepsilon_i \quad \text{and} \quad \varepsilon_i \sim N(0, \sigma^2) \quad \text{i.i.d.}$$

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) \quad \text{note} \quad \bar{x} = \mu_x$$

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$E(\hat{\beta}) = E\left[\frac{\sum_{i=1}^n Y_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}\right]$$

$$= \frac{\sum_{i=1}^n E[y_i] (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \frac{\sum_{i=1}^n (\alpha + \beta x_i) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{let } \gamma = \alpha - \beta \bar{x}$$

$$= \frac{\sum_{i=1}^n (\gamma + \beta (x_i - \bar{x})) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

collecting terms and cancelling we arrive at:

$$E(\hat{\beta}) = \frac{\gamma \sum_{i=1}^n (x_i - \bar{x}) + \sum_{i=1}^n \beta (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\gamma \sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} + \beta$$

focus on the definition of the sums of x

$$\begin{aligned} & \sum_{i=1}^n (x_i - \bar{x}) \\ &= [\sum x_i] - n\bar{x} = 0 \\ & \text{since ... } n\bar{x} = \sum x_i \end{aligned}$$

\

So we get a collapse of the expected value (kind of a let down):

$$E(\hat{\beta}) = \beta$$

We can extend this to multiple dimensions in X

$$P(Y|X) = N(f(X), \sigma^2)$$

Assume $f(X)$ is linear in X

$$f(X) = \sum_{i=1}^N W_i X_i + W_0 = W^T X$$

What is this? .. this is $y = mx + b$ in multiple dimensions.

Where we denote by X an extended vector with an additional entry

$$X = [1, X_1, X_2, \dots, X_N], W = [W_0, W_1, W_2, \dots, W_N]$$

Note in this case N is not the sample size. It is the number of independent variables we are trying to model.

Maximizing the log likelihood on a training set $\{(X^l, Y^l) | l = 1, \dots, L\}$

$$W^* = \operatorname{argmax}_W \sum_{l=1}^L \log P(Y^l | X^l; W)$$

Where

$$P(Y | X) = N(f(X), \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(Y - f(X))^2}{2\sigma^2} \right\}$$

Note in this case L IS the sample size. It is the number of samples of data we are processing.

Maximizing the log likelihood on a training set $\{(X^l, Y^l) | l = 1, \dots, L\}$

$$W^* = \operatorname{argmax}_W \sum_{l=1}^L \log P(Y^l | X^l; W)$$

Where

$$P(Y | X) = N(f(X), \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(Y - f(X))^2}{2\sigma^2} \right\}$$

Note in this case L IS the sample size. It is the number of samples of data we are processing.

Perform the derivative calculation:

$$\begin{aligned} \frac{\partial \sum_l (y - f(x; W))^2}{\partial w_i} &= \sum_l 2(y - f(x; W)) \frac{\partial (y - f(x; W))}{\partial w_i} \\ &= \sum_l -2(y - f(x; W)) \frac{\partial f(x; W)}{\partial w_i} \end{aligned}$$

Set the derivative to zero,

$$\sum_{l=1}^L (Y^l - W^T X^l) X^l = 0$$

Can be written as an alternative matrix representation

$$\mathbf{XY} - \mathbf{XX}^T \mathbf{W} = \mathbf{0}$$

Where X is data matrix, with each column being a training example, Y is the vector of target variables.

$$\mathbf{X} = [X^1, X^2, \dots, X^L], \mathbf{Y} = [Y^1, Y^2, \dots, Y^L]^T$$

$$\text{Solving } \mathbf{XY} - \mathbf{XX}^T \mathbf{W} = \mathbf{0} \quad W = (XX^T)^{-1} XY$$

Where

$$X^+ = (XX^T)^{-1} X$$

is called the Pseudo Inverse.

This leads us to the following final relation for MLE on Linear Regression ... showing the linear regression using the MLE for the optimization to find the proper fit parameters is the following elegant equation.

$$\arg \max_{\beta} \prod_{i=1}^n \text{Normal}(y_i | \mu = (X\beta)_i, \sigma) = (X^T X)^{-1} X^T y$$

We will dig into this further and also learn about MAP - which is the Bayesian path for this math.

Random Finds - for exposition

$$\text{Error} = \sum_{i=1}^n (\text{real.output} - \text{prediction})^2$$

Suppose we have k predictor variables x_1, \dots, x_k and a dependent variable y . We consider the linear relation:

$$\xi_i(\beta_0, \beta_1, \dots, \beta_k) = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki}$$

Often we do not explicitly write ξ as a function of the parameters β_0, \dots, β_k , and write instead:

$$\xi_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki}$$

The parameters $\beta_0, \beta_1, \dots, \beta_k$ are called (regression) coefficients. In particular, β_0 is called the (regression) intercept and β_1, \dots, β_k are (regression) slope coefficients.

The term simple linear regression is often used to cover the special case of $k = 1$. If there is more than one predictor, i.e., $k > 1$, the term multiple linear regression is common.

Based on the predictions of a parameter vector $\langle \beta_0, \beta_1, \dots, \beta_k \rangle$, we consider the residual sum of squares as a measure of prediction error:

$$\text{RSS}_{\langle \beta_0, \beta_1, \dots, \beta_k \rangle} = \sum_{i=1}^k [y_i - \xi_i(\beta_0, \beta_1, \dots, \beta_k)]^2$$

We would like to find the best parameter values (denoted traditionally by a hat on the parameter's variable: $\hat{\beta}_i$) in the sense of minimizing the residual sum of squares:

$$\langle \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k \rangle = \arg \min_{\langle \beta_0, \beta_1, \dots, \beta_k \rangle} \text{RSS}_{\langle \beta_0, \beta_1, \dots, \beta_k \rangle}$$

The prediction corresponding to the best parameter values is denoted by $\hat{\xi} \in \mathbb{R}^n$ and called the best linear predictor:

$$\hat{\xi}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \dots + \hat{\beta}_k x_{ki}$$

It is also possible, and often convenient, to state the linear regression model in terms of matrix

operations. Traditionally, we consider a so-called predictor matrix X of size $n \times (k+1)$, where n is the number of observations in the data set and k is the number of predictor variables. The predictor matrix includes the values for all predictor variables and it also includes an “intercept column” $(X^T)_0$ for which $X_{i0} = 1$ for all $1 \leq i \leq n$ so that the intercept β_0 can be treated on a par with the other regression coefficients.

Using the predictor matrix X , the linear predictor vector ξ is:

$$\xi = X\beta$$

This is the formulation of the Matrices ...

$$X = \underbrace{\begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & x_{23} & \cdots & x_{2k} \\ \vdots & \ddots & & & & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \cdots & x_{nk} \end{bmatrix}}_{n \times k}$$

$$\hat{\beta} = \overbrace{\begin{bmatrix} \hat{\beta}_0 & \hat{\beta}_1 & \cdots & \hat{\beta}_k \end{bmatrix}}^{k \times 1}^T$$

$$\hat{y} = \overbrace{\begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \vdots \\ \hat{y}_k \end{bmatrix}}^{n \times 1}$$

So let's show how this works out from the algebraic form of the regression to the matrix form of the pseudo inverse ... a with a little bit more detail.

Let X be the $n \times (k+1)$ regression matrix for a linear regression model with k predictor variables for a data set y with n observations. The solution for OLS regression

$$\hat{\beta} = \langle \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k \rangle = \arg \min_{\beta} \sum_{i=1}^k (y_i - (X\beta)_i)^2$$

is given by:

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Proof. With n observations, the vector ξ of predicted values for given coefficient vector β is:

$$\xi = X\beta$$

Where (recall) we have the ‘b’ encoded in the X vector definition

$$\xi = X\beta = X_0\beta_0 + X_N\beta_N$$

and

More explicitly, this means that:

$$\begin{aligned}\xi_1 &= \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \dots + \beta_k X_{1k} \\ \xi_2 &= \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \dots + \beta_k X_{2k} \\ &\dots \\ \xi_n &= \beta_0 + \beta_1 X_{n1} + \beta_2 X_{n2} + \dots + \beta_k X_{nk}\end{aligned}$$

The OLS estimator is obtained (like in the special case of simple linear regression) by minimizing the residual sum of squares (RSS). The RSS for the multiple linear regression model is

$$Q = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_k X_{ik})^2$$

To find the minimum of Q we calculate the first partial derivative of Q for each β_j :

$$\begin{aligned}\frac{\partial Q}{\partial \beta_0} &= 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_k X_{ik})(-1) \\ \frac{\partial Q}{\partial \beta_1} &= 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_k X_{ik})(-X_{i1}) \\ \frac{\partial Q}{\partial \beta_2} &= 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_k X_{ik})(-X_{i2}) \\ \frac{\partial Q}{\partial \beta_k} &= 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_k X_{ik})(-X_{ik})\end{aligned}$$

For the minimum $\hat{\beta}$ the derivative of each equation must be zero:

$$\begin{aligned}
\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} - \dots - \hat{\beta}_k X_{ik}) &= 0 \\
\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} - \dots - \hat{\beta}_k X_{ik}) X_{i1} &= 0 \\
\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} - \dots - \hat{\beta}_k X_{ik}) X_{i2} &= 0 \\
&\dots \\
\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} - \dots - \hat{\beta}_k X_{ik}) X_{ik} &= 0
\end{aligned}$$

Alternatively, we can use matrix notation and combine the above equations into the following form:

$$X^T y - X^T X \hat{\beta} = 0$$

Rearranging this, the following expression is known as normal equations:

$$X^T X \hat{\beta} = X^T y$$

Just for illustration, the system of normal equations in expanded matrix notation is:

$$\left[\begin{array}{cccc}
n & \sum_{i=1}^n X_{i1} & \dots & \sum_{i=1}^n X_{ik} \\
\sum_{i=1}^n X_{i1} & \sum_{i=1}^n X_{i1}^2 & \dots & \sum_{i=1}^n X_{i1} X_{ik} \\
\dots & \dots & \dots & \dots \\
\sum_{i=1}^n X_{ik} & \sum_{i=1}^n X_{ik} X_{i1} & \dots & \sum_{i=1}^n X_{ik}^2
\end{array} \right] \left[\begin{array}{c}
\hat{\beta}_0 \\
\hat{\beta}_1 \\
\vdots \\
\hat{\beta}_k
\end{array} \right] = \left[\begin{array}{c}
\sum_{i=1}^n y_i \\
\sum_{i=1}^n X_{i1} y_i \\
\dots \\
\sum_{i=1}^n X_{ik} y_i
\end{array} \right]$$

The estimator $\hat{\beta}$ can be obtained by rearranging again:

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Equation 11

R-squared (R^2) Score

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Introduction: R-squared is a statistical measure that represents the proportion of variance in the dependent variable explained by the independent variables.

Description: It is a goodness-of-fit measure that ranges from 0 to 1, with values closer to 1 indicating a better fit.

Importance in ML: R-squared is used to evaluate the performance of regression models. It provides insights into how well the model explains the variance in the target variable.

Equation 12

Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Introduction: Mean Squared Error is a common loss function used to measure the average squared difference between predicted and actual values.

Description: It calculates the squared difference for each observation and then averages them to provide a metric of model accuracy.

Importance in ML: MSE is used to evaluate the performance of regression models. Lower MSE values indicate a better fit, making it essential for identifying the optimal model.

Equation 13

Mean Squared Error with L2 Regularization (MSE + L2 Reg)

$$\text{MSE}_{\text{regularized}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Introduction: Regularized MSE adds a penalty to the loss function to prevent overfitting by discouraging overly complex models.

Description: The regularization term $\lambda \sum \beta_j^2$ helps keep model parameters small, which leads to simpler models.

Importance in ML: L2 regularization is key in reducing overfitting by controlling the complexity of the model, ensuring that the model generalizes well to new data.

Equation 14

Eigenvectors and Eigenvalues

$$Av = \lambda v$$

Introduction: Eigenvectors and eigenvalues are fundamental concepts in linear algebra used to understand the structure of matrices.

Description: They represent the directions along which linear transformations act by only scaling the vectors, without changing their direction.

Importance in ML: Eigenvectors and eigenvalues are used in dimensionality reduction techniques like PCA (Principal Component Analysis), which helps reduce the number of features while preserving essential information.

Equation 15

Entropy

$$\text{Entropy} = - \sum_i p_i \log_2(p_i)$$

Introduction: Entropy is a measure of uncertainty or randomness in a dataset.

Description: It quantifies the impurity in a dataset, making it a key concept in information theory and decision trees.

Importance in ML: Entropy is used in decision tree algorithms to decide the best split at each node by measuring the purity of a dataset. Lower entropy indicates a more homogenous group of samples.

Equation 16

K-Means Clustering

$$\arg \min_S \sum_{k=1}^K \sum_{x \in S_k} \|x - \mu_k\|^2$$

Introduction: K-Means is an unsupervised learning algorithm used for clustering data into K distinct groups.

Description: It minimizes the sum of squared distances between data points and the centroid of their assigned cluster.

Importance in ML: K-Means is a fundamental clustering technique used in exploratory data analysis and segmentation tasks. It helps discover patterns and relationships in unlabeled data.

Equation 17

Kullback-Leibler (KL) Divergence

$$D_{KL}(P\|Q) = \sum_{x \in \chi} P(x) \log \frac{P(x)}{Q(x)}$$

Introduction: KL Divergence is a measure of how one probability distribution diverges from a second, reference probability distribution.

Description: It is commonly used to measure the difference between two probability distributions.

Importance in ML: KL Divergence is used in machine learning for loss functions in models like variational autoencoders. It quantifies how well the learned distribution approximates the true data distribution.

17.1 Introduction to Kullback-Leibler (KL) Divergence

The **Kullback-Leibler (KL)** divergence, also called relative entropy, measures the discrepancy between two probability distributions, P (the true distribution) and Q (the approximation). It quantifies the inefficiency of using Q to encode data from P .

For discrete distributions, the KL divergence is defined as:

$$D_{KL}(P\|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

For continuous distributions, it is expressed as:

$$D_{KL}(P\|Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

Key properties of KL divergence: - It is not symmetric: $D_{KL}(P\|Q) \neq D_{KL}(Q\|P)$. - It is non-negative: $D_{KL}(P\|Q) \geq 0$, with equality only when $P = Q$.

17.2 Example Calculation

We compute $D_{KL}(P\|Q)$ for two discrete distributions: - $P(x_1, x_2) = [x_1, 1-x_1]$ - $Q(x_1, x_2) = [0.5, 0.5]$

The KL divergence formula becomes:

$$D_{KL}(P\|Q) = x_1 \log \frac{x_1}{0.5} + (1 - x_1) \log \frac{1 - x_1}{0.5}$$

We evaluate this for $x_1 = \{0.1, 0.2, 0.3, 0.4, 0.5\}$.

17.3 Results Table

The table below shows the calculations for D_{KL} , $P(x_1)$, $P(x_2)$, and their respective log terms.

Value of x_1	D_{KL}	$P(x_1)$	$\log\left(\frac{P(x_1)}{Q(x_1)}\right)$	$P(x_2)$	$\log\left(\frac{P(x_2)}{Q(x_2)}\right)$
0.1	0.3681	0.1	-1.609	0.9	0.5878
0.2	0.1927	0.2	-0.9163	0.8	0.4700
0.3	0.08228	0.3	-0.5108	0.7	0.3365
0.4	0.02014	0.4	-0.2231	0.6	0.1823
0.5	0.0000	0.5	0.0000	0.5	0.0000

17.4 Observations

1. D_{KL} decreases as x_1 approaches 0.5, where $P = Q$.
2. The terms $P(x_i) \log \frac{P(x_i)}{Q(x_i)}$ reflect the penalty for mismatches between P and Q .

Equation 18

Log Loss

$$-\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

Introduction: Log Loss, or logistic loss, is a loss function used for binary classification tasks.

Description: It measures the performance of a classification model whose output is a probability value between 0 and 1.

Importance in ML: Log Loss is crucial in evaluating classification models where the output is a probability. Lower log loss indicates a more accurate model, especially for probabilistic predictions.

Equation 19

Support Vector Machine (SVM) Objective

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w \cdot x_i - b))$$

Introduction: SVM is a supervised learning model used for classification and regression tasks.

Description: It finds the hyperplane that best separates different classes by maximizing the margin between them.

Importance in ML: SVMs are effective for high-dimensional spaces and are used in classification problems where the decision boundary is non-linear.

Equation 20

Singular Value Decomposition (SVD)

$$A = U\Sigma V^T$$

Introduction: SVD is a matrix factorization technique used in linear algebra.

Description: It decomposes a matrix into three other matrices, revealing the intrinsic structure of the data.

Importance in ML: SVD is used in dimensionality reduction, data compression, and noise reduction. It is also a core algorithm behind recommendation systems.

Equation 21

Lagrange Multiplier

$$\max f(x) ; g(x) = 0$$

$$L(x, \lambda) = f(x) - \lambda * g(x)$$

Introduction: The Lagrange multiplier is a method for finding the local maxima and minima of a function subject to equality constraints.

Description: It introduces a new variable, λ , which helps in optimizing a function while considering the constraint.

Importance in ML: Lagrange multipliers are used in optimization problems with constraints, such as training machine learning models with regularization terms.

Equation 22

The Human Equation

HUMAN

Introduction:

Description:

Importance in ML:

Appendix

Matrix Formulation for Linear Separability

To check for linear separability, we can use a matrix formulation that involves the following:

1. Data Matrix: Let X be the matrix of feature vectors, where each row corresponds to a data point. The size of X will be $n \times d$, where n is the number of data points and d is the number of features.
2. Label Vector: Let y be the vector of labels, where each entry represents the class label of a data point (typically +1 or -1 for binary classification). The size of y will be $n \times 1$.
3. Linearly Separable Condition: The data points are linearly separable if there exists a vector w and a bias term b such that:

$$y_i (w^T x_i + b) > 0, \quad \forall i = 1, 2, \dots, n$$

This condition means that all data points from one class are on one side of the hyperplane and all data points from the other class are on the opposite side.

Matrix Representation of Linear Separability In matrix form, we can express the linear separability condition as:

$$y \cdot (Xw + b\mathbf{1}) > 0$$

Where: - X is the data matrix of size $n \times d$, - y is the label vector of size $n \times 1$, - w is the weight vector of size $d \times 1$, - b is the bias term, and - $\mathbf{1}$ is a vector of ones of size $n \times 1$.

SVM's Role in Linear Separability For SVM specifically, the goal is to maximize the margin between the two classes while ensuring that the classes remain linearly separable. The margin is the distance between the closest points from each class (known as support vectors) to the hyperplane. SVM tries to find the hyperplane that maximizes this margin, ensuring that the data is not only separable but also well-separated by the hyperplane.

Mathematically, this is represented as:

$$\text{Maximize } \frac{2}{\|\mathbf{w}\|}$$

Subject to the constraints:

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad \forall i$$

Where: - y_i is the class label for the i -th point (either +1 or -1), - \mathbf{x}_i is the feature vector for the i -th data point, - \mathbf{w} is the weight vector, - b is the bias term.

This ensures that the points are classified correctly with the maximum margin.

Conclusion: Matrix Formulation To summarize, linear separability in machine learning is about finding a hyperplane that can separate the classes of data in feature space. This is often formulated mathematically as:

$$\mathbf{y} \cdot (X\mathbf{w} + b1) > 0$$

For a binary classification problem, a classifier like SVM finds the optimal hyperplane that separates the classes with the largest margin, ensuring that the data is linearly separable.

Epilogue - Guidance and Test Materials

Applying Musk Rules to Writing Your Book on the “25 Key Equations in Machine Learning”

1. Always Question the Requirements

- **Book Scope and Audience:** Reevaluate the core purpose and audience of the book. Who is your ideal reader? Do you really need to cover 25 equations, or would a different number be more impactful? Is there any unnecessary content that isn't truly essential to your goal of effective teaching and reference?
- **Key Questions to Ask:** Do each of the sections contribute to the core learning objectives? Could some be combined, shortened, or presented differently to be clearer and more engaging?

2. Try Very Hard to Delete Parts or Processes

- **Delete Unnecessary Content:** Go through each chapter and decide if every subsection or detail is truly required. Are there parts of the explanations, historical contexts, or code examples that could be condensed or cut without losing value?
- **Focus on Core Concepts:** Sometimes, less is more. Identify and focus on the key insights of each equation. For example, rather than going deep into every historical development of the normal distribution, focus on how the normal distribution is applied directly in machine learning.

3. Simplify or Optimize, But Not Too Early

- **Simplification After Completion:** In the initial phase, write freely to capture all your thoughts and ideas. Don't worry too much about the length or complexity. Once you have a draft, then simplify—remove verbose explanations, simplify complex code

snippets, or summarize complex derivations to keep it engaging.

- **Iteration in Drafts:** Write each chapter in an exploratory way first, then optimize the language and reduce complexity in subsequent drafts.

4. Move Faster

- **Accelerate Writing with LLMs and Tools:** Use tools like large language models (LLMs) to help generate content, summarize, and write initial drafts faster. You've already indicated using Mathpix, LaTeX, and R Markdown—lean into these tools to keep your process fast and flexible.
- **Set Short-Term Goals:** Use milestones to maintain momentum. Instead of working on all 25 chapters at once, work on drafting 5 chapters in a week. Use techniques like “pomodoro” to keep the pace fast without burnout.

5. Finally: Automate

- **Automate Repetitive Tasks Last:** Once content is ready, automate formatting and typesetting. Given that you are using R Markdown and LaTeX, continue leveraging those tools for repetitive typesetting tasks, but do this after you've streamlined the core content.
- **Automation for Consistency:** Automate consistency checks for formulas and code. This could be done with scripts that verify equation formatting, syntax checking for code snippets, or tools that flag unformatted variables.

Specific Strategies to Speed Up the Book Creation:

1. Outline Consolidation:

- Review your current outline and condense or merge sections that overlap. Consider if all 25 equations deserve equal weight. You could provide detailed explanations for 10-15 equations while giving more general summaries for the rest if they are less foundational.

2. Minimize and Prioritize:

- Focus on the most impactful content. Given the target audience, ask yourself: do they need a historical deep dive, or is the practical application what matters most? Emphasize applications, examples, and visuals over long theoretical discussions unless absolutely necessary.

3. Visuals and Examples:

- To make learning engaging, visuals and practical examples are key. Automate the creation of code snippets and figures that are repetitive. Tools like R and Python

scripts can automatically generate the visuals needed for each equation, saving considerable time.

4. Leverage Existing Content:

- Since you are creating this as a reference book, consider using existing open-source implementations of the concepts as examples, rather than creating everything from scratch. This will save time while still providing value.

5. Peer Review and Feedback Early:

- Incorporate early feedback to streamline revisions. Share your draft with a small group of readers and ask for their feedback on clarity and structure. This will help you catch issues early, saving time during the editing phase.

6. Use a Versioned Approach:

- Set smaller, achievable versions (like v0.3, v0.4, etc.) where each iteration refines content. For example, v0.3 might be the initial full draft, v0.4 could focus on reducing complexity, and v0.5 could involve formatting and visual consistency. This will help you move in structured steps rather than getting bogged down in perfecting each detail upfront.

Moving Faster, in Practice:

- **Sprint Write Sections:** Work on chapters in sprints. Set a specific time limit for each chapter, such as three days per equation. The goal is to create momentum and avoid stagnation.
- **Lean on a Support System:** Delegate part of the process. If possible, outsource parts like copy-editing, formatting, or even generating diagrams. The faster you move past the basic draft stage, the better you can refine and make the book truly impactful.

Conclusion:

The Musk Rules encourage challenging the purpose of everything you're including, focusing on essential content, and using tools and automation wisely, but only after the core content is strong. By questioning, deleting, simplifying, moving fast, and automating thoughtfully, you can accelerate the process of creating your book without compromising quality. Focus on the main principles that will make your book a genuinely valuable reference—clear, effective, and engaging explanations, visualizations, and examples—while avoiding unnecessary details that slow you down.

How do these recommendations feel in the context of your goals for this book? Would you like help refining a specific section or strategy for moving faster?

22.1 TEMP — R code - Python code - test area

22.2 testing Python Integration (temp section - to remove later)

Create a variable `x` in the Python session:

```
x = [1, 2, 3]
```

Access the Python variable `x` in an R code chunk:

```
py$x
```

```
## [1] 1 2 3
```

Create a new variable `y` in the Python session using R, and pass a data frame to `y`:

```
py$y <- head(cars)
```

Print the variable `y` in Python:

```
print(y)
```

```
## {'speed': [4.0, 4.0, 7.0, 7.0, 8.0, 9.0], 'dist': [2.0, 10.0, 4.0, 22.0, 16.0, 10.0]}
```

22.3 Define Numbers in R

Let's define the numbers we will use in both R and Python:

```
# Define a vector of numbers in R
define_numbers <- c(1, 2, 3, 4, 5)
```

22.4 R Code Block

This is a simple R code block that calculates the sum of 5 numbers:

```
# R code to sum 5 numbers
r_sum <- sum(define_numbers)
print(paste("The sum of the numbers in R is:", r_sum))

## [1] "The sum of the numbers in R is: 15"
```

22.5 Python Code To Pass *

```
py$n2 <- define_numbers
```

22.6 Python Code Block

This is a simple Python code block that calculates the sum of 5 numbers:

```
# Python code to calculate the sum of numbers defined in R

# Define a list of numbers
numbers = [1, 2, 3, 4, 5]

# Retrieve the numbers passed from R using reticulate's 'py' object
total_sum = sum(numbers)
total_sum2 = sum(n2)

# Print the result
print(f"The sum of the numbers 'total sum' is: {total_sum}")

## The sum of the numbers 'total sum' is: 15
```

```
print(f"The sum of the numbers 'n2' is: {total_sum2}")
```

```
## The sum of the numbers 'n2' is: 15.0
```

Appendix II - test area

Code Listing 22.1: Test Code Block Listing

```
# Your R code here  
summary(cars)
```

```
##      speed          dist  
##  Min.   : 4.0   Min.   :  2.00  
##  1st Qu.:12.0   1st Qu.: 26.00  
##  Median :15.0   Median : 36.00  
##  Mean    :15.4   Mean   : 42.98  
##  3rd Qu.:19.0   3rd Qu.: 56.00  
##  Max.    :25.0   Max.   :120.00
```


List of Figures

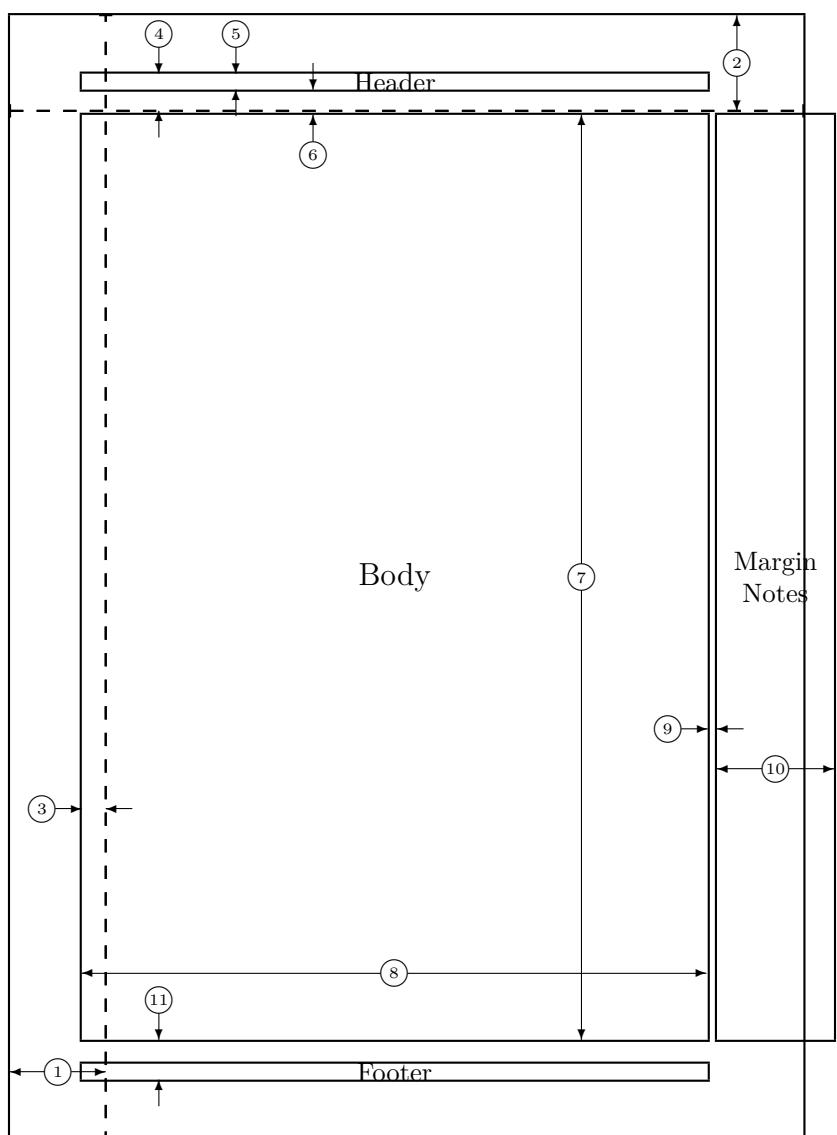
1.1	Gradient Descent on Polynomial Function - Iteration vs Function Value	7
1.2	Gradient Descent on Polynomial Function - Trace the Path on the Polynomial	8
2.1	Summary of Normal Distribution with PDF and CDF	16
2.2	Carl Friedrich Gauss	16
2.3	Pierre-Simon Laplace	16
3.1	Default caption for all figures	28
4.1	Summary of Alternate Sigmoid Like Curves	39
5.1	Default caption for all figures	46
5.2	Default caption for all figures	47
6.1	Default caption for all figures	71
7.1	Venn diagram representing the intersection and complements of sets A and B.	79
7.2	Probability Graph Ω bridged by commutative intersection $P(A \cap B) = P(B \cap A)$	84
7.3	Probability Graph displaying inherent symmetry around Ω	85
7.4	Marginals and sums network of Traditional Bayes' Space	88
7.5	Red Grid for Bayes Calc (edit later)	91
8.1	MLE Flow Chart	130
9.1	Default caption for all figures	139
9.2	Default caption for all figures	141

Code Block Listings

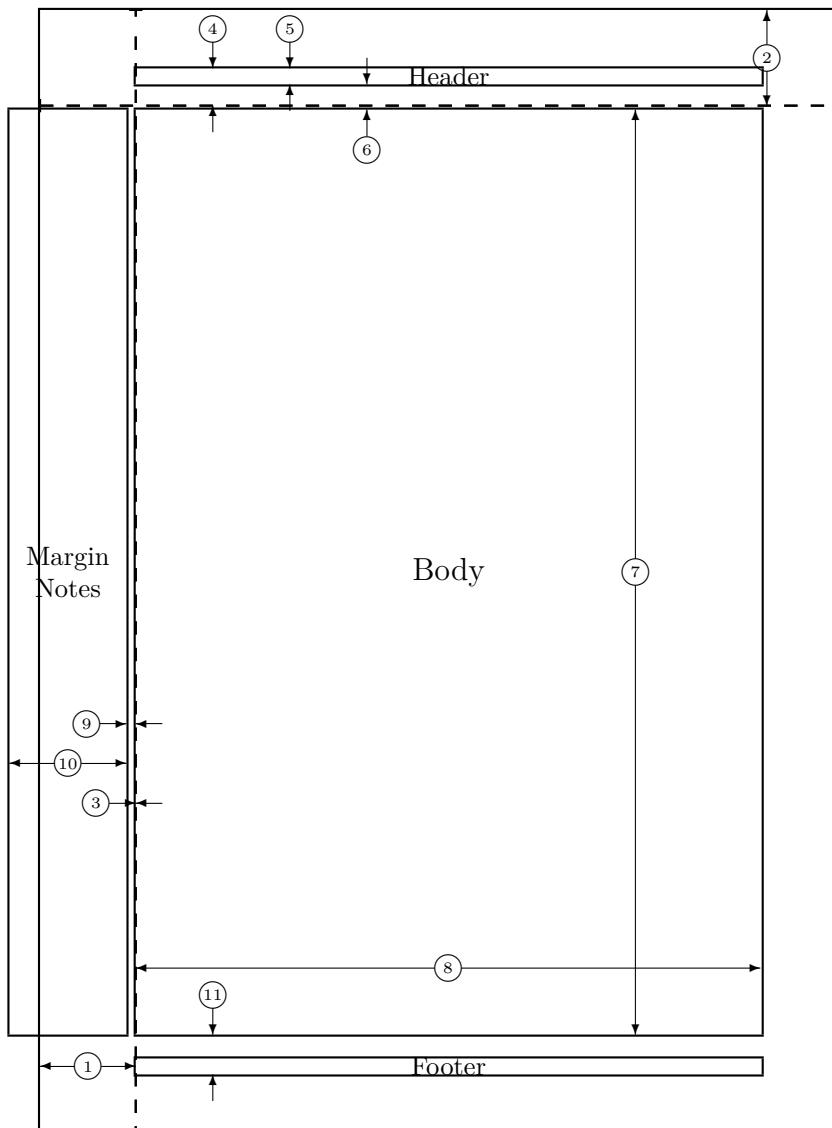
1.1	Gradient Descent: Over a Polynomial	6
1.2	Gradient Descent: Objective Function vs. Iteration	6
1.3	Gradient Descent: Plot the Path of Descent	8
22.1	Test Code Block Listing	189

Galley Sheet

22.7 GALLEY Sheet info



1 one inch + \hoffset	2 one inch + \voffset
3 \oddsidemargin = -18pt	4 \topmargin = -28pt
5 \headheight = 12pt	6 \headsep = 19pt
7 \textheight = 696pt	8 \textwidth = 471pt
9 \marginparsep = 7pt	10 \marginparwidth = 88pt
11 \footskip = 30pt	\marginparpush = 7pt (not shown)
\hoffset = 0pt	\voffset = 0pt
\paperwidth = 597pt	\paperheight = 845pt



- | | |
|--|--|
| 1 one inch + \hoffset | 2 one inch + \voffset |
| 3 \evensidemargin = 0pt | 4 \topmargin = -28pt |
| 5 \headheight = 12pt | 6 \headsep = 19pt |
| 7 \textheight = 696pt | 8 \textwidth = 471pt |
| 9 \marginparsep = 7pt | 10 \marginparwidth = 88pt |
| 11 \footskip = 30pt
\hoffset = 0pt
\paperwidth = 597pt | \marginparpush = 7pt (not shown)
\voffset = 0pt
\paperheight = 845pt |

22.8 GALLEY Sheet info - ends

References

