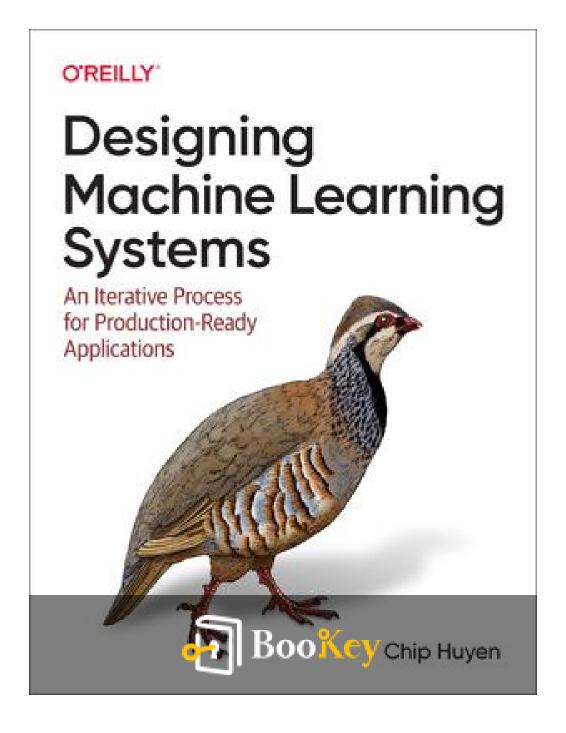
Designing Machine Learning Systems PDF (Limited Copy)

Chip Huyen







Designing Machine Learning Systems Summary

Building practical and scalable ML applications effectively.

Written by Books OneHub





About the book

In "Designing Machine Learning Systems," Chip Huyen takes readers on a transformative journey that demystifies the complex world of machine learning deployment, emphasizing that it is not just about algorithms and models but about creating robust, scalable, and maintainable systems. This book provides a comprehensive framework for building machine learning systems from the ground up, integrating practical insights and real-world case studies that highlight the challenges and best practices in the field. By blending theoretical principles with hands-on guidance, Huyen empowers readers—from budding data scientists to seasoned engineers—to navigate the entire lifecycle of machine learning projects, ensuring they can effectively bridge the gap between research and production and drive impactful outcomes in their organizations. With its clear structure and engaging narrative, this guide invites you to unlock the potential of machine learning by mastering its design principles.





About the author

Chip Huyen is a prominent figure in the field of machine learning and artificial intelligence, renowned for his expertise in designing and implementing robust machine learning systems. With a solid educational background, including a Master's degree from Stanford University, Huyen has contributed significantly to both academia and industry. He has worked on various machine learning projects at high-profile tech companies and is an advocate for effective machine learning practices. In addition to his technical insights, Huyen is known for his ability to communicate complex concepts in a clear and accessible manner, making him a sought-after speaker and educator. His book, "Designing Machine Learning Systems," reflects his deep understanding of the challenges and best practices in developing scalable and efficient machine learning solutions.







ness Strategy













7 Entrepreneurship







Self-care

(Know Yourself



Insights of world best books















Summary Content List

Chapter 1: When and When not to Use Machine Learning

Chapter 2: Understanding Machine Learning Systems

Chapter 3: Designing ML Systems in Production

Chapter 4: Summary

Chapter 5: Mind vs. Data

Chapter 6: Data Sources

Chapter 7: Data Formats

Chapter 8: Data Processing and Storage

Chapter 9: Summary

Chapter 10: Sampling

Chapter 11: Labeling

Chapter 12: Class Imbalance

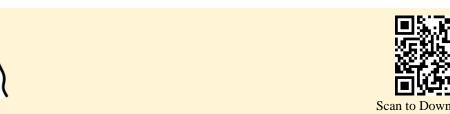
Chapter 13: Data Augmentation

Chapter 14: Summary

More Free Book

Chapter 15: Learned Features vs. Engineered Features

Chapter 16: Common Feature Engineering Operations



Chapter 17: Data Leakage

Chapter 18: Engineering Good Features

Chapter 19: Summary

Chapter 20: Model Selection

Chapter 21: Model Training

Chapter 22: Model Offline Evaluation

Chapter 23: Summary



Chapter 1 Summary: When and When not to Use Machine Learning

In the first chapter of "Designing Machine Learning Systems" by Chip Huyen, readers are introduced to the practical considerations and principles that underpin the deployment of machine learning (ML) in real-world scenarios. The chapter sets out to demystify the hype surrounding ML and provide a nuanced understanding of its applicability, highlighting both its strengths and limitations.

Firstly, it asserts that while ML is a powerful tool for a myriad of tasks, it is essential to ascertain whether ML is the necessary or cost-effective solution for a given problem. This evaluation begins with a clear definition of machine learning: an approach that learns complex patterns from existing data to make predictions on unseen data. Each element of this definition is explored to shed light on when ML is appropriate.

1. **Learning Capability**: For an ML system to be effective, it must possess the capacity to learn from data. Unlike relational databases that require explicit programming to define relationships, ML systems can derive patterns independently from the data provided. An example is predicting rental prices based on listing characteristics, whereby the system learns the correlation between inputs and outputs through training.



- 2. **Complex Patterns**: ML shines in scenarios where underlying patterns are too complex to articulate directly. Simple problems (like matching zip codes to states) can be solved via straightforward algorithms, whereas intricate relationships, such as those influencing rental prices, benefit from ML's capability to ascertain these patterns autonomously.
- 3. **Existence of Patterns**: The effectiveness of ML hinges on the availability of discernible patterns within the data. Predictive modeling for stock prices thrives because of identifiable patterns, whereas attempting to predict a fair die's outcomes would yield no meaningful results due to a lack of patterns.
- 4. **Data Availability**: ML systems depend heavily on the availability of data. They are trained on historical data to draw patterns and make predictions. In scenarios where relevant data is absent or insufficient, ML solutions are not feasible.
- 5. **Predictive Context**: The utility of ML encompasses predictive problems it excels in framing questions that necessitate forecasting future events or behaviors. This re-framing broadens the scope of various problems, making them suitable for ML implementation.
- 6. **Generalizability of Patterns**: For the predictions made by an ML model to be actionable, the unseen data must share similar characteristics



with the training data. There is an inherent risk involved; a model trained on outdated data may falter when applied to current scenarios.

The chapter further discusses optimal conditions for ML applications. Problems that are repetitive, scalable, and subject to dynamic changes are more amenable to ML solutions. This adaptability allows the models to refine their predictions continuously based on new data inputs.

Conversely, the text cautions against the misuse of ML in scenarios where it may cause ethical dilemmas, be inefficient compared to simpler solutions, or where the consequences of a single miscalculation could be severe.

Additionally, it emphasizes the potential to decompose complex problems into smaller, manageable components, where ML can effectively solve subproblems that do not warrant a full-scale ML approach.

In terms of practical applications, the explosion of ML usage across both enterprise and consumer applications is highlighted. From personalized recommendations employed by platforms like Amazon and Netflix to ML's role in predictive typing on smartphones, the technology has embedded itself in various facets of daily life. The chapter also notes the unique demands of enterprise applications, which often prioritize accuracy and operational efficiency over rapid response times.

ML applications extend to diverse sectors, addressing challenges such as





fraud detection, pricing optimization, demand forecasting, customer acquisition cost reduction, churn prediction, automated customer support, brand monitoring, and even healthcare diagnostics. Each use case exemplifies how ML can generate value by either optimizing processes or enhancing decision-making.

In summary, Chapter 1 underscores the strategic evaluation required prior to embarking on an ML project, articulating when ML is suitable and the multifaceted conditions that bolster its effectiveness. By grounding discussions in actionable insights and real-world examples, Chip Huyen equips readers with a foundational understanding necessary for successfully integrating machine learning into their applications.





Critical Thinking

Key Point: Evaluating the Necessity of Machine Learning
Critical Interpretation: As you journey through life, consider the
important lesson of evaluating whether a complex solution is really
necessary for the problems you face. Just as Chip Huyen emphasizes
in his exploration of machine learning, not every situation warrants the
sophisticated capabilities of ML. You might find yourself striving for
the newest technology or the latest trend, but like in designing
effective ML systems, it's crucial to first ask yourself if the problem at
hand truly requires such an intricate approach. This realization can
lead you to simplify your challenges, prioritize efficiency, and focus
on solutions that are not only effective but also cost-effective. By
discerning the right tools for the right problems, you can navigate life
with clarity and purpose, ensuring that your efforts yield the greatest
impact with the resources you have.





Chapter 2 Summary: Understanding Machine Learning Systems

Understanding machine learning (ML) systems is crucial for effective design and deployment, especially as the field evolves from research to production. A fundamental distinction must be made between academia and industry practice. While ML knowledge often stems from coursework and research, real-world challenges in deploying these systems diverge significantly from theoretical frameworks, necessitating a different approach to design and operation.

- 1. **Objectives and Stakeholders**: Unlike research, which typically prioritizes model performance, production environments involve a myriad of stakeholders, each with divergent objectives. For instance, in a restaurant recommendation system, ML engineers may push for a complex model to maximize user engagement, while sales teams might prioritize profits from advertising, and product managers may focus on minimizing latency to boost orders. These conflicting goals demand a collaborative effort to balance everyone's interests and find a model that achieves a satisfactory compromise, illustrating the complexity of real-world ML deployment.
- 2. **Computational Priorities**: The emphasis shifts from training to inference in production settings. Research often focuses on fast model training and high throughput, while practical applications prioritize low



latency and fast inference. Latency, the time taken from receiving a query to producing a result, is a critical factor, as demonstrated by research indicating that slight increases can significantly affect user engagement metrics and overall business outcomes. In contrast, throughput—the number of queries processed over time—takes precedence in research environments, welcoming methods like batching that increase efficiency at the cost of responsiveness.

- 3. **Data Dynamics**: The nature of data varies drastically between research and production. Research datasets are typically clean and static, designed to facilitate benchmarking and experimentation. In contrast, production data is often messy, dynamic, and subject to constant change, raising challenges such as bias, privacy concerns, and the need for ongoing data management. This requires professionals to develop systems adept at adapting to an evolving data landscape, accounting for real-time data influxes as well as regulatory issues.
- 4. Fairness and Ethical Considerations: In production, fairness becomes a pressing concern, yet many researchers overlook it during the preliminary stages of model building, focusing instead on achieving high accuracy. However, biases embedded in algorithms have real-world implications, such as unfair treatment in lending or hiring practices. Ignoring these ethical dimensions can lead to widespread discrimination. Recent surveys reveal that a minuscule share of organizations actively work on addressing



algorithmic bias, signaling a gap that needs urgent attention.

5. **Interpretability**: The question of how to interpret ML decisions is essential. Unlike traditional software, where outputs can often be traced back to discrete logic, ML systems, especially those relying on complex algorithms, can function as black boxes. Users, whether managers or end-users, require explanations for decisions to foster trust and confidence. There is an evident lack of focus on explainability in research, which stands in stark contrast to the necessity for it in industry applications. Without effective interpretability, organizations risk deploying models that users do not understand or trust.

As the landscape of machine learning continues to evolve, the gap between research and production will demand innovative practices and frameworks that can bridge these worlds effectively. Consequently, there is a growing trend toward emphasized collaborative efforts between technical development and ethical considerations, integrating diverse perspectives while leveraging data and algorithms to create fair and efficient outcomes.

In summary, successful deployment of ML systems in production depends on recognizing and adapting to the inherent differences between theoretical research practices and the intricate realities of operational environments, where varying objectives, data dynamics, and ethical considerations all play a pivotal role in shaping system design and functionality. Organizations





must prioritize flexibility, fairness, and interpretability to navigate this complex landscape effectively.

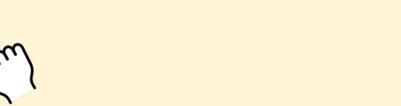


Critical Thinking

More Free Book

Key Point: Balancing diverse stakeholder objectives is essential in real-life applications.

Critical Interpretation: Imagine you're part of a team tasked with revamping a local restaurant's marketing strategy. Just as in the world of machine learning systems, your group consists of members from different departments—each with their own priorities. The data scientists are eager to build a sophisticated recommendation algorithm that captures every user nuance, but the sales team is planting their foot firmly on maximizing revenue through targeted ads, while the product managers stress the importance of quick and efficient interactions with customers. This scenario mirrors the need for collaboration and finding common ground among various interests, underscoring that in life—like in ML design—success often hinges on your ability to unify contrasting perspectives and forge productive compromises. By embracing this principle, you can foster harmonious teamwork and create solutions that satisfy everyone involved, ultimately leading to a more prosperous outcome.



Chapter 3: Designing ML Systems in Production

In the rapidly evolving field of machine learning (ML), the engineering challenges faced in system design have seen significant advancements in just a few years. For instance, when the groundbreaking BERT model was introduced in 2018, it was often deemed too large and complex for practical use, boasting an impressive 340 million parameters and a size of 1.35 GB. However, it wasn't long before BERT and its derivatives became integral to nearly every English search on Google, showcasing the pace of progress in ML system deployment.

The design of ML systems involves a structured approach to defining their various components—such as interfaces, algorithms, data, infrastructure, and hardware—with the overarching goal of satisfying specific requirements. To ensure successful system development, it is crucial to establish these requirements upfront, which can differ by use case. Nonetheless, there are four foundational characteristics that most ML systems should embody: reliability, scalability, maintainability, and adaptability.

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey



Why Bookey is must have App for Book Lovers



30min Content

The deeper and clearer interpretation we provide, the better grasp of each title you have.



Text and Audio format

Absorb knowledge even in fragmented time.



Quiz

Check whether you have mastered what you just learned.



And more

Multiple Voices & fonts, Mind Map, Quotes, IdeaClips...



Chapter 4 Summary: Summary

In the journey of deploying Machine Learning (ML) systems, Step 6 emphasizes the importance of business analysis. The performance of ML models must be rigorously evaluated against established business objectives to extract valuable insights. These insights are crucial as they guide stakeholders in determining which projects may warrant further investment and which should be discontinued due to lack of productivity.

This chapter outlines the intricate nature of bringing an ML system into a production environment. It highlights the distinctions between ML initiatives in research settings and those within traditional software engineering frameworks. Such systems are ambitious undertakings, characterized by their complexity, which is composed of varied components and necessitates collaboration among diverse stakeholders. The contexts in which these systems operate can serve a wide array of tasks, whether targeted at consumers or enterprises, each presenting unique challenges and requirements.

As the landscape of ML continues to evolve, so do the tools and best practices that support ML systems. Acknowledging that it is impractical to cover every possible aspect of ML in production, this chapter aims to provide a foundation that is broadly applicable across various tasks. The intention is to help readers navigate potential obstacles and foster





preparedness in evaluating and planning ML usage in their projects. There remains an open invitation for feedback on areas that may have been overlooked.

The complexities of ML systems can be distilled into simpler building blocks. With a high-level overview now established, the next chapters will delve into the fundamental components of these systems, starting with data engineering. For readers who may find various challenges conceptual, subsequent examples will serve to provide clarity and make these complexities more tangible.

- 1. It's crucial to analyze model performance against specific business goals to derive actionable insights, allowing for data-driven project management decisions.
- 2. ML systems in production differ significantly from research projects or traditional software engineering, indicating a need for specialized approaches and methodologies.
- 3. The evolving nature of ML technologies and methodologies necessitates that users stay informed and adaptable to embrace best practices.
- 4. By understanding that ML systems are built from simpler components, stakeholders can more easily navigate their complexities and foster collaboration for successful implementation.



Chapter 5 Summary: Mind vs. Data

The significance of data in developing machine learning (ML) systems cannot be overstated; it is foundational to their success and effectiveness in real-world applications. The field of data engineering plays a crucial role in this regard, focusing on the processes of collecting, handling, and processing the ever-growing volumes of data essential for ML systems. For practitioners already versed in data engineering, a direct dive into advanced topics like sampling and labeling for training data would be advantageous.

Understanding the relationship between the 'mind'—represented by algorithms and intelligent design—and 'data' is vital as discussions around this dichotomy have stirred significant debate among experts. Over recent years, leading academics have engaged in robust discussions about the interplay between these two forces in ML. Although some advocate for the dominance of data, arguing its critical role in driving performance improvements, others caution against placing too much emphasis on it at the expense of innovative algorithm design.

Notable figures such as Dr. Judea Pearl posit a more skeptical view on data-centric approaches, advocating for the importance of intelligent design over mere quantity of data. He highlighted his concerns about the sustainability of a data-heavy paradigm and predicted that professionals adhering strictly to this approach may find themselves obsolete in the





coming years. In contrast, influential researchers like Professor Christopher Manning underline the risks posed by massive data processed with simplistic algorithms, asserting that such models can lead to subpar results.

The debate largely revolves around the sufficiency of finite data. While ample data is indeed essential for training effective ML systems, it is the qualitative aspects of data and the computational underpinnings that foster better machine learning performance. Esteemed figures like Dr. Monica Rogati emphasize that high-quality data serves as the bedrock of data science, which includes machine learning. Therefore, an organization aiming to enhance its products or processes through data science must prioritize both the improvement of data quality and its volume.

The recent evolution of ML models illustrates the increasing reliance on data, evident in the exponential growth of datasets utilized in training models. For example, where models like GPT-2 leveraged 10 billion tokens, GPT-3 escalated this requirement to 500 billion tokens within a year. This upward trend signifies that as the complexity and capabilities of models increase, so too does the necessity for larger and richer datasets.

The current understanding across both academic and industrial landscapes is clear: regardless of future debates on the primacy of data versus intelligent design in ML, the undeniable truth remains that high-quality, abundant data is crucial for the development of effective and powerful machine learning





systems.





Critical Thinking

Key Point: The Vital Role of Quality Data in Decision Making Critical Interpretation: As you navigate through life, consider how the principles of machine learning can reflect your approach to decision-making. Just as high-quality data is essential in developing effective ML systems, you, too, can enhance your choices by prioritizing the quality of information you gather. In your daily endeavors, whether at work or in personal matters, seek out reliable data, insights, and experiences before making informed decisions. By valuing quality over quantity, you empower yourself to make choices that are not only well-founded but also lead to more significant outcomes in your journey. This approach not only enriches your understanding but also helps you build a more robust foundation for achieving your goals, just like an ML model thrives on the right data.





Chapter 6: Data Sources

In Chapter 6 of "Designing Machine Learning Systems," Chip Huyen explores the complexities and considerations involved in sourcing data for machine learning (ML) systems. A pivotal point emphasized is that while the volume of data has greatly increased in recent years, simply having more data does not guarantee improved model performance. In fact, low-quality data, which may be outdated or inaccurately labeled, can significantly hinder model efficacy.

- 1. **Understanding Data Sources**: ML systems rely on diverse data sources, each with unique characteristics and processing requirements. Recognizing the origins and nature of your data can enhance its utility. Various data types include user input, system-generated data, internal databases, and third-party data, all of which serve distinct purposes and require tailored processing approaches.
- 2. **User Input Data**: This type of data is explicitly provided by users and can take many forms such as text, images, or files. However, user input is

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey

Fi

ΑŁ



Positive feedback

Sara Scholz

tes after each book summary erstanding but also make the and engaging. Bookey has ling for me.

Fantastic!!!

I'm amazed by the variety of books and languages Bookey supports. It's not just an app, it's a gateway to global knowledge. Plus, earning points for charity is a big plus!

ding habit o's design al growth

José Botín

Love it! Wonnie Tappkx ★ ★ ★ ★

Bookey offers me time to go through the important parts of a book. It also gives me enough idea whether or not I should purchase the whole book version or not! It is easy to use!

Time saver!

Masood El Toure

Bookey is my go-to app for summaries are concise, ins curated. It's like having acc right at my fingertips!

Awesome app!

**

Rahul Malviya

I love audiobooks but don't always have time to listen to the entire book! bookey allows me to get a summary of the highlights of the book I'm interested in!!! What a great concept !!!highly recommended! Beautiful App

* * * * *

Alex Wall

This app is a lifesaver for book lovers with busy schedules. The summaries are spot on, and the mind maps help reinforce wh I've learned. Highly recommend!



Chapter 7 Summary: Data Formats

In contemporary applications of machine learning, understanding the diverse sources and formats of data is crucial, especially for tasks such as recommendation systems, which leverage user behavior data to enhance relevance. Companies often gather extensive information, encompassing media activities and browsing habits, categorized into various demographic profiles, facilitating insights like the correlation of brand preferences. However, difficulties arise in data storage due to the varying access patterns from different sources, coupled with the need for both economical and rapid accessibility.

- 1. **Data Serialization and Storage Formats**: Once acquired, storing data involves serializing it into formats suitable for persistence. Serialization is the process of transforming a data structure or object state into a storable format that can be reconstructed later. Various serialization formats, each with specific characteristics, exist; these include text-based formats like JSON and CSV, and binary formats like Parquet and Avro. Choosing the right format hinges on criteria such as human readability, efficiency, and how the data will be accessed.
- 2. **Understanding JSON**: One of the most prevalent formats is JSON (JavaScript Object Notation), known for its simplicity and readability. Defined by key-value pairs, JSON is adaptable, supporting structured data,



such as an object with nested fields, and unstructured blobs of text. Its wide support across programming languages makes it a go-to choice for many developers.

- 3. **Row-major vs. Column-major Formats**: Data storage formats can be broadly classified into row-major and column-major. CSV is a common row-major format, where consecutive elements of a row are stored adjacently, making it efficient for row-based data access but potentially less effective for column-based tasks. In contrast, Parquet, a column-major format, stores consecutive elements of a column together. This distinction affords advantages depending on the task; row-major formats excel during data writes, while column-major formats enhance performance during reads, especially when accessing specific features.
- 4. **Performance with Pandas and NumPy**: The nuances of these formats have broader implications in programming libraries such as Pandas and NumPy. Pandas is optimized for columnar data representation, which means operations on a DataFrame (a fundamental structure in Pandas) are significantly faster when accessing data by columns rather than by rows. In comparison, NumPy's ndarray defaults to row-major storage, allowing for efficient access patterns that align with its inherent data structure.
- 5. **Text vs. Binary Storage** The choice between text and binary formats also plays a significant role in data efficiency. While text files, like CSV



and JSON, are human-readable and straightforward to manipulate, they tend to consume more storage than their binary counterparts. For instance, storing a number like one million in a text file requires 7 bytes, while a binary representation takes just 4 bytes, showcasing the compactness of binary files. Conversion from text to binary formats, as exemplified by transforming a CSV dataset into Parquet, can lead to substantial reductions in file size, optimizing storage and performance.

In summary, understanding the intrinsic properties of various data formats and their implications is vital for designing efficient machine learning systems. Decisions related to data storage, serialization, and access patterns can greatly influence the performance of machine learning applications, ultimately determining their effectiveness in leveraging large datasets to produce insightful outcomes.





Chapter 8 Summary: Data Processing and Storage

In Chapter 8 of "Designing Machine Learning Systems" by Chip Huyen, the narrative delves into the critical aspects of data processing and storage, focusing on how these elements shape the efficacy of machine learning (ML) applications.

- 1. **Data Formats and Storage Efficiency**: The chapter highlights the efficiency of data storage formats, presenting a comparative discussion between CSV and Parquet formats. For instance, while a CSV file may occupy 14MB, the same data in Parquet format reduces to 6MB, and AWS recommends Parquet for its performance benefits, noting it is up to twice as fast to unload and can save up to six times the storage space on platforms like Amazon S3.
- 2. **Types of Database Processing**: Two primary forms of data processing are detailed: transactional processing and analytical processing.

 Transactional databases, designed for processing real-time transactions, adhere to ACID properties—Atomicity, Consistency, Isolation, and Durability—which are essential for ensuring successful transaction management in applications like food ordering and money transfers. In contrast, analytical databases, optimized for trend analysis and data aggregation, are geared towards answering complex queries such as average ride prices over a specified period.



- 3. Evolution of Database Terminology. The chapter notes the outdated nature of the terms OLTP (OnLine Transaction Processing) and OLAP (OnLine Analytical Processing), citing the technological advancements that allow modern databases, like CockroachDB and Apache Iceberg, to bridge the functional gap between transactional and analytical needs. The emergence of a decoupled storage and processing model further signifies a departure from traditional paradigms, allowing data to be stored centrally while query processing is mediated through distinct optimization layers.
- 4. **Data Availability and Processing Speed** The text elaborates on how the term "online" has evolved, now encompassing not just internet connectivity but also real-time data processing. For businesses, quick data access is increasingly vital, compelling the need for real-time (or near real-time) data processing solutions—especially crucial for both transactional and analytical tasks.
- 5. Understanding ETL Processes: The ETL (Extract, Transform, Load) process emerges as a foundational framework for data preparation in ML. Extraction involves sourcing data, which might come in varied formats and conditions. The transformation phase engages in significant processing, from cleaning to standardizing diverse data inputs, and finally, the loading stage determines how and where to store the processed information. This structured approach is essential in shaping the quality and utility of the data



used in ML systems.

- 6. Structured vs. Unstructured Data: The distinctions between structured and unstructured data are well-articulated. Structured data conforms to a predetermined schema, facilitating straightforward analysis, but poses challenges when schema changes are necessary. Unstructured data, while offering flexibility and easier storage options, requires additional effort to derive meaningful insights. The chapter emphasizes that the choice between these formats has substantial implications for data management practices, as well as the architecture of repositories, where structured data typically resides in data warehouses and unstructured data is housed within data lakes.
- 7. The Shift from ETL to ELT: Finally, the transition from the traditional ETL approach to the more contemporary ELT (Extract, Load, Transform) reflects the rapid evolution in data management owing to the proliferation of sources and the increased volume of data. By allowing for the initial storage of raw data without stringent schema constraints, organizations can remain agile in data processing. However, this practice also raises concerns about storage costs and data retrieval efficiency, signaling a potential return to structured approaches as cloud technologies evolve and standardize infrastructures.

In conclusion, Chapter 8 of Chip Huyen's book provides a comprehensive





overview of essential data processing principles, illustrating how they are pivotal in building robust machine learning systems. This detailed analysis not only underscores the technological advancements in data management but also highlights the ongoing need for flexible and efficient approaches to handle the dynamic landscape of data.





Critical Thinking

Key Point: The Importance of Efficient Data Storage Formats Critical Interpretation: Imagine you are in a world where the complexity of your daily choices—from what to eat for dinner to how to allocate your finances—is dictated by how efficiently you can access relevant information. Just as the chapter explains the advantages of using Parquet over CSV for data storage in machine learning applications, you too can apply this lesson to your life by prioritizing efficiency in your decision-making processes. By streamlining the way you gather and organize personal information—say, using apps to consolidate your grocery lists or budgeting tools—you're effectively minimizing time spent sifting through overwhelming options. This conscious choice not only simplifies your life, making decision-making faster and less stressful, but also empowers you to act on your ambitions with clarity and effectiveness. Your journey towards a more organized life can mirror the advancements in data management, allowing you to harness the power of efficiency to unlock greater potential in your personal and professional endeavors.





Chapter 9: Summary

In Chapter 9 of "Designing Machine Learning Systems," the author, Chip Huyen, delves into the crucial role of data in the development of intelligent systems. A prevailing misconception persists that advanced algorithms can compensate for a lack of extensive data. However, the effectiveness of prominent machine learning (ML) systems such as AlexNet, BERT, and GPT clearly highlights that the recent advancements in ML are heavily dependent on access to vast datasets. Consequently, it is imperative for ML practitioners to be adept at managing and processing large volumes of data.

To better navigate the complexities of handling data, Huyen introduces essential principles of data engineering, which he believes every aspiring ML professional should master. These principles encompass critical aspects such as sourcing data from various origins, selecting appropriate data formats, and managing both structured and unstructured data. By equipping readers with these foundational skills, Huyen aims to prepare them for the challenges of processing seemingly daunting data in real-world production scenarios.

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey



Read, Share, Empower

Finish Your Reading Challenge, Donate Books to African Children.

The Concept



This book donation activity is rolling out together with Books For Africa. We release this project because we share the same belief as BFA: For many children in Africa, the gift of books truly is a gift of hope.

The Rule



Your learning not only brings knowledge but also allows you to earn points for charitable causes! For every 100 points you earn, a book will be donated to Africa.

Chapter 10 Summary: Sampling

In the realm of machine learning, the importance of data is often emphasized, yet the warning against potential biases in data cannot be overstated. Biases can arise from various sources, including the processes of collecting, sampling, and labeling data. Historical data may carry inherent human biases, and models trained on such data risk perpetuating these biases. Thus, while data is indispensable, it should be approached with a healthy degree of skepticism.

Sampling is a crucial phase in the machine learning workflow; it is frequently underestimated in traditional ML training. Sampling is employed at various stages throughout a machine learning project, whether it involves creating training data from real-world datasets, generating splits for training, validation, and testing, or monitoring events within an ML system. At times, sampling proves to be necessary due to limitations in accessing complete data or the infeasibility of processing immense datasets due to time, computational power, or financial constraints. In other scenarios, sampling fosters efficiency by allowing quicker and cheaper experimentation, such as running initial tests on a smaller data subset to gauge a model's potential before a full-scale evaluation.

Understanding the different sampling techniques is essential to mitigate possible biases and to select methods that enhance sampling efficiency.



There are two primary categories of sampling methodologies: non-probability sampling and random sampling.

- 1. Non-Probability Sampling encompasses methods where data selection does not adhere to probability criteria. Common techniques include:
- Convenience Sampling, where readily available data is selected for its accessibility.
- **Snowball Sampling**, which involves selecting new samples based on existing ones by leveraging connections between entities in a dataset.
- **Judgment Sampling**, where experts determine which samples to include based on their insights.
- **Quota Sampling**, where samples are chosen based on predefined quotas without randomization.

Despite its convenience, non-probability sampling often leads to selection biases, making it unsuitable for creating dependable ML models. It is prevalent in scenarios such as language modeling and sentiment analysis, where datasets are often collected for their ease of accessibility rather than their representativeness. For instance, sentiment analysis might rely heavily on user-generated content from platforms like Amazon or IMDB, neglecting broader demographics.

2. On the other hand, Random Sampling includes methods designed to give every sample an equal chance of being selected. Examples of these methods



are:

- **Simple Random Sampling** allows for straightforward implementation but can overlook rare categories, which may not appear in selected samples, leading to inadequate representation of certain classes.
- **Stratified Sampling** improves on this by partitioning the overall population into distinct groups (strata) and ensures proportional representation during sampling. However, it is not without limitations, particularly when it comes to multi-label tasks.
- **Weighted Sampling** assigns probabilities to samples based on their significance, allowing for enhanced representation of underrepresented categories. This technique helps address discrepancies between sampled data and real-world distributions.
- 3. **Importance Sampling** is a powerful technique used when direct sampling from a desired distribution (P(x)) is impractical, and an alternative distribution (Q(x)) is used instead. This method allows for efficient sampling and estimation by weighting obtained samples appropriately, which is especially useful in scenarios like reinforcement learning where evaluating actions can be costly.
- 4. **Reservoir Sampling** addresses challenges associated with streaming data, enabling the selection of a fixed number of samples from an ongoing data stream without prior knowledge of the total volume of incoming data. The algorithm ensures that every incoming sample has an equal chance of



being included in the final selection, allowing for memory-efficient and randomized sampling that does not require the total dataset to fit in memory.

In conclusion, while sampling is an essential component of machine learning workflows, it requires careful consideration of the various biases and methodologies at play. Understanding and employing effective sampling strategies ensures that models are developed on representative data, minimizing bias and enhancing the reliability of results. Thus, a judicious approach to sampling is vital for the advancement of machine learning methodologies.





Chapter 11 Summary: Labeling

In the domain of machine learning (ML), labeling data is a crucial yet challenging task. Despite the potential for unsupervised learning, most ML applications currently rely on supervised models that require quality labels. The efficacy of these models hinges on both the volume and accuracy of the labeled data. While some tasks naturally yield labels—like determining whether a user clicks an advertisement—many fields confront difficulties in obtaining reliable labels due to various constraints.

Acquiring hand labels is fraught with complications. The first challenge is cost; expert annotators, for instance, are expensive and often in limited supply. Hand labeling can also breach privacy regulations, making it necessary to keep sensitive data confidential. Furthermore, the process is inherently slow: the time required for annotation can scale linearly with the amount of data, impacting the model's ability to adapt to new requirements. For instance, updating a sentiment analysis model to include new emotional categories necessitates relabeling existing data, which consumes additional time and resources.

Moreover, a phenomenon termed label multiplicity poses significant concern. When relying on multiple annotators across varied expertise levels, discrepancies in labeling can arise, leading to confusion over which labels should be used for training. Such variances highlight the necessity for clear





problem definitions and well-structured guidelines for annotators to minimize disagreements. Additionally, it is vital to maintain data lineage—tracking the origins of labeled data—enabling teams to identify biases and rectify issues stemming from mislabeled samples.

To tackle the challenges of labeling data, several techniques have been developed, including weak supervision, semi-supervised learning, transfer learning, and active learning. Weak supervision utilizes heuristics to assign labels, relying on predefined rules without needing extensive hand-labeled datasets. Tools like Snorkel facilitate generating labels through heuristics, enabling organizations to annotate vast datasets quickly while maintaining privacy. The idea is to create labeling functions that encode expert knowledge and apply them efficiently across data samples.

Semi-supervised learning capitalizes on a small set of labeled data to generate additional labels. By leveraging structural assumptions about the data, this approach employs methods like self-training, which involves making predictions on unlabeled data based on initial labeled datasets. Perturbation-based methods assume that slight variations in the input should not alter the label, allowing researchers to augment training datasets effectively.

Transfer learning involves adapting models trained for one task (often with abundant data) to perform another task that might be data-scarce. This



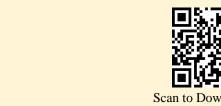


approach has gained traction due to its cost-effectiveness and ability to boost performance in downstream applications. For instance, language models pre-trained on extensive datasets can be fine-tuned for specific objectives, facilitating the development of robust ML systems even with limited resources.

Active learning emerges as a strategy to enhance the efficiency of data labeling. It enables models to select which unlabeled samples should be annotated, focusing on those that will maximize learning benefits.

Techniques such as uncertainty sampling and query-by-committee allow models to identify the most informative data points to label, thereby optimizing the overall training process.

In summary, each of these labeling techniques offers unique advantages in overcoming the inherent difficulties in acquiring high-quality labeled data. Organizations can strategically apply these methods to enhance their machine learning systems while ensuring cost-efficiency, quality, and adaptability in their models.



Critical Thinking

Key Point: Embrace creativity in overcoming challenges.

Critical Interpretation: Just as machine learning models face hurdles in labeling data, you will encounter obstacles in your own life. The key takeaway from this chapter is the importance of being resourceful and adaptable. Instead of letting difficulties like lack of resources or time constraints hinder your progress, think of innovative solutions—like utilizing your existing knowledge or collaborating with others to navigate the challenges you face. This mindset encourages you to leverage creativity, whether it's finding alternatives to achieve your goals, learning to ask for help when needed, or rethinking your approach in response to changing circumstances. By embracing this adaptable spirit in your personal and professional endeavors, you'll enhance your ability to tackle any problem that comes your way, transforming obstacles into opportunities for growth.





Chapter 12: Class Imbalance

In Chapter 12 of "Designing Machine Learning Systems" by Chip Huyen, the discussion revolves around key principles of active learning, class imbalance, and methods to handle these challenges in machine learning models.

- 1. **Active Learning**: Active learning is emphasized as a powerful approach where models select uncertain examples for labeling, either from a stationary pool of unlabeled data or from real-time incoming data. This adaptability allows models to learn effectively in changing environments. Active learning is especially beneficial when data is continuously evolving, enhancing the model's capability to adapt.
- 2. Class Imbalance: Class imbalance poses significant challenges in machine learning, particularly in classification tasks where the distribution between classes is drastically uneven. For instance, lung cancer detection datasets may contain a minuscule percentage of positive cases compared to normal instances. This imbalance can hinder model learning for minority

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey



unlock your potencial

Free Trial with Bookey







Scan to download



funds for Blackstone's firs overcoming numerous reje the importance of persister entrepreneurship. After two successfully raised \$850 m

Chapter 13 Summary: Data Augmentation

Data augmentation is an essential set of techniques aimed at increasing the volume and diversity of training data used in various machine learning applications. Initially adopted primarily in scenarios with limited datasets, like medical imaging, recent advancements demonstrate its utility even when ample data is available. The enhanced robustness against noise and adversarial attacks achieved through data augmentation has led to its integration as a standard practice in both computer vision and natural language processing (NLP) tasks.

One of the principal forms of augmentation involves simple label-preserving transformations. In computer vision, this entails common modifications such as cropping, flipping, rotation, and erasing parts of an image while retaining its original label. This approach allows for meaningful variations of the same object. For instance, rotating an image of a dog preserves its identity, rendering the output still valid for training. Prominent machine learning frameworks such as PyTorch and Keras facilitate such augmentations efficiently. A notable implementation aspect mentioned by Krizhevsky et al. indicates that many of these transformations can be computed concurrently with model training, leading to substantial computational efficiency.

In NLP, a similar approach can be employed by randomly substituting words in sentences with synonyms or closely related terms. This modification can



significantly expand the training dataset by generating multiple variations of similar meanings. For example, modifying the phrase "I'm so happy to see you" to "I'm so glad to see you" or "I'm very happy to see you" exemplifies how a single sentence can produce diverse training outputs.

The second category of data augmentation, known as perturbation, also maintains label integrity but introduces noise or deceptive alterations to models' training data. This method capitalizes on the fact that neural networks are sensitive to changes in input data. Perturbations can mislead models into misclassifying data points, a tactic used in adversarial attacks. Studies have shown that even minor pixel alterations can lead to significant misclassification rates in well-known datasets like CIFAR-10 and ImageNet. Consequently, adding noisy samples to the training set can heighten model recognition of its decision boundaries' weaknesses, thereby bolstering overall reliability against adversarial attempts. Although less common in NLP due to the inherent fragility of language, perturbation techniques—such as BERT's random token replacement—have demonstrated slight performance improvements, validating their role in enhancing model robustness.

Data synthesis forms the third category of data augmentation, which is particularly valuable in environments where data collection is slow, expensive, or comes with privacy concerns. While a comprehensive synthesis of training data remains largely aspirational, there are tangible





strategies to produce synthetic data that can meaningfully contribute to model training. In NLP, templates can serve as a powerful means to generate training queries efficiently. For example, a template such as "Find me a [CUISINE] restaurant within [NUMBER] miles of [LOCATION]" can yield diverse, contextually relevant queries which can substantially augment the original dataset.

In computer vision, data synthesis can involve generating new data points by blending existing examples while preserving their labels. Techniques like "mixup," which combines two labeled examples to create a new instance, have been shown to enhance generalization, diminish overfitting, and increase robustness to adversarial challenges. Pioneering approaches utilizing generative methods, such as CycleGAN, have also illustrated their potential to improve model performance significantly in complex tasks like CT segmentation.

Overall, leveraging data augmentation, whether through simple transformations, perturbation techniques, or synthesis strategies, can result in substantial improvements in the performance and reliability of machine learning models across various applications. This multidimensional enhancement underscores its growing importance as an invaluable tool in the design of machine learning systems.



Chapter 14 Summary: Summary

The essence of machine learning hinges significantly on the quality and curation of training data, which remains a crucial foundation for the efficacy of modern ML algorithms. Regardless of the sophistication of the algorithms employed, the performance is inherently linked to the quality of the training data. Bad training data leads to poor algorithm performance; hence, investing time and resources into crafting and refining this data is imperative for meaningful learning outcomes.

Once training data is established, the next logical step involves feature extraction, a process essential for training machine learning models effectively. While some argue that large models may require vast amounts of data to perform optimally—making smaller datasets less applicable—it is vital to experiment with datasets of varying sizes to discern their respective impacts on model performance.

In the realm of machine learning, multilabel tasks are significant, where a single example can possess multiple labels, indicating the complexity and richness of data interpretations. Additionally, when labeling is straightforward, it often negates the necessity for domain expertise, underscoring the importance of clarity in the labeling process.

Several methodologies and studies have emerged to tackle the challenges



associated with training data collection and labeling. For instance, the work on Snorkel presents a framework for rapid training data creation through weak supervision, showcasing innovative ways to generate labeled datasets without exhaustive manual efforts. Furthermore, the literature on cross-modal data programming illustrates how these methods can expedite medical machine learning initiatives.

Effective strategies to enhance model performance also include co-training, which is vital for situations involving labeled and unlabeled data, ensuring a streamlined approach to maximizing dataset utility. To confront issues like class imbalance—a common phenomenon in datasets that can skew model training—researchers have explored various strategies, particularly around resampling techniques and focal loss methodologies. These approaches facilitate better handling of imbalanced datasets, thus improving the robustness of machine learning models.

In conclusion, the quality of training data is paramount in machine learning; thus, ensuring its adequacy through various innovative strategies is essential for building models that deliver reliable and high-performing results. Investing in effective data creation and curation techniques is crucial not only for enhancing algorithmic performance but also for advancing the overall advancement of machine learning applications. This focus on data quality will lay the groundwork for more sophisticated and efficient machine learning systems in the future.



Chapter 15: Learned Features vs. Engineered Features

In the realm of machine learning systems, one critical but often overlooked issue is data leakage, which can significantly undermine the integrity of a model when placed in production. This chapter emphasizes the importance of detecting and avoiding data leakage while also delving into the principles of effective feature engineering, highlighting the balance between learned features and engineered features.

1. Understanding Feature Engineering: While deep learning has revolutionized the way features can be learned from data, the need for feature engineering remains pertinent, especially in contexts where manual extraction of features is required. Many students express skepticism regarding the necessity of feature engineering with the advancements in deep learning, which is sometimes referred to as feature learning. They are correct that deep learning can automate the extraction of numerous features, yet it has not yet reached a point where all critical features can be independently generated.

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey



ness Strategy













7 Entrepreneurship







Self-care

(Know Yourself



Insights of world best books















Chapter 16 Summary: Common Feature Engineering Operations

In the process of designing machine learning systems, particularly in model development, feature engineering plays a pivotal role. This crucial practice involves selecting and extracting relevant features from data to improve model performance. For complex tasks, the number of features can escalate to millions, especially in domains such as video recommendations or fraud detection, where domain expertise is essential.

Feature engineering encompasses various operations critical for transforming raw data into useful information. Here are several key techniques that are widely employed:

1. Handling Missing Values One of the initial challenges faced with real-world data is missing values. Missing data can be classified into three categories: Missing Not At Random (MNAR), where the absence of a value is related to that value; Missing At Random (MAR), where missingness is due to another observed variable; and Missing Completely At Random (MCAR), where missingness occurs without relation to any observed data. To manage these, practitioners can either delete data (which risks losing valuable information) or impute missing values with suitable alternatives (like the mean, median, or a default value). Each method has its advantages and potential pitfalls, including bias introduction or data loss.



- 2. **Scaling**: In many machine learning models, features can vary significantly in scale, which can lead to skewed predictions. For instance, a model might misinterpret the relative significance of age versus income because of their differing scales. To counter this, feature scaling is employed, often normalizing values to a common range, such as [0, 1], or standardizing them to have zero mean and unit variance. Proper scaling not only enhances model performance but also helps prevent data leakage, a common issue associated with improper handling of training and inference datasets.
- 3. **Discretization**: Continuous variables can sometimes mislead model interpretations, especially with minor differences in values. Discretization simplifies this by converting continuous features into discrete buckets. For example, income levels can be categorized into ranges, which makes it easier for models to learn from broader groupings instead of individual data points. Choosing appropriate boundaries through techniques like basic quantiles or domain knowledge is crucial for effective discretization.
- 4. **Encoding Categorical Features**: Categorical variables present another challenge, as their potential diversity can lead to encoding issues when new categories emerge over time. While static categories can be easily represented, production environments often feature dynamic categories, which may necessitate strategies like the hashing trick—assigning hashed



values to categories to manage unforeseen entries without overwhelming the model with vast and unbounded options.

- 5. **Feature Crossing**: This strategy involves combining two or more features to capture complex relationships effectively, particularly non-linear associations. For example, merging marital status and number of children into a new feature could expose critical insights for purchase predictions. However, the downside is that feature crossing can exponentially increase the feature space, necessitating a careful balance between complexity and overfitting.
- 6. **Positional Embeddings**: Introduced in the context of transformer models, positional embeddings provide a way to encode the position of tokens within a sequence. This is particularly important in models where input sequences are processed in parallel, as it ensures understanding of token order—essential for accuracy in tasks like language modeling. The embeddings can be learned or fixed through mathematical functions, determining how positional information is conveyed to the model.

In summary, feature engineering is a multifaceted process comprised of various techniques aimed at optimizing data representation for machine learning models. By effectively handling missing values, scaling features, discretizing variables, encoding categories, utilizing feature crossing, and applying positional embeddings, practitioners can significantly enhance

More Free Book



model accuracy and efficacy, ultimately leading to better predictions and decision-making.

Concept	Description
Feature Engineering	Pivotal process in ML model development, involving selection and extraction of relevant features to improve performance.
Handling Missing Values	Addressing missing data categorized as MNAR, MAR, or MCAR through deletion or imputation methods, each with potential pitfalls.
Scaling	Normalizing or standardizing features to ensure relative significance is properly interpreted and to prevent data leakage.
Discretization	Converting continuous variables into discrete buckets to simplify model learning and improve interpretability.
Encoding Categorical Features	Using techniques like the hashing trick to manage dynamic categories in production environments without overwhelming the model.
Feature Crossing	Combining features to capture complex relationships, although it can increase feature space and risk overfitting.
Positional Embeddings	Encoding token positions in sequences for models like transformers, crucial for understanding context and order in tasks.
Conclusion	Effective feature engineering enhances model accuracy and decision-making through optimized data representation.





Critical Thinking

Critical Interpretation: As you navigate your daily challenges, consider how essential it is to identify and harness the most relevant aspects of your life, just like feature engineering in machine learning. Think of your experiences, relationships, and skills as data points waiting to be optimized. By focusing on what truly matters and refining the way you present those attributes to the world—whether through enhancing your communication skills, learning from past experiences, or understanding your emotional responses—you can dramatically amplify your potential. Just as in model development, where careful feature selection can lead to transformative insights and decisions, so too can you glean greater clarity and direction by

understanding and prioritizing the key factors that influence your

become the best version of yourself.

journey. Embrace this process of self-engineering, and watch as you





Chapter 17 Summary: Data Leakage

In the realm of machine learning, data leakage is a significant challenge that can drastically affect model performance, particularly when models show impressive results during testing but fail in real-world scenarios. This phenomenon occurs when aspects of the labels inadvertently influence the features used for predictions, which can create a misleadingly optimistic evaluation during training. Understanding data leakage is essential for developing robust machine learning systems, and it's crucial to identify its common causes and best practices.

One illustrative example of data leakage is seen in the scenario of predicting cancer from CT scans. When building a model trained on data from one hospital, it performed well on that hospital's test data but faltered when applied to data from another hospital. The reason behind this discrepancy was that the first hospital used a specific scanning machine for patients suspected of having cancer, which led the model to incorporate the effects of scanning technology rather than the actual disease indicators. Such incidents underscore the subtleties of data leakage, emphasizing the need for vigilance in model training and evaluation.

The risks associated with data leakage are not exclusive to novices; even experienced practitioners can fall victim to its pitfalls. For example, during a Kaggle competition, participants that managed to exploit leaks in the



provided data emerged as the winners, showcasing how data leakage can skew competitive results.

To mitigate the risk of data leakage, it is crucial to recognize its common sources:

- 1. **Time-Correlation Issues**: Randomly splitting time-series data can lead to future information skewing the model's predictions. To maintain a clean temporal integrity, training datasets should be split chronologically.
- 2. **Scaling Procedures**: Scaling features based on overall dataset statistics can inadvertently leak information about test data. The correct approach is to split the data first and then calculate the scaling parameters based solely on the training split.
- 3. **Handling Missing Data**: Filling in missing values using global statistics from the entire dataset can introduce leakage. Instead, statistics should be derived from the training data exclusively.
- 4. **Data Duplication Concerns**: Insufficient handling of duplicated data can result in identical examples being present in both training and validation/test sets. Always check for duplicates prior to splitting datasets.
- 5. Group Leakage: Strong correlations among grouped data can create



leakage if samples from the same group appear in different splits.

Understanding how data is organized is essential to prevent this type of leakage.

6. **Collection Process Reflectivity**: The way data is collected can lead to nuances that influence model performance, exemplified by the scanning procedures used in medical imaging. Thorough knowledge of data sourcing and processing can help mitigate risks.

Detecting data leakage requires a proactive approach throughout the entire machine learning lifecycle. Continuous monitoring of feature correlations and implementing ablation studies—removing features to assess their impact on model performance—are integral practices. Such analyses help detect features that may harbor unseen leakage when considered collectively or in conjunction with others. The significance of domain expertise in guiding these investigations cannot be overstated, as a deeper understanding of the data sources and processes involved can greatly enhance model reliability.

In conclusion, awareness and recognition of data leakage are paramount in building effective machine learning systems. By diligently identifying and controlling for potential sources of leakage during data collection, preprocessing, and model evaluation, practitioners can significantly bolster the integrity and performance of their machine learning endeavors.





Chapter 18: Engineering Good Features

In Chapter 18 of "Designing Machine Learning Systems" by Chip Huyen, the author discusses critical aspects of feature engineering that can greatly influence the performance of machine learning models. The content outlines essential principles and considerations when selecting and managing features and provides insights into assessing feature importance and generalization.

Firstly, one should monitor the addition of new features to a model. If a newly added feature significantly enhances model performance, one must evaluate whether it genuinely contributes valuable information or is merely compounding data leakage related to the output labels. Data leakage can occur when insights from the test dataset unwittingly inform the training process, such as when brainstorming new features or tuning hyperparameters based on test performance.

Generally, although adding more features can improve model performance, this does not hold universally. A multitude of features can introduce risks,

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey



Why Bookey is must have App for Book Lovers



30min Content

The deeper and clearer interpretation we provide, the better grasp of each title you have.



Text and Audio format

Absorb knowledge even in fragmented time.



Quiz

Check whether you have mastered what you just learned.



And more

Multiple Voices & fonts, Mind Map, Quotes, IdeaClips...



Chapter 19 Summary: Summary

In Chapter 19 of "Designing Machine Learning Systems" by Chip Huyen, the intricate art of feature engineering is emphasized as a critical component for the success of machine learning (ML) systems in production. The chapter highlights a fundamental tradeoff between generalization and specificity when engineering features. For example, while creating a feature such as IS_RUSH_HOUR—indicating peak traffic times—provides a broader understanding, it sacrifices the specificity offered by the HOUR_OF_THE_DAY feature. This underscores the importance of balancing generalizable features with specific ones, as relying solely on the former can result in the loss of essential data.

The chapter further outlines that success in ML depends heavily on the quality of features, prompting organizations to allocate resources toward effective feature engineering. It conveys that the process is complex and suggests that hands-on experimentation is invaluable to understanding how different features impact model performance. Engaging with the feature engineering practices of successful teams in competitions like Kaggle can also provide insights and techniques that contribute to refining this discipline.

Recognizing the collaborative nature of feature engineering, Huyen addresses the necessity of incorporating subject matter experts, who may not



be engineers, into the workflow. To facilitate this, it is crucial to develop processes that enable diverse contributions.

The text subsequently presents a compilation of best practices essential for feature engineering that all practitioners should adopt. These include:

- 1. Splitting data chronologically into train, validation, and test sets to ensure meaningful evaluation.
- 2. Oversampling data only after the initial split to avoid bias.
- 3. Scaling and normalizing data appropriately to prevent data leakage.
- 4. Utilizing statistics exclusively from the training set for scaling features and managing missing values.
- 5. Maintaining an understanding of the data collection and processing journey to monitor its lineage.
- 6. Assessing and acknowledging the importance of features to the model's outcomes.
- 7. Opting for features that demonstrate strong generalization capabilities.
- 8. Discarding features that lose relevance over time.

Importantly, the chapter reiterates that the transition to model training does not conclude the feature engineering process, as data handling and feature refinement are ongoing tasks. Continuous learning from new incoming data is crucial for improving models post-deployment. This emphasis on relentless data and feature management suggests a dynamic and iterative





approach essential for the longevity and efficacy of machine learning applications.

Overall, Huyen's insights serve as a foundational guide to understanding the complexities of feature engineering, stressing that a strong feature set is the bedrock of effective machine learning solutions.

More Free Book

Chapter 20 Summary: Model Selection

In selecting the most suitable machine learning model for a specific problem, a strategic approach is essential due to practical constraints in time and computational resources. Model selection is a thoughtful process, considering multiple potential ML and non-ML solutions. The decision-making process involves understanding the underpinning principles of machine learning models, which are fundamentally functions that convert inputs into outputs.

- 1. **Understanding Models**: A model functions as a transformation that processes input data into corresponding predictions. In machine learning, particularly supervised learning, the model is derived from data where inputs and outputs are provided. Identifying an appropriate model involves both understanding the loss function—essentially the objective function used to measure model performance—and the learning procedure that aids in parameter optimization.
- 2. **Objective Functions**: The objective function plays a pivotal role as it quantifies the error or loss during the training process. This function varies depending on the problem and may incorporate multiple approaches for supervised learning scenarios. For instance, methods like Root Mean Squared Error (RMSE) or Cross Entropy are common metrics used to evaluate the performance of models depending on the nature of the output



(scalar or distribution). Understanding and appropriately adjusting the objective function can significantly influence the learning process and the quality of the output parameters.

- 3. **Learning Procedures**: The process by which models find the best set of parameters, usually through iterative approaches like gradient descent, is crucial for training. By applying optimization techniques such as Momentum, Adam, or RMSProp, practitioners can manage how parameters adjust during the learning process. Each learning algorithm often comes with assumptions and characteristics that influence the effectiveness of the optimization.
- 4. **Framing the Problem**: Identifying the type of ML problem—be it classification, regression, or a complex system that integrates several objectives—determines the path forward. Whether it involves binary classification to categorize emails or using regression to predict house prices, appropriately framing the problem enables the selection of suitable algorithms and techniques.
- 5. Handling Multiple Objectives: When a problem encompasses conflicting goals, a systematic approach is required. For example, maximizing user engagement on a recommendation system must also consider the ethical implications of content quality. Various methods, like combining objective functions or training multiple models to focus on



different aspects of the goal, can lead to more maintainable and effective solutions.

- 6. **Embracing Flexibility in Algorithms**: The algorithm landscape is diverse and evolves, with both traditional methods and cutting-edge neural networks playing significant roles in real-world applications. While deep learning may offer significant advancements in several areas, classic ML algorithms remain relevant, particularly in production environments where explainability and computational efficiency are paramount. Therefore, it's crucial to understand the assumptions and requirements of each model to select the best-fit algorithm for the specific problem.
- 7. **Constructing Ensembles**: After deploying an individual model, improving performance can be achieved through ensemble techniques. By combining predictions from multiple models—using methods such as bagging, boosting, or stacking—one can often enhance prediction accuracy. The less correlated the base learners, the greater the potential for improvement.
- 8. **Implementing AutoML**: Automation in model selection and hyperparameter tuning can significantly ease the burden on practitioners. Techniques ranging from basic hyperparameter searches to complex neural architecture searches can optimize performance while minimizing the need for extensive manual intervention.



9. **Phased Approach to Development**: It's beneficial to adopt a phased approach to ML model development, beginning with heuristics or simple models before moving to complex systems that incorporate more elaborate data handling and optimization strategies. This gradual buildup allows for validation and iterative improvement, ensuring models are robust and effective over time.

In conclusion, selecting and refining a machine learning model is an intricate process characterized by strategic decision-making, iterative experimentation, and a lucid understanding of underlying principles. Each step, from defining the problem through evaluating performance, contributes to the development of a successful machine learning solution tailored to specific needs.

Key Section	Summary
1. Understanding Models	Models transform input data into predictions. Understanding loss functions and learning procedures is essential for model selection.
2. Objective Functions	Objective functions quantify training error and differ by problem type; common metrics include RMSE and Cross Entropy.
3. Learning Procedures	Learning involves optimizing parameters through methods like gradient descent and requires understanding different optimization techniques.
4. Framing the Problem	Identifying the type of ML problem informs algorithm selection, whether it's classification or regression.



Key Section	Summary
5. Handling Multiple Objectives	Conflicting goals require systematic methods, like combining objectives or training multiple models for better solutions.
6. Embracing Flexibility in Algorithms	A diverse algorithm landscape necessitates knowing the strengths of both traditional and modern ML methods to suit specific problems.
7. Constructing Ensembles	Ensemble techniques improve model performance by combining predictions from multiple models, leveraging their diversity.
8. Implementing AutoML	Automation in model selection and hyperparameter tuning simplifies the modeling process, enhancing performance with less manual effort.
9. Phased Approach to Development	A gradual, phased development starts with simple models, leading to more complex systems, allowing validation and improvement.
Conclusion	Selecting and refining ML models is a complex process requiring strategic decisions, experiments, and understanding of principles.





Critical Thinking

Key Point: Understanding Models as Transformations

Critical Interpretation: Imagine standing at the edge of a vast forest,
each tree representing a decision yet to be made. In your life, just like
in machine learning, the choices you face can transform your path
from uncertainty to clarity. By understanding the 'models'—the
frameworks and principles that govern your decisions—you can hone
your ability to process the inputs you receive, be they information,
emotions, or opportunities. Recognizing that every model serves to
translate your experiences into actionable outcomes empowers you to
be intentional in your choices. Just as a well-trained machine learning
model optimizes its predictions, you too can refine your
decision-making process through self-awareness and learning from
past experiences, creating a tangible impact on the trajectory of your
life.





Chapter 21: Model Training

In the realm of machine learning, as projects progress and the need for enhanced model performance arises, it's essential to transition from simpler models to more complex architectures. This shift often necessitates experimenting with various models and architectures, with a particular focus on techniques that aid in managing multiple machine learning models, especially in scalable environments. This means exploring distributed training methods, suitable for handling larger, resource-intensive datasets that may exceed memory limits.

A key aspect of distributed training is its ability to facilitate model training across multiple machines, utilizing data parallelism as the predominant strategy. In this approach, datasets are distributed across various machines for simultaneous processing, thereby accelerating the training process. However, challenges such as gradient accumulation can arise. Synchronous stochastic gradient descent (SSGD) requires all machines to synchronize their gradients, which can lead to delays if any machine underperforms. Conversely, asynchronous SGD (ASGD) updates weights separately based

Install Bookey App to Unlock Full Text and Audio

Free Trial with Bookey

Fi

ΑŁ



Positive feedback

Sara Scholz

tes after each book summary erstanding but also make the and engaging. Bookey has ling for me.

Fantastic!!!

I'm amazed by the variety of books and languages Bookey supports. It's not just an app, it's a gateway to global knowledge. Plus, earning points for charity is a big plus!

ding habit o's design al growth

José Botín

Love it! Wonnie Tappkx ★ ★ ★ ★

Bookey offers me time to go through the important parts of a book. It also gives me enough idea whether or not I should purchase the whole book version or not! It is easy to use!

Time saver!

Masood El Toure

Bookey is my go-to app for summaries are concise, ins curated. It's like having acc right at my fingertips!

Awesome app!

**

Rahul Malviya

I love audiobooks but don't always have time to listen to the entire book! bookey allows me to get a summary of the highlights of the book I'm interested in!!! What a great concept !!!highly recommended! Beautiful App

* * * * 1

Alex Wall

This app is a lifesaver for book lovers with busy schedules. The summaries are spot on, and the mind maps help reinforce wh I've learned. Highly recommend!



Chapter 22 Summary: Model Offline Evaluation

In the realm of machine learning (ML), reproducibility poses a significant challenge. Variations in frameworks and hardware can lead to non-deterministic outcomes, making it arduous to repeat experiments unless all aspects of the experimental environment are known. The current approach of treating ML models as black boxes necessitates running numerous experiments to identify optimal configurations. It is anticipated that as the field evolves, we will develop a deeper understanding of various models to reduce the dependency on extensive experimentation.

Evaluating the performance of machine learning models is crucial yet complex, especially for businesses integrating AI into their operations. Companies often struggle with understanding if their models are effective, as exemplified by one that deployed a model to detect intrusions in surveillance drones but lacked metrics to assess its performance. Without robust evaluation frameworks, it becomes challenging to pinpoint the best solutions, which can hinder management buy-in for ML projects.

In an ideal scenario, evaluation methods in the development and production stages would align; however, acquiring ground truths is not feasible in production environments. Users' feedback can help approximate these truths for recommendation systems, but biases exist. In complex tasks where direct evaluation is impossible, continuous monitoring becomes essential to gauge



shifts in model performance.

To structure model evaluation, establishing baselines is fundamental. Various baseline comparisons highlight how performance metrics alone can be misleading. For instance, evaluating models against random outputs, simplistic heuristics, or the most common class provides critical context. Recognizing that a model delivering a high F1 score can still perform poorly against a random baseline or an established heuristic reinforces the necessity of contextual evaluation.

Several evaluation methods contribute to a holistic understanding of model performance beyond mere accuracy, emphasizing robustness and fairness. Perturbation tests assess a model's resilience to input noise, while invariance tests examine the fairness of predictions across different sensitive inputs. Directional expectation tests ensure that changes in key features result in expected changes in predictions, while model calibration is essential for aligning probability predictions with actual outcomes.

Confidence measurement allows determination of when predictions are robust enough to share with users, reducing potential user distrust.

Slice-based evaluation highlights the critical importance of analyzing model performance across different demographic or behavioral subgroups, helping to identify biases and enhancing overall model effectiveness.





Analysis of model evaluation through these diverse lenses emphasizes the necessity of not only understanding overall performance but also investigating how models operate across various critical subgroups. This approach, coupled with tools for slicing data effectively, provides insights that drive improved model performance and ensures accessibility across distinct user populations.

In summary, the journey towards effective machine learning model evaluation comprises multiple core principles:

- 1. The importance of reproducibility and understanding environmental factors.
- 2. Evaluation metrics should be contextualized against baselines, providing meaningful insights into model performance.
- 3. Robust evaluation must encompass methods that ensure model fairness and resistance to bias.
- 4. Continuous monitoring and adaptation in production environments are essential, as is the need for incorporating user feedback.
- 5. Slice-based evaluation reveals performance discrepancies across various demographic segments, guiding strategic improvements and mitigating biases.

Embracing these principles lays the groundwork for building effective, trustworthy, and successful machine learning systems.





Chapter 23 Summary: Summary

In this chapter, the author delves into an integral aspect of machine learning (ML) projects that many practitioners find particularly enjoyable: the process of developing, training, and evaluating ML models. While this phase can be thrilling, it is also riddled with challenges. One significant hurdle is ensuring that models can perform effectively in large, distributed systems that comprise hundreds of millions or even billions of parameters. Achieving optimal functionality in such systems requires specialized engineering expertise.

An essential part of the model development process is the meticulous tracking and versioning of various experiments. Though many practitioners recognize its importance, managing this aspect often feels tedious.

Additionally, understanding how well a model will function in a production environment, particularly when solely relying on training data, presents difficulties. Despite thorough evaluations in a controlled setting, the true performance of a model remains uncertain until it is deployed in a real-world context.

The chapter concludes with a sneak peek into the next phases of the ML project cycle. The following chapter will focus on model deployment, while the subsequent one will cover strategies for the ongoing monitoring and evaluation of models post-deployment.



- 1. Machine learning model development is an enjoyable yet challenging aspect of the project cycle, especially when managing large distributed systems.
- 2. Effective tracking and versioning of experiments are necessary for validating models, despite often feeling cumbersome.
- 3. Evaluating model performance solely based on training data creates uncertainty regarding its real-world effectiveness.
- 4. A model's true performance can only be assessed after deployment, underscoring the importance of this subsequent step in the ML project cycle.

This chapter serves as a critical bridge between theoretical understanding and practical application, highlighting the continuous need for iteration and optimization in the realm of machine learning.

