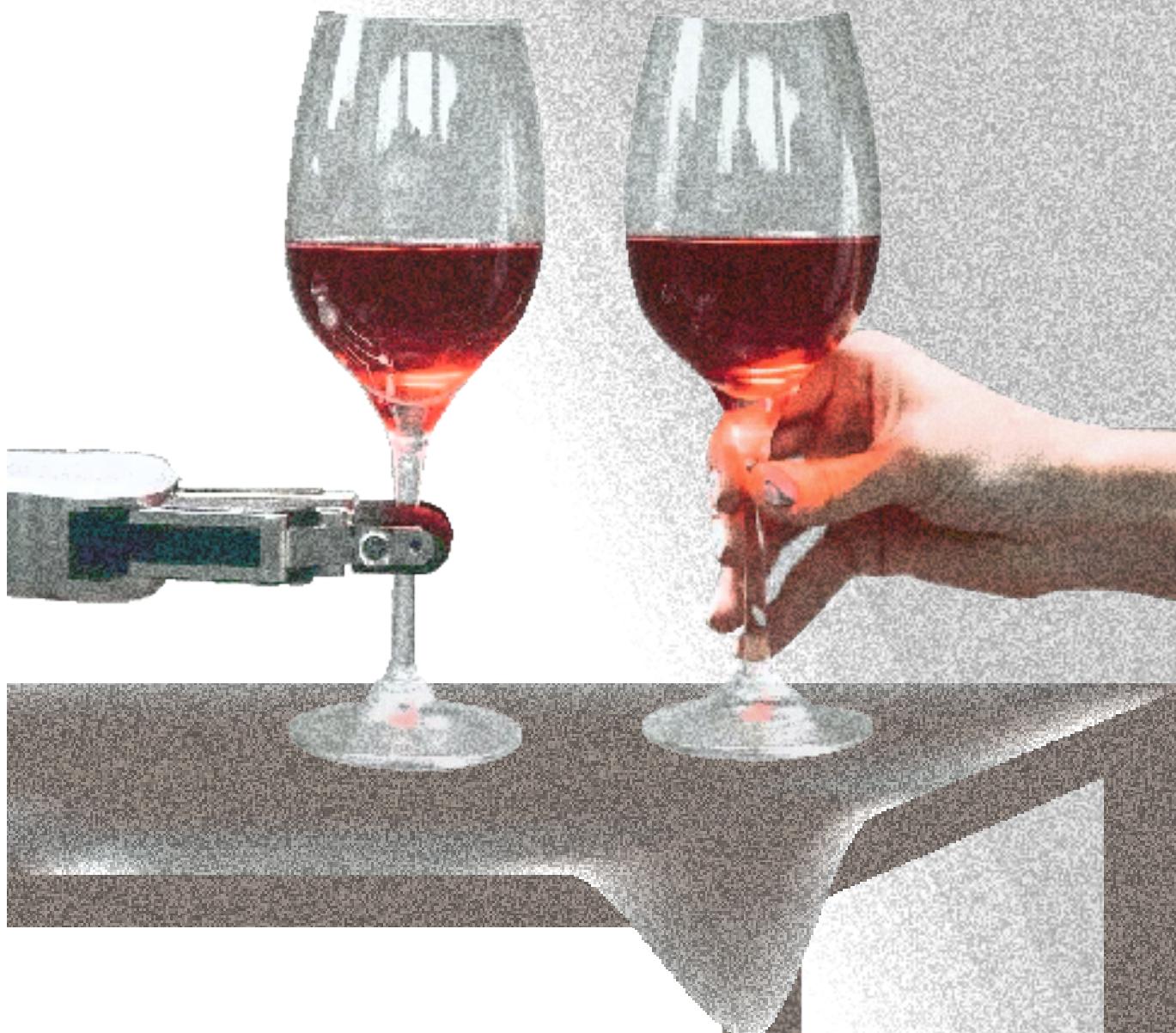


# FULL SELF-DRIVING, SKYNET AND OTHER AI MYTHS

BRAD FLAUGHER



# **Full Self-Driving, Skynet and Other Artificial Intelligence Myths**

**Modern decision making with deep machine learning models**

Brad Flaugher

February 9, 2023

*We have now accumulated sufficient evidence to see that whatever language the central nervous system is using, it is characterized by less logical and arithmetical depth than what we are normally used to.*

– John von Neumann, 1958 [1]

# Preface

This book is an attempt to organize my brain on paper. I have been thinking about decision making with computers, and the computer's relationship to the brain for too long now, and before I get dementia (I'm 35) I would like to put a stake in the ground and tell everyone where they can stick their thoughts about the future.

I am simultaneously very hopeful for the future and horrified by the way that people talk about it. I think many so-called "Data Scientists" are full of shit when they talk, and some of them know it but would like to keep their jobs. I think that many investors are snowed by fancy language and toy demos that they extrapolate into a magical future. Users of digital products are woefully ignorant about how the products work and how they fit into a larger ecosystem. I'm amused, annoyed and would like to clear things up for everyone, including myself.

I am a self-employed programmer and financially secure, these are my credentials. I have no corporate master and I am not raising money to create the future of AI, so have a limited bias towards AI boosterism. I write programs close enough to the bleeding-edge that I can see the present as fast as almost anyone. Because I am simultaneously a contractor (<https://inoxsoft.com/>), investor (<https://ventures.nextfab.com/about/>) and teacher (<https://bradflaughher.com/bootcamp.html>) I get the opportunity to see many projects, and I hope I can effectively comment on some common misconceptions and errors in thinking I see in many businesspeople, developers and investors.

So let's go on this journey together, I'll try and explain to you (with words, code and examples) how modern AI works, and when you should be scared (self-driving cars without the appropriate infrastructure) or when you should not be scared (please shut up about Skynet becoming self-aware). I'll also try to do case studies in as many arenas that I think are interesting. If you have notes for me you can reach me at [brad@bradflaughher.com](mailto:brad@bradflaughher.com).

Please skim this book, please skip around and read some stories. I think the first three chapters are important for everyone to know. The rest of the stories can be chosen a la carte. Enjoy.

*Brad Flaugher*

PS Thank you Angela Altamirano for the cover wonderful cover art. Thank you to my early supporters, William Dickson, Daniel Hoevel, Patrick Maloney, Jerrod Howlett and Jeff Zinser. A special thank you to Louis Cid and the Harvard Club of New York for giving me an early forum. Finally thank you to Dr. Péter Érdi for putting me under his wing many years ago and Dr. Elisa Esposito for her constant support.

PPS This book uses the wonderful Free and Open Source Kaobook LaTeX project (<https://github.com/fmarotta/kaobook>) and many Free and Open Source AI models for text and image generation like Stable Diffusion (<https://huggingface.co/spaces/stabilityai/stable-diffusion>) and Mann-E ([https://huggingface.co/mann-e/mann-e\\_4\\_rev-0-1](https://huggingface.co/mann-e/mann-e_4_rev-0-1)) and some closed models like GPT-3, ChatGPT and Dall-E from OpenAI (<https://openai.com/>).

# Contents

Preface	iii
Contents	iv
<b>1 History: The End of Good Old-Fashioned Artificial Intelligence</b>	<b>1</b>
1.1 The AI Textbook in 1997 . . . . .	1
1.2 Does AI Need to Know Grammar to Translate? . . . . .	1
1.3 Explicit Rules and Codified Human Knowledge . . . . .	2
1.4 IBM Tries Every Possible Chess Move . . . . .	4
1.5 Statistical Analysis of Handwriting is the Way of the Future . . . . .	4
1.6 Less Programmer Intelligence and More Data Intelligence . . . . .	5
1.7 From Explicit Rules to a Black Box, and Beyond . . . . .	5
1.8 Key Takeaways . . . . .	6
<b>2 The Regression Theory of Everything</b>	<b>7</b>
2.1 Let's Avoid Knowledge Representation! . . . . .	7
2.2 A Simple Neural Network is also a Linear Regression . . . . .	8
2.3 Dummy Variables for Dummies (Wonkish) . . . . .	8
2.4 Try That Again With 33,640,010 Parameters . . . . .	10
2.5 Multicollinearity and the End of Science . . . . .	13
2.6 Let's Test Some Random Inputs! Feature Importance and Explainability . . . . .	13
2.7 The Universal Machine Learning Workflow . . . . .	14
2.8 A Machine Learning Engineer is a Data Janitor . . . . .	16
2.9 Key Takeaways . . . . .	17
<b>3 Creativity and Decision Making with Deep Learning Models</b>	<b>19</b>
3.1 Theories of Creativity . . . . .	19
3.2 Creative Uses of Power Tools . . . . .	20
3.3 Garbage In, Garbage Out . . . . .	21
3.4 Garbage In, New Perspective Out? . . . . .	21
3.5 Concept Drift and the End of Usefulness . . . . .	22
3.6 The Impossibility of Fairness . . . . .	23
3.7 Transfer Learning Everywhere . . . . .	23
3.8 Industrial-Scale Plagiarism . . . . .	25
3.9 Working Together is Best . . . . .	25
3.10 Key Takeaways . . . . .	26
<b>4 Model Cards and Case Studies</b>	<b>28</b>
4.1 Fluid Dynamics . . . . .	28
4.2 Shakespearean Text Generation . . . . .	29
4.3 Lithium Mining Classifier . . . . .	29
4.4 Chess Playing Model . . . . .	29
4.5 Go Playing Model . . . . .	29
4.6 Large Stock Order Detector . . . . .	29
4.7 Stock Trading Bot . . . . .	29
4.8 Sperm Counter . . . . .	29
4.9 Handwriting Recognizer . . . . .	29
4.10 Liver Drug Discovery . . . . .	29
4.11 Online Dating Matcher . . . . .	29

4.12	Online Ad Server Classifier . . . . .	29
4.13	Tennis Playing Bot . . . . .	29
4.14	Hate Speech Classifier . . . . .	30
4.15	Contract Review Bot . . . . .	30
4.16	Facial Recognition . . . . .	30
4.17	Smartwatch Danger Classifier . . . . .	30
4.18	CCTV Threat Classifier . . . . .	30
4.19	War Robot . . . . .	30
4.20	Harvard Acceptance Bot . . . . .	30
4.21	Simple Credit Score . . . . .	30
4.22	Social Credit Score . . . . .	30
4.23	Artificial General Intelligence Chatbot . . . . .	30
<b>5</b>	<b>Self-Driving With Statistics</b>	<b>31</b>
5.1	Self-Driving Horses . . . . .	31
5.2	Semiautonomy is Stupid . . . . .	31
5.3	Trolley Problems . . . . .	31
5.4	Concept Drift (Reprise) . . . . .	31
5.5	Multicollinearity (Reprise) . . . . .	31
5.6	Explaining The Unexplainable in Court . . . . .	31
5.7	A Train is a Self-Driving Car, Right? . . . . .	32
<b>6</b>	<b>Unplugging Skynet</b>	<b>33</b>
6.1	AI and Human Safety . . . . .	33
6.2	Useful Incompatability . . . . .	33
6.3	Checks and Balances . . . . .	33
<b>7</b>	<b>Revolutionary for Whom?</b>	<b>34</b>
7.1	The Battle of the Assistants . . . . .	34
7.2	Employees That Are Better Than You . . . . .	34
7.3	Slow on the Uptake . . . . .	34
7.4	Free-Rider Problems . . . . .	34
7.5	When You Can't Tell The Difference . . . . .	34
7.6	Sucker Traps . . . . .	35
7.7	Dead Inside . . . . .	35
7.8	This Book is a Case Study . . . . .	35
<b>8</b>	<b>Errors and Omissions</b>	<b>36</b>
	<b>Bibliography</b>	<b>37</b>

# History: The End of Good Old-Fashioned Artificial Intelligence

1

*"Programming will be obsolete. I believe the conventional idea of "writing a program" is headed for extinction, and indeed, for all but very specialized applications, most software, as we know it, will be replaced by AI systems that are trained rather than programmed. In situations where one needs a "simple" program (after all, not everything should require a model of hundreds of billions of parameters running on a cluster of GPUs), those programs will, themselves, be generated by an AI rather than coded by hand."* Matt Welsh, 2023 [2]

## 1.1 The AI Textbook in 1997

Dr. Elaine Rich's textbook on artificial intelligence (AI), published in the 1980s, was a groundbreaking work that helped to establish many of the core concepts and techniques in the field of AI. However, the rapid advancements in AI over the past few decades have led to many of the chapters in this textbook becoming obsolete.

One of the main reasons for this is the prevalence of deep learning, big data, and large-scale statistical models in modern AI. These techniques have largely replaced the symbolic, rule-based approach to AI that was emphasized in the textbook, making many of the chapters on knowledge representation and expert systems less relevant.

Additionally, the explosion of data and the availability of powerful computing resources have made it possible to apply machine learning techniques at a scale that was previously unimaginable. This has led to the development of highly effective machine learning models that can handle complex tasks such as image and speech recognition with a high degree of accuracy, making many of the chapters on simpler machine learning techniques such as decision trees<sup>1</sup> and linear regression less relevant. [3]<sup>2</sup>

We'll discuss this history and a few examples from the "early days" of AI to help us understand where we are headed. We'll start with machine translation, then discuss chess and finally neural networks, which will be the focus of the rest of this book.

## 1.2 Does AI Need to Know Grammar to Translate?

Noam Chomsky is a linguist and philosopher who has made significant contributions to the field of linguistics with his theory of universal grammar. Chomsky believes that all human languages share a common underlying structure, and that this structure is innate to humans. He proposes that this innate structure is the result of a "language acquisition device" present in the human brain, which allows us to learn and produce language. Chomsky

1.1 The AI Textbook in 1997 . . . . .	1
1.2 Does AI Need to Know Grammar to Translate? . . . . .	1
1.3 Explicit Rules and Codified Human Knowledge . . . . .	2
1.4 IBM Tries Every Possible Chess Move . . . . .	4
1.5 Statistical Analysis of Handwriting is the Way of the Future . . . . .	4
1.6 Less Programmer Intelligence and More Data Intelligence . . . . .	5
1.7 From Explicit Rules to a Black Box, and Beyond . . . . .	5
1.8 Key Takeaways . . . . .	6

[1]: Although mathematically neural networks are decision trees <https://arxiv.org/abs/2210.05189>

[3]: Rich et al. (2009), *Artificial Intelligence*

[2]: the book is now in its third edition and unlikely to be updated as Dr. Rich as retired <https://www.cs.utexas.edu/~ear/>

also argues that the structure of language is largely independent of its content, and that the ability to produce and understand language is a fundamental aspect of human nature. His theory has been influential in the field of linguistics and has sparked much debate and research on the nature of language and its relationship to the human mind.

For English speakers or anyone who has learned English as a second language you'll have many examples of special cases, irregular verbs, bad English and former street slang that became good and proper over time. For programmers this is a nightmare, how can we codify human knowledge in a timely fashion? If we tried to write the rules of the English language in code (which many have tried to do) the rules themselves might change before we were finished writing them.

Explicitly translating languages through code is a difficult task because it requires a thorough understanding of the grammar, vocabulary, and syntax of both languages, as well as the nuances and subtleties of their respective cultures<sup>3</sup>. Simply coding rules for how to translate words or phrases from one language to another is not sufficient, as there are often multiple valid translations for a given phrase depending on the context in which it is used.

A more effective approach to translation is to use statistical techniques that rely on a large corpus of translated data, such as Canadian laws<sup>4</sup>. This type of data-driven approach involves training a machine learning model on a large dataset of translations, allowing it to learn the patterns and relationships between the languages. The model can then use this knowledge to make educated translations of new phrases or sentences, taking into account the context in which they are used.

While this approach is not perfect, it has proven to be highly effective in machine translation and can produce accurate translations even for languages that are very different from each other. The use of a large dataset of translations also allows the model to learn from the mistakes and variations present in real-world translations, further improving its accuracy.

### 1.3 Explicit Rules and Codified Human Knowledge

When we "teach" a computer to perform a task by explicitly writing down all of the rules of that task, we are really codifying human understanding.<sup>5</sup> When we codify human understanding we write down every rule that we know explicitly. For small tasks we can do this with 100 percent accuracy, and only minor headache on the part of the software developer.

For example, let's write a boring function to tell you the number of days for a given month.

3: For programmers this is a nightmare, how can we codify human knowledge in a timely fashion? If we tried to write the rules of the English language in code (which many have tried to do) the rules themselves might change before we were finished writing them.

4: They're in French AND English, which is useful data that we can use to correlate phrases and transform English to French and vice-versa.

5: Programming this way makes some software development totally boring, I almost switched my major in college to math after considering what a life would look like manually writing rules for handling "edge cases" for the rest of my natural life.

```
def days_in_month(year, month):
    if month in [1, 3, 5, 7, 8, 10, 12]:
        return 31
    elif month in [4, 6, 9, 11]:
        return 30
    elif month == 2:
        if (year % 4 == 0 and year % 100 != 0) or year % 400 == 0:
            return 29
        else:
            return 28
    else:
        return "Invalid month"
```

Writing code can be a tedious and repetitive task, especially when it comes to debugging and testing. It can be especially frustrating when you're working on a large project and you're trying to track down a specific bug that's causing the program to crash. Testing code can also be boring, as it often involves running the same tests over and over again to ensure that the code is working correctly.

Additionally, writing code can be boring because it requires a lot of concentration and focus. It can be easy to get lost in the details and lose track of time, especially if you're working on a complex problem. It can also be challenging to come up with creative solutions to problems, and it can be frustrating when your code doesn't work as expected.

While writing and testing code can be rewarding and fulfilling, it can also be a tedious and boring process. It requires a lot of patience, persistence, and attention to detail, and it can be easy to get frustrated and lose motivation. However, with practice and perseverance, it is possible to overcome these challenges and find enjoyment in the process of writing and testing code.

AI has traditionally operated by explicitly codifying human knowledge into machine-readable formats by doing the boring job of coding. This approach, which I'm calling "codified human knowledge" relies on humans to carefully structure and organize information in a way that can be understood by the AI system. The AI system then uses this structured knowledge to make decisions and perform tasks.<sup>6</sup>

However, recent advances in AI have largely ignored the knowledge representation problem and instead have focused on using statistical techniques and neural networks to automatically learn patterns and relationships in data. This approach, known as "deep learning," involves training large neural networks on vast amounts of data, allowing the AI system to make educated classifications and transformations of data without explicit human guidance.

Deep learning has proven to be highly effective in a variety of applications, such as image and speech recognition, and has contributed to the rapid progress we have seen in AI in recent years. However, the reliance on large amounts of data and the lack of transparency in these models can make it difficult to understand how they are making decisions, which can be a concern in certain applications (hence the title of this book).



Figure 1.1: "a frustrated programmer writing boring rules on his computer" made with Stable Diffusion 2.1

6: Some smart people think that AI spells the "End of Programming" <https://cacm.acm.org/magazines/2023/1/267976-the-end-of-programming/fulltext>, I think the headline is clickbait but the general idea that programming is changing along with advances in AI is true. We are all data scientists now!

## 1.4 IBM Tries Every Possible Chess Move

Deep Blue was a revolutionary computer developed by IBM that was specifically designed to play chess at the highest level. It was programmed with a vast database of chess knowledge and was able to analyze millions of positions per second.

Garry Kasparov was the reigning world chess champion at the time, and he was considered to be one of the greatest players in history. He was beaten by Deep Blue, which used rules based GOFAI (Good Old Fashioned AI) to essentially calculate every possible move and project that move to the end of the game, then choose the best position by brute force. Despite Kasparov's best efforts, he was no match for Deep Blue's computational power. In the end, Deep Blue emerged victorious, defeating Kasparov in a historic match that changed the world of chess forever.

Deep Blue was a turning point in the development of AI, but Deep Blue's methods (namely calculating every possible outcome of a Chess game to determine the best move) was not suitable for many of the world's problems. It turns out that Chess is fun, but the world is not like chess.<sup>7</sup> The "real"<sup>8</sup> future of AI was being developed elsewhere, using statistics and a toy model of the brain to solve a very practical problem for banks.

## 1.5 Statistical Analysis of Handwriting is the Way of the Future

It was the early 1990s and Yann LeCun was a researcher at Bell Labs in New Jersey. At the time, the process of reading and processing checks was a tedious and time-consuming task that was done manually by bank employees. LeCun saw the potential for using artificial intelligence to automate this process, and he began experimenting with using convolutional neural networks (CNNs) to recognize patterns in images of checks.

At the time, CNNs were a relatively new type of neural network that had been developed in the 1980s for image recognition tasks. They were inspired by the structure of the human visual system, and were able to process images in a way that was similar to how the human brain does.

LeCun's work was groundbreaking, and he was able to achieve impressive results using CNNs to process checks. By 1993, he had developed a system that was able to read and process checks with a high degree of accuracy, significantly reducing the amount of time and effort that was required to process checks manually.<sup>9</sup>

LeCun's work on using CNNs for check processing was a major milestone in the field of artificial intelligence, and it laid the foundation for the development of many other applications of CNNs in the years that followed. Today, CNNs are widely used in a variety of applications, including facial recognition, image classification, and natural language processing.<sup>10</sup>



Figure 1.2: "Garry Kasparov setting a computer on fire" made with Stable Diffusion 2.1

7: if you would like to read an example of the simplest chess engine that I could imagine that is written in a similar "codified" way as Deep Blue check out <https://github.com/thomasahle/sunfish/blob/master/sunfish.py>

8: We might go back, and try again to more explicitly code everything up, and in some cases we still need to, but from the author's perspective, we live in a deep learning/neural network world.

9: CNN digit OCR models are frequently featured in beginner training tutorials for deep learning libraries like PyTorch and Tensorflow, check one out [https://github.com/jonkrohn/DLTFpT/blob/master/notebooks/alexnet\\_in\\_tensorflow.ipynb](https://github.com/jonkrohn/DLTFpT/blob/master/notebooks/alexnet_in_tensorflow.ipynb)

10: check out Yann LeCun demonstrating a convolutional neural network in 1993 at [https://www.youtube.com/watch?v=FwFduRA\\_L6Q](https://www.youtube.com/watch?v=FwFduRA_L6Q).

## 1.6 Less Programmer Intelligence and More Data Intelligence

I think it's useful to separate the knowledge in the AI problem-space into two groups. The data and the programmer together make the programs that we use every day, and for the rest of this book I'll try and separate the discussion of the smarts of each to help us better understand the world.  
11

Early AI relied heavily on a human programmer to design, write, and debug computer programs. Good programmers needed domain expertise, problem-solving skills, logical thinking, and the ability to learn and adapt to new programming languages and technologies.

We are now in the age of big data, and everyone knows that "data is gold". Statistical AI methods that are now most prevalent rely on extracting meaningful insights and knowledge from large datasets. This involves using statistical and analytical methods to discover patterns and trends in data, and using this information to inform business decisions or solve problems.

Throughout this book I'll discuss the interaction between programmer and data, and what can go wrong. Working with big data and statistics at a large scale has given AI tremendous ability, but has made understanding and testing models infinitely more difficult. It is your authors belief that understanding the nuance of this interaction between programmer and data is essential to understanding modern AI.

## 1.7 From Explicit Rules to a Black Box, and Beyond

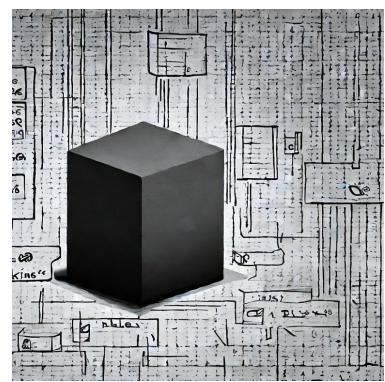
Artificial intelligence (AI) has come a long way since its inception, and the way that it makes decisions has changed significantly over time. In the early days of AI, explicit rules were used to tell the AI system what to do in certain situations. These rules were often written by humans and encoded into the system, and the AI would follow them to make decisions.

However, with the advent of deep learning, we have started to rely on a statistical understanding of the truth for AI to make decisions. Deep learning is a type of machine learning that involves training artificial neural networks on large datasets. These neural networks are able to learn patterns and relationships in the data, and can use this knowledge to make decisions.

The use of deep learning has led to the development of powerful AI that is able to perform tasks that were previously thought to be impossible for a machine to handle. For example, deep learning has led to the development of AI systems that are able to recognize faces, translate languages, and even beat humans at complex games like chess and Go.



**Figure 1.3:** "a group of computer programmers striking outside of Microsoft's offices with placards saying 'rule-based programming is boring'" made with Dall-E 2



**Figure 1.4:** "from explicit rules, to a black box and beyond" made with Stable Diffusion 2.1

While deep learning has led to significant advances in AI, it has also made it harder to debug and understand how the AI system is making its decisions. With explicit rules, it was relatively easy to understand why the AI made a particular decision. However, with deep learning, it is often difficult to understand exactly how the AI arrived at its decision. This can make it challenging to troubleshoot problems with the AI system and to ensure that it is making decisions that are fair and unbiased.

AI has come a long way since its early days, and the way that it makes decisions has changed significantly over time. While explicit rules were once used to tell the AI what to do, we now rely on a statistical understanding of the truth for AI to make decisions. This has led to the development of powerful AI that is able to perform a wide range of tasks, but it has also made it harder to debug and understand how the AI is making its decisions.<sup>12</sup>

## 1.8 Key Takeaways

- ▶ **AI is a misleading term**, we should instead talk about rules-based programming (GOFAI) and deep learning so we don't confuse ourselves, our partners and our users.<sup>13</sup>
- ▶ **Good Old-Fashioned AI (rules-based AI) is hard to create and maintain.** We used to use it for chess engines (like Deep Blue) and translation, but now programmers favor using statistical machine learning techniques.
- ▶ Good Old-Fashioned AI is not well suited to many problems, like machine translation and handwriting detection.
- ▶ Modern Deep Learning techniques use data to train models instead of humans explicitly writing rules, and the **deep learning models are often as good as the data they are trained with**.<sup>14</sup>

12: From here on out I'll try and use "GOFAI" (Good Old-Fashioned Artificial Intelligence) to describe the rules-based approach, and use "Deep Learning" or "Machine Learning" to talk about modern neural network and statistics-based techniques

13: Avoid the imprecise words "Algorithms" and "AI" and instead consider the more precise "rules-based AI" AKA "GOFAI" to describe situations where programmers explicitly write rules that they design and "deep learning" or "machine learning" where programmers create models using a dataset and let the machine figure out the rules via a neural network.

14: More *intelligent* data and maybe less intelligent programmers, or maybe programmers who are better trained in statistics and complex systems theory instead of "classical" computer science.

# The Regression Theory of Everything

# 2

*"AI Scientists disagree as to whether these language networks possess true knowledge or are just mimicking humans by remembering the statistics of millions of words. I don't believe any kind of deep learning network will achieve the goal of AGI [Artificial General Intelligence] if the network doesn't model the world the way the brain does. Deep learning networks work well, but not because they solved the knowledge representation problem. They work well because they avoided it completely, relying on statistics and lots of data instead. How deep learning networks work is clever, their performance impressive, and they are commercially valuable. I am only pointing out that they don't possess knowledge and, therefore, are not on the path to having the ability of a five-year-old child." Jeff Hawkins, 2022 [4]*

## 2.1 Let's Avoid Knowledge Representation!

The knowledge representation problem in AI is the challenge of how to formally represent knowledge in a way that a computer can understand and reason about. This typically involves creating a set of symbols, rules, and structures that can be used to represent concepts, relationships, and other types of information. The goal is to create a representation that is both expressive enough to capture all relevant aspects of the domain, and computationally tractable enough to allow for efficient reasoning and inference. There are many different approaches to knowledge representation, including logic-based, semantic networks, frames, and ontologies, each with their own strengths and weaknesses.

Deep learning techniques handle knowledge representation differently than traditional symbolic AI methods. Unlike symbolic AI, which relies on explicit and hand-coded representations of knowledge, deep learning techniques learn to represent knowledge implicitly through the use of neural networks.

In deep learning, knowledge is represented in the form of the weights of the neural network. These weights are learned through training on a large dataset and they capture the underlying relationships and patterns in the data. The neural network can then use these learned weights to make predictions, classifications or generate new data.

Deep learning models can handle large and complex datasets, and can automatically extract features from the data without the need for manual feature engineering. This makes them particularly well-suited for tasks such as image and speech recognition, natural language processing, and other areas where large amounts of data are available. However, they are not as good at explicating how they arrived at a decision, which can be a disadvantage.

2.1 Let's Avoid Knowledge Representation! . . . . .	7
2.2 A Simple Neural Network is also a Linear Regression . . . . .	8
2.3 Dummy Variables for Dummies (Wonkish) . . . . .	8
2.4 Try That Again With 33,640,010 Parameters . . . . .	10
2.5 Multicollinearity and the End of Science . . . . .	13
2.6 Let's Test Some Random Inputs! Feature Importance and Explainability . . . . .	13
2.7 The Universal Machine Learning Workflow . . . . .	14
2.8 A Machine Learning Engineer is a Data Janitor . . . . .	16
2.9 Key Takeaways . . . . .	17



**Figure 2.1:** "mdjrn-v4 a disembodied computer being force-fed data like a duck looking somewhat like a sci-fi character 8k" made with Mann-E

In summary, deep learning techniques handle knowledge representation by learning the underlying patterns and relationships in the data through the use of neural networks, which can then be used for prediction, classification, and generation tasks. In GOFAL, knowledge is held by the programmer and explicitly coded into rules, deep learning methods instead use data to guess the best rules from the relationships present in the dataset.

## 2.2 A Simple Neural Network is also a Linear Regression

A neural network can be mathematically equivalent to a regression or a decision tree under certain conditions.

A neural network is a machine learning model composed of layers of interconnected artificial neurons, which are designed to process and analyze data. They can be used for a wide range of tasks, such as image and speech recognition, natural language processing, and prediction.

A regression is a statistical method used to predict a continuous variable based on one or more input features. A linear regression, for example, is a simple neural network with one input layer, one output layer and no hidden layers. In this case, the weights of the network are the coefficients of the linear equation and the network is equivalent to a linear regression model.

A decision tree is a tree-based model used for classification and prediction tasks. It consists of a series of if-then rules that are used to make decisions based on the input data. A neural network with one input layer, one output layer and one hidden layer<sup>1</sup> is equivalent to a decision tree. The network implements a piecewise linear function which can represent the decision boundaries of a decision tree. So, under certain conditions, a neural network can be mathematically equivalent to a linear regression or a decision tree. These conditions include having one input and one output layers, and having a specific activation function in the case of a decision tree.[5]

## 2.3 Dummy Variables for Dummies (Wonkish)

Chapter Summary: It's all numbers, man. Machine learning techniques require that we turn everything (images, text, sound) into numbers and shove them into the model in the same way we use dummy variables in a simple regression. If you are satisfied with this, please skip this section. If you would like to learn a bit about the details and see some code examples, please keep reading. This section is necessarily technical, but should be approachable for anyone who has taken a college statistics class.<sup>2</sup>

Dummy variables are used in regression analysis to include categorical variables in a model. Categorical variables are variables that take on a finite number of distinct values, such as "red", "green", "blue" or "yes", "no". Since

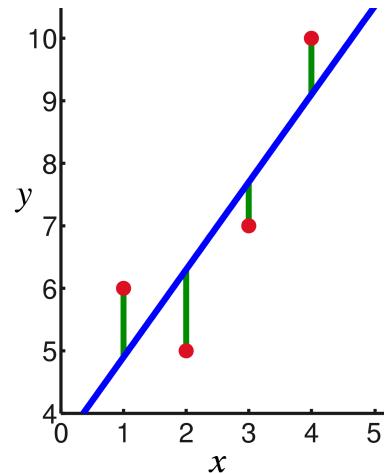


Figure 2.2: A simple linear regression, the red points are the training data, and the blue line is the regression line. If you don't understand this please read [https://en.wikipedia.org/wiki/Regression\\_analysis](https://en.wikipedia.org/wiki/Regression_analysis)

1: One with a ReLU activation, Neurons in a neural network can have different activation functions, if you don't know what it is and don't want to read about it on wikipedia, that's OK. [https://en.wikipedia.org/wiki/Activation\\_function](https://en.wikipedia.org/wiki/Activation_function)

[5]: Aytekin (2022), *Neural Networks are Decision Trees*

2: Thanks to Paul Krugman for popularizing (to me at least) the *wonkish* classifier, I just mean it's a little nerdy. But it's important.

these variables cannot be directly included in a regression model, as they are not numerical, they need to be transformed into numerical variables.

The process of creating dummy variables is also known as one-hot encoding. It involves creating a new binary variable for each category of the original variable. For example, if you have a categorical variable "color" with three categories: "red", "green", "blue", you would create three binary variables: "color\_red", "color\_green", "color\_blue". Each binary variable would take a value of 1 if the original variable is equal to the category, and 0 otherwise.

When using dummy variables in a regression, it is important to remember to include only  $n-1$  binary variables, where  $n$  is the number of categories in the original variable. This is because including all  $n$  binary variables would result in perfect multicollinearity, which is when two or more independent variables are perfectly correlated. One of the binary variables can be dropped to avoid this problem.

Dummy variables are used in regression analysis to include categorical variables in a model. The process of creating dummy variables involves creating a new binary variable for each category of the original variable and one-hot encoding it. It is important to remember to include only  $n-1$  binary variables, to avoid perfect multicollinearity.

The creation of dummy variables in a regression is analogous to preprocessing image, text and other data for a neural network for deep learning. This preprocessing is important as it ensures that the data is in a format that can be easily understood and processed by the network. The preprocessing steps for numbers, text, and images are slightly different.

For numbers:

- ▶ Normalization: It is common to normalize the input data by scaling it to have a mean of 0 and a standard deviation of 1. This helps to ensure that all input features have similar scales and prevents any one feature from dominating the network's computations.
- ▶ Imputation: Handling missing data is important, as it can negatively impact the model's performance. Common imputation techniques include replacing missing values with the mean, median, or mode of the feature.

For text:

- ▶ Tokenization: Text data must first be converted into a numerical format that can be understood by the network. This is typically done by tokenizing the text into individual words or n-grams and then encoding them as integers or real-valued vectors. A one-hot encoding exactly like the dummy variable method used in regression is also frequently used.<sup>3</sup><sup>4</sup>
- ▶ Stop-words removal: The most common words in any language like "a", "an", "the", etc. that do not contain much meaning are called stop-words, they are often removed to reduce the dimensionality of the data.

Listing 2.1: Mapping text to numbers.

```
'movies': 99,
'after': 100,
'think': 101,
'characters': 102,
'watch': 103,
'two': 104,
.films': 105,
'character': 106,
'seen': 107,
'many': 108,
'being': 109
```

3: Sometimes text is just mapped to a number! Shocking, but it works. See how it is taught in the TensorFlow tutorials [https://www.tensorflow.org/text/guide/word\\_embeddings](https://www.tensorflow.org/text/guide/word_embeddings)

4: GPT-3 uses byte-level Byte Pair Encoding (BPE) tokenization and has a vocabulary size of 50,257.

- ▶ Stemming/Lemmatization: Words that have the same meaning can be stemmed or lemmatized to reduce the vocabulary size and increase the chances of generalization.
- ▶ Vocabulary Size: Each model must choose a vocabulary size or the maximum number of tokens that it will analyze. This may cause misspellings, slang or typos to be discarded in analysis.

For images:

- ▶ Converting to RGB or Greyscale: Each image is analyzed by its pixel color value, every point on an image will either have 3 color values (red, green, blue) or one single value (on a white/black scale) if the image is analyzed in greyscale.
- ▶ Convolutions<sup>5</sup>: Pixel values are analyzed in groups that are defined by the model, since individual pixel values are only colors (or greyness) they must be combined together by the model to detect patterns like faces and stop signs. The method of convolution is defined by the model itself.
- ▶ Resizing: neural network can only accept images of a fixed size, so resizing the image to match the network's requirements is important.
- ▶ Normalization: It is common to normalize the pixel values to be in the range of 0-1 or -1 to 1. This will help the model converge faster.
- ▶ Data Augmentation: To increase the amount of data and prevent overfitting, common data augmentation techniques such as flipping, rotation, and cropping can be applied to the images.

5: If you would prefer a visualization of a neural network check out this excellent video by Dennis Dmitriev <https://www.youtube.com/watch?v=3JQ3hYko51Y>.

In summary, preprocessing is an important step in training a neural network, as it ensures that the data is in a format that can be easily understood and processed by the network. The preprocessing steps for numbers, text, and images involve normalization, imputation, tokenization, stop-words removal, stemming/lemmatization, resizing and data augmentation.

## 2.4 Try That Again With 33,640,010 Parameters

In the first section of this chapter we introduced the idea that a simple neural network is mathematically equivalent to a regression. This is true and a useful way to think about neural networks and deep learning.<sup>6</sup>

In practice, neural networks frequently trained with millions of parameters. Here is an example of the trainable parameters of a simple image classifier. This example is a dense neural network to classify handwritten digits, and is not yet as sophisticated as the one used by Yann LeCun and company in the 1990s

In these next few examples, don't worry if you don't understand the layer types or know what batch normalization means. The point I would like to make is that neural networks are often created with millions of trainable parameters, once you agree with me regarding this point we will discuss the implications of this in more depth.

6:

$$y = Wx + b \quad (2.1)$$

There's your simple formula for a single-cell neural network and regression, in practice we'll change  $W$ ,  $x$  and  $b$  to matrices and introduce nonlinear activation function  $f$  so,

$$y = f(Wx + b) \quad (2.2)$$

if you please. If equations scare you, don't worry about it for now.

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 64)	50240
batch_normalization (BatchNo	(None, 64)	256
dense_1 (Dense)	(None, 64)	4160
batch_normalization_1 (Batch	(None, 64)	256
dense_2 (Dense)	(None, 64)	4160
batch_normalization_2 (Batch	(None, 64)	256
dropout (Dropout)	(None, 64)	0
dense_3 (Dense)	(None, 10)	650
Total params: 59,978		
Trainable params: 59,594		
Non-trainable params: 384		

59,594 parameters! That's not too bad for a simple neural network, but in practice the networks often get even bigger. Here is a *real* example that will make your brain a little hot.

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 28, 28, 64)	640
conv2d_1 (Conv2D)	(None, 28, 28, 64)	36928
max_pooling2d (MaxPooling2D) (None, 14, 14, 64)	0	
batch_normalization (BatchNo (None, 14, 14, 64)	256	
conv2d_2 (Conv2D)	(None, 14, 14, 128)	73856
conv2d_3 (Conv2D)	(None, 14, 14, 128)	147584
max_pooling2d_1 (MaxPooling2 (None, 7, 7, 128)	0	
batch_normalization_1 (Batch (None, 7, 7, 128)	512	
conv2d_4 (Conv2D)	(None, 7, 7, 256)	295168
conv2d_5 (Conv2D)	(None, 7, 7, 256)	590080
conv2d_6 (Conv2D)	(None, 7, 7, 256)	590080
max_pooling2d_2 (MaxPooling2 (None, 3, 3, 256)	0	
batch_normalization_2 (Batch (None, 3, 3, 256)	1024	

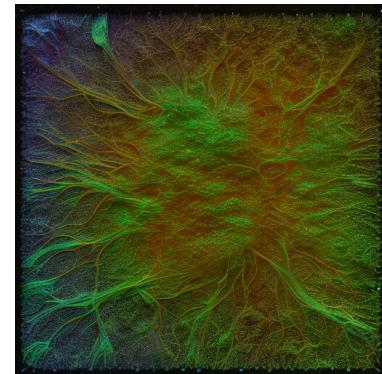


Figure 2.3: "mdjrny-v4 layers and layers of neurons looking like a gooey lasagna, but the data has green blood 8k" made with Mann-E

conv2d_7 (Conv2D)	(None, 3, 3, 512)	1180160
conv2d_8 (Conv2D)	(None, 3, 3, 512)	2359808
conv2d_9 (Conv2D)	(None, 3, 3, 512)	2359808
max_pooling2d_3 (MaxPooling2D)	(None, 2, 2, 512)	0
batch_normalization_3 (Batch Normalization)	(None, 2, 2, 512)	2048
conv2d_10 (Conv2D)	(None, 2, 2, 512)	2359808
conv2d_11 (Conv2D)	(None, 2, 2, 512)	2359808
conv2d_12 (Conv2D)	(None, 2, 2, 512)	2359808
max_pooling2d_4 (MaxPooling2D)	(None, 1, 1, 512)	0
batch_normalization_4 (Batch Normalization)	(None, 1, 1, 512)	2048
flatten (Flatten)	(None, 512)	0
dense (Dense)	(None, 4096)	2101248
dropout (Dropout)	(None, 4096)	0
dense_1 (Dense)	(None, 4096)	16781312
dropout_1 (Dropout)	(None, 4096)	0
dense_2 (Dense)	(None, 10)	40970
<hr/>		
Total params: 33,642,954		
Trainable params: 33,640,010		
Non-trainable params: 2,944		
<hr/>		

When a neural network has millions of trainable parameters and deep layers of neurons<sup>7</sup>, it can be difficult to explain to a human what each of those parameters represents or how they contribute to the overall function of the network. This is because the interactions between the different layers and neurons can be complex and non-linear, making it challenging to understand the specific role of each parameter. The parameters model complex interactions between the inputs, making it difficult to understand their specific function. Furthermore, in deep neural networks, the high number of layers can lead to high level of abstraction, meaning that the individual neurons and their weights have little interpretability.

When a neural network has millions of interacting parameters, it can lead to mathematical chaos<sup>8</sup>, which is a phenomenon where small changes in the initial conditions of the network can lead to vastly different outputs. This is because the interactions between the large number of parameters can create non-linear relationships that are sensitive to small changes. This can

7: In case you are wondering, the human brain has about 86 billion neurons that do the "thinking". Physical neurons are more complicated than the neurons used in deep learning and have supporting structures like astrocytes that do some computation as well. <https://en.wikipedia.org/wiki/Astrocyte>

8: Mathematical chaos is a real thing, it refers to the behavior of certain dynamic systems that are highly sensitive to initial conditions. This sensitivity leads to seemingly random and unpredictable behavior, even though the underlying equations governing the system are deterministic (meaning they are transparent and known to us). [https://en.wikipedia.org/wiki/Chaos\\_theory](https://en.wikipedia.org/wiki/Chaos_theory)

make it difficult to predict the behavior of the network, as small changes in the input or the parameters can lead to unexpected and seemingly random outputs.<sup>9</sup>

## 2.5 Multicollinearity and the End of Science

Data science is a horrible term because it implies that data scientists are scientists and that the work they do is scientific when in reality data scientists are not scientists and the work they do is not scientific. Data scientists use mathematics, statistics, and computer science to analyze data, but it is not scientific in the traditional sense. Data scientists do not use the scientific method and do not conduct experiments or develop theories. Data science is more akin to engineering or data analytics than actual scientific research.

The main goal of a scientist is to gain knowledge and understanding of the natural world through research, experimentation, and data analysis. Scientists make models of the world to test and explain. So-called data scientists make models too but a model with layers of interacting parameters, trained under chaotic conditions is almost impossible to explain. Multicollinearity<sup>10</sup> is just one issue, but deep learning techniques are practically incompatible with the scientific method. With deep learning methods it can be difficult to determine which variable had the most impact on the outcome of the model. Coefficients of the model to be unstable and unreliable, which can make it difficult to explain the model in a meaningful way.

Explaining the inner-workings of a neural network with millions of interacting parameters is difficult for many reasons. Most of the complexity is due to the number of parameters, which can make it difficult to understand how the model works and why it makes certain decisions. The interactions between the parameters are often non-linear and can be difficult to understand, as the model can be sensitive to small changes in the parameters, which is the definition of mathematical chaos. This can make it difficult to explain why the model is making certain decisions, as the interactions between the parameters are not always clear.

## 2.6 Let's Test Some Random Inputs! Feature Importance and Explainability

With a huge complex system of neurons that combine inputs in novel ways, it becomes very hard to understand which inputs are the most important for the system as a whole. To better understand deep learning models, data scientists use randomized or averaged inputs to a model to test the feature importance of a deep learning neural network. This type of testing is done to determine which features are most influential in the overall output of the network. This is done by randomly or averaging the input values for each feature and then running the model to see how the output is affected. For example, if a neural network is used to detect objects in an image, the data

9: Chaotic outputs in a large deterministic system are OK in many domains, and there are controls that can be put in place to keep AI "safe" but deep learning by itself will not produce those controls.

10: "Multicollinearity is when two or more variables in a study are closely related to each other. Think of it like dating: imagine you are trying to figure out what makes someone a good partner. You might look at things like how much they make, how attractive they are, and how kind they are. But, these things are all related to each other. For example, attractive people might make more money, and kind people might be more attractive. So, it's hard to know which one is really important for being a good partner, because they are all related. That's like multicollinearity in a study." By ChatGPT given the prompt "explain multicollinearity to a teenager using a dating example"

scientist may randomize the color of the objects to see how the model's accuracy is affected. If the accuracy drops significantly, they can infer that color is an important feature in the network.

Randomized or averaged inputs to a model can also be used to determine if a particular feature is necessary for the network to function correctly. For example, if the output of the network is not as accurate when a particular feature is randomized or averaged, then the data scientist can infer that this feature is important for the network's performance. By using this method, data scientists can gain insights into which features are important and which can be removed from the model to improve performance.

## 2.7 The Universal Machine Learning Workflow

The Universal Machine Learning Workflow[6] is an important chapter in a technical guide for data scientists by the author of the most popular machine learning framework in the world<sup>11</sup> that outlines the *Universal* workflow for machine learning projects. I think this chapter should be understood by everyone using, investing in and creating machine learning models.

Before Chollet gets into the details of model building, he chooses to begin with a note on ethics:

*"You may sometimes be offered ethically dubious projects, such as "building an AI that rates the trustworthiness of someone from a picture of their face." First of all, the validity of the project is in doubt: it isn't clear why trustworthiness would be reflected on someone's face. Second, such a task opens the door to all kinds of ethical problems. Collecting a dataset for this task would amount to recording the biases and prejudices of the people who label the pictures. The models you would train on such data would merely encode these same biases into a black-box algorithm that would give them a thin veneer of legitimacy. In a largely tech-illiterate society like ours, "the AI algorithm said this person cannot be trusted" strangely appears to carry more weight and objectivity than "John Smith said this person cannot be trusted," despite the former being a learned approximation of the latter. Your model would be laundering and operationalizing at scale the worst aspects of human judgement, with negative effects on the lives of real people."*

*Technology is never neutral. If your work has any impact on the world, this impact has a moral direction: technical choices are also ethical choices. Always be deliberate about the values you want your work to support." [6]*

Chollet uses the outline below to explain the *Universal* workflow. I'll summarize the workflow and my own notes for a nontechnical audience here:

1. Define the Task
  - a) Collect a Dataset <sup>12</sup>
  - b) Understand Your Data
  - c) Choose a Measure of Success <sup>13</sup>
2. Develop a Model

[6]: Chollet (2022), *Deep learning with python, second edition*

11: It's called TensorFlow and was one of the most popular and "Most Loved" Technologies of 2022 <https://survey.stackoverflow.co/2022/#most-popular-technologies-misc-tech>



Figure 2.4: "the face of a trustworthy person" made with Stable Diffusion 2.1. (Hey, it's a white lady!)

12: This is often the hardest part. Data is often labeled by hand or stolen from another domain. See [https://en.wikipedia.org/wiki/Transfer\\_learning](https://en.wikipedia.org/wiki/Transfer_learning)

13: Accuracy isn't everything! We might prefer a less accurate model that avoids false-positives in some scenarios.

- a) Prepare the Data <sup>14</sup>
  - b) Choose an Evaluation Protocol
  - c) Beat a Baseline <sup>15</sup>
  - d) Develop a model that overfits
  - e) Regularize and Tune Your Model
3. Deploy the Model
- a) Explain Your Work to Your Stakeholders and Set Expectations <sup>16</sup>
  - b) Ship an Inference Model
  - c) Monitor Your Model in the Wild
  - d) Maintain Your Model

Of all of the steps in the workflow, "Defining The Task" and "Explaining The Work To Shareholders and Setting Expectations" are where the most miscommunication occurs.

In defining the task, machine learning engineers are often given impossible problems to solve, and because they want to keep paying their mortgage they solve another related problem instead and allow stakeholders to jump to their own conclusions. By clearly understanding what a model is predicting and how the data is collected some of this misunderstanding can be avoided.

I will discuss Large Language Models later, but they all do the same thing. They predict the next words in a sentence, just like the keyboard on your iPhone. They come up with amazing text but by fundamentally understanding what they are predicting you gain insight into their limitations. They are also all trained on the text publicly available on the internet, this includes scientific sources, but also fan-fiction and anime, so the idea that a model trained on this data could be relied on to predict anything truthful is preposterous.

Many models work like magic and many users assume models are predicting the future, when they are really correlating based on their past data. Even simple models of creditworthiness might have a similar problem. While building a creditworthiness model for a bank, due to lack of data a machine learning engineer might (on purpose or inadvertently) create a model that predicts whether someone is a smoker instead of whether they are creditworthy. Because smoking and poverty are correlated, maybe the model "works", but it doesn't do what stakeholders think it does.

Setting expectations is another hellscape of misaligned incentives.<sup>17</sup> Good machine learning engineers and data scientists are supposed to educate and advise stakeholders about the limitations of their own models. Read Chollet's advice on this topic below:

*The expectations of non-specialists towards AI systems are often unrealistic. For example, they might expect that the system "understands" its task and is capable of exercising human-like common sense in the context of the task. To address this, you should consider showing some examples of the failure modes of your model (for instance, show what incorrectly classified samples look like, especially those for which the misclassification seems surprising).*

14: This is discussed in section 2.3 above, everything gets turned into numbers.

15: Does your model predict better than a totally random model and/or coin-flip?

16: This is almost never done. And as Chollet explains this is both difficult and requires mutual understanding and domain knowledge from the part of "the business" and the Machine Learning Engineer.

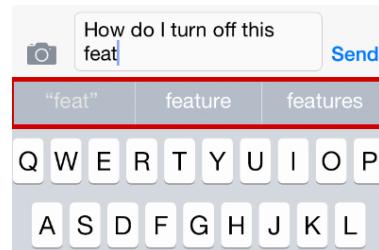


Figure 2.5: Fundamentally LLMs are using the same techniques as a predictive keyboard on your phone. This is why Yann LeCun says even though they are impressive some like ChatGPT are not particularly innovative. <https://www.youtube.com/watch?v=ULbpPHjiSBg>

17: The CEO of OpenAI wrote this and still raised ten billion dollars: *"ChatGPT is incredibly limited, but good enough at some things to create a misleading impression of greatness. it's a mistake to be relying on it for anything important right now. it's a preview of progress; we have lots of work to do on robustness and truthfulness."* <https://twitter.com/sama/status/1601731295792414720>

*They might also expect human-level performance, especially for processes that were previously handled by people. Most machine learning models, because they are (imperfectly) trained to approximate human-generated labels, do not nearly get there. You should clearly convey model performance expectations. Avoid using abstract statements like "The model has 98 percent accuracy" (which most people mentally round up to 100 percent), and prefer talking, for instance, about false negative rates and false positive rates. You could say, "With these settings, the fraud detection model would have a 5 percent false negative rate and a 2.5 percent false positive rate. Every day, an average of 200 valid transactions would be flagged as fraudulent and sent for manual review, and an average of 14 fraudulent transactions would be missed. An average of 266 fraudulent transactions would be correctly caught." Clearly relate the model's performance metrics to business goals.*

*You should also make sure to discuss with stakeholders the choice of key launch parameters- for instance, the probability threshold at which a transaction should be flagged (different thresholds will produce different false negative and false positive rates). Such decisions involve trade-offs that can only be handled with a deep understanding of the business context.[6]*

One does not need to be a psychologist to understand that these conversations almost never happen. The workflow of machine learning is universal, we are all doing fancy regressions and we all deal with the same wild expectations, shitty data and misalignment with business.

## 2.8 A Machine Learning Engineer is a Data Janitor

Since the machine learning workflow is *universal*, it is fairly straightforward to automate. There are a plethora of offerings from companies big and small who offer AutoML tools that anyone with basic Excel skills can use<sup>18</sup>. Since modern AI is "all a regression" these tools do the regression for you, users just need to feed them the data and the AutoML tool finds out the relationship. On the surface it seems like the job of the machine learning engineer has become obsolete before it even became a proper discipline.

Despite the existence of AutoML tools, machine learning engineers (and data scientists) do some actual work. A 2016<sup>19</sup> Crowd Flower survey of 16,000 data scientists showed that on average we spend our time doing the following:

- ▶ **60 percent of a machine learning engineer's time is spent cleaning and organizing data**, we are data janitors! If you had great data ready in an Excel file you would have been able to fit your own model without us. We need to organize the mess of data that exists in the world so it is in a nice format to be fit in a model. For example it is well known that almost all large language models use the same common crawl dataset (<https://commoncrawl.org/>), but how it is organized and weighted can change the quality of the outputs dramatically. The organization of the same dataset can lead to the same dataset either



Figure 2.6: "mdjrny-v4 a handsome businessperson explaining the business context to a scientist wearing a white coat over coffee 8k" made with Mann-E

18: See this awesome-AutoML github project for an up-to-date list of commercial AutoML tools <https://github.com/windmaple/awesome-AutoML#commercial-products>

19: I know this study is old, but I have found that the findings still hold when talking to data scientists and/or machine learning engineers that I train and employ.

producing a nice predictive keyboard or something approaching ChatGPT.

- ▶ **19 percent of a machine learning engineer's time is spent collecting data**, we often start making models without a huge dataset, so data needs to be collected in order to make the first model, and after that first model (to match people for our new dating app, let's say) is deployed we rely on users to provide us with data for future models. Collecting new data is a creative endeavour sometimes, and often involves human effort. Amazon's Mechanical Turk<sup>20</sup> service has been leveraged for this purpose for years, and OpenAI has successfully deployed outsourced talent to solve some of the trickiest problems facing large language models<sup>21</sup>.
- ▶ **9 percent of a machine learning engineer's time is spent fitting models**, we need to fit models, there is some art to choosing the right architecture and modeling techniques. An AutoML tool can do this reasonably well too, but considering some fancy training servers cost \$28 per hour or more<sup>22</sup>, sometimes it is useful to have an expert guess the best model architecture and save time in training, even if that expert is getting paid \$249,000 per year<sup>23</sup>.
- ▶ **4 percent of a machine learning engineer's time is spent refining algorithms** many machine learning engineers spend almost no time doing this. Other researchers (corporate-academic types who write scientific papers) spend a lot of time doing this, it averages out to four percent, but doesn't represent the average day of a machine learning engineer.

Overall, the main job of most machine learning engineers is to clean up and collect a ton of data, and then feed that data into the same machine learning model that everyone else is using, and then rinse and repeat.<sup>24</sup>

## 2.9 Key Takeaways

- ▶ **Deep learning models are fundamentally large unscientific regressions** they are trained to create a function that maps input data to output data.
- ▶ **Deep learning models are chaotic systems containing millions of interacting parameters** they are not designed to be explained or created in a way that their weights can be used for scientific analysis. They find reasonable answers and don't care how they get there. Multicollinearity (understanding the relationship of an input and output) and feature importance (understanding which inputs are most important) are only understandable with a high level of statistical error.
- ▶ **Small changes in inputs of a deep learning model may dramatically change the outputs** deep learning models are complex deterministic systems that can exhibit chaotic behavior. Their inner workings are functionally unknowable and practically impossible to test.
- ▶ **Machine learning engineers spend most of their time collecting and organizing data**, because deep learning models often share

20: "Getting the right training data was another challenge. ImageNet was a collection of one hundred thousand labeled images that required significant human effort to generate, mostly by grad students and Amazon Mechanical Turk workers."  
<https://arstechnica.com>

21: Read more about how OpenAI used Kenyan workers to help label toxic content <https://time.com/6247678/openai-chatgpt-kenya-workers/>.

22: See the latest pricing from AWS here <https://aws.amazon.com/sagemaker/pricing/>

23: See the latest machine learning engineer salaries at levels.fyi <https://www.levels.fyi/t/software-engineer/focus/ml-ai>

24: Here is a discussion on reddit asking "When was the last time you wrote a custom neural net?" and supports my understanding of the world [https://www.reddit.com/r/MachineLearning/comments/yto34q/d\\_when\\_was\\_the\\_last\\_time\\_you\\_wrote\\_a\\_custom/](https://www.reddit.com/r/MachineLearning/comments/yto34q/d_when_was_the_last_time_you_wrote_a_custom/)

common architecture, getting good data to train with is the best thing that an engineer can do to train good models. In practice only a handful of corporate-academic types are experimenting with new and exciting architectures.

- ▶ **Deep learning models can have impressive and useful outputs, but the creators of models should be encouraged to highlight their failures and limitations.** Machine learning engineers might be more keen to highlight failures and limitations if they are encouraged to do so by their users, managers and investors.

# Creativity and Decision Making with Deep Learning Models

3

## "AI Policy

*I expect you to use AI (ChatGPT and image generation tools, at a minimum), in this class. In fact, some assignments will require it. Learning to use AI is an emerging skill, and I provide tutorials in Canvas about how to use them. I am happy to meet and help with these tools during office hours or after class*

*Beware of the limits of ChatGPT:*

- *If you provide minimum effort prompts, you will get low quality results. You will need to refine your prompts to get good outcomes. This will take work.*
- *Don't trust anything it says. If it gives you a number or a fact, assume it is wrong unless you either know the answer or can check in with another source. You will be responsible for any errors or omissions provided by the tool. It works best for topics you understand.*
- *AI is a tool, but one that you need to acknowledge using. Please include a paragraph at the end of any assignment that uses AI explaining what you used the AI for and what prompts you used to get the results. Failure to do so is in violation of academic honesty policies*
- *Be thoughtful about when this tool is useful. Don't use it if it isn't appropriate for the case or circumstance."*

*Dr. Ethan Mollick, 2023 - Syllabus for class at The Wharton School at the University of Pennsylvania*

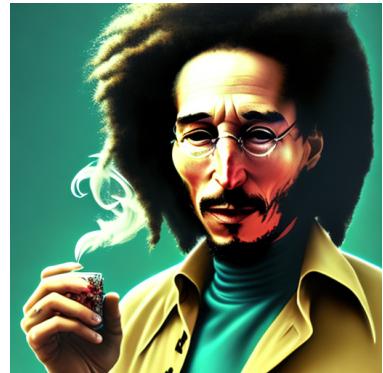
## 3.1 Theories of Creativity

Machine learning models use data (from the past) to discover rules and make classifications. Because of the way they are constructed these classifications, suggestions, artworks are by definition derivative or "having parts that originate from another source". I won't get into a philosophical discussion on what the nature of creativity is, but it's worth considering how using deep learning models biases us towards the past, but also could give us insights from other domains.

It can be argued that creativity is simply the combination of existing works. This is because many new ideas and innovations are often inspired by and built upon existing concepts. For example, a new form of music may be created by combining elements from different genres. Similarly, a new technology may be created by combining and improving upon existing technologies.

It can also be argued that creativity involves much more than just combining existing works. Creativity is not just about recombining existing ideas, but also about coming up with completely new and original concepts. This

3.1 Theories of Creativity . . . . .	19
3.2 Creative Uses of Power Tools . . . . .	20
3.3 Garbage In, Garbage Out . . . . .	21
3.4 Garbage In, New Perspective Out? . . . . .	21
3.5 Concept Drift and the End of Usefulness . . . . .	22
3.6 The Impossibility of Fairness . . . . .	23
3.7 Transfer Learning Everywhere . . . . .	23
3.8 Industrial-Scale Plagiarism	25
3.9 Working Together is Best . . . . .	25
3.10 Key Takeaways . . . . .	26



**Figure 3.1:** "mdjrny-v4 Steve Jobs Smoking weed with Bob Marley 8k" made with Mann-E. It's 100% derivative, but it's art too (I guess).

requires a unique perspective and a deep understanding of the subject matter, as well as the ability to think outside the box.

Deep learning models of speech, when heavily used, may slow down the evolution of language. AI art models may slow down "progress" in art, whatever that means. Deep learning models of disease trained on data from 1980 may be irrelevant to today's diseases. That said these same models may give us interesting insights in new domains, models trained on beautiful paintings may be put to use in a new domain (like designing beautiful interiors) and that model could give new insight to interior designers, models of the interaction of ants could be put to use in designing cities and so on and so forth. AI cuts many ways, it makes us faster but makes us more reliant on the past, models can be used across domains but should be used intentionally and transparently when possible. Each use opens up a new world of possibilities for users, and sometimes a new headache for intellectual property lawyers. We'll discuss all of these topics in this chapter.

## 3.2 Creative Uses of Power Tools

Power tools, even simple ones like a drill are complex machines that take human input and transform it using static rules. The pressure the user of a drill puts on the bit and the speed at which they pull the trigger all deterministically affect the output. Just because a tool is a deterministic machine doesn't mean it is unable to produce creative works.

What is happening as we use generative tools like ChatGPT or image-to-text models is that the "creative act" has been relocated. The creative act is now the prompt you give the tool, the user's input. And someone can still be good at using AI, just like someone can be good at using any piece of software.

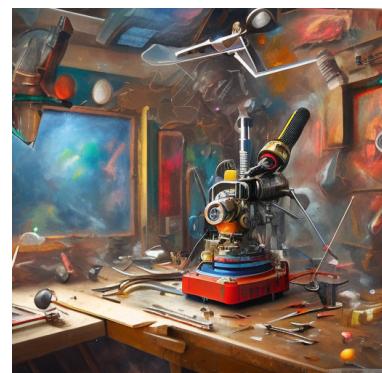
Modern AI is a complex system of algorithms, data, and analytics that can be used to solve complex problems. AI systems can learn from data, identify patterns, and make predictions about the future. AI systems are typically used to automate and assist human decision making. AI systems are programmed with specific objectives and goals, and the user input decides the output. For example, an AI system could be programmed to solve a mathematical problem and the user input would determine the parameters of the problem and the output would be the solution.

A power drill is also a complex deterministic system. The user input is limited to the type of drill bit, the speed of the drill, and the direction of the drill. The output is determined by these inputs, as the drill will only drill in the direction and speed determined by the user. The user also has to choose the correct drill bit to ensure the drill can do the job correctly and safely.

Both modern deep learning models and a power drill are complex but ultimately deterministic systems. The user input, however limited it may be, ultimately controls the output of the system. Both systems require a



**Figure 3.2:** "Batman dressed as a cowboy" made with Mann-E. It's 100% derivative as well, but better than any Google image searches I did. It looks amazing! I also might get sued if I use it even though the creators of Mann-E tell me I have the rights to it.



**Figure 3.3:** "mdjrnny-v4 a mikita drill being used in an artist's studio to make a colorful artwork 8k" made with Mann-E

user to understand how to use them and to make the correct input to get the desired output.

### 3.3 Garbage In, Garbage Out

The "Power Tool" of AI is a deterministic, static and unchanging system and the rules of that system are determined by the data that the model is shown. If a model is trained on shitty data, it will produce shitty results. End of book...

... maybe not. It's worth thinking about this for a moment. It is often said that modern AI can "generalize and make informed inferences given new data". If deep learning models are really just a big regression, these models will always come up with an answer, but if the world changes, these models will still be projecting complicated averages of their past data into the future.

So, let's separate AI's decision making into two extremes; *Creative Decision Making* and *Critical Decision Making*. The stakes are very low in a world of *Creative Decision Making* and who gives a shit if the AI is all a regression, and it just mushes together the limited data that it's seen. In a creative context you can also ask an AI interesting questions, so long as you don't solely rely on its output without checking the facts first<sup>1</sup>. Even if "Garbage In, Garbage Out" holds, garbage can still be helpful for a creative process.

Note that this book is called "Full-Self Driving, Skynet..", a self-driving car and a nuclear-bomb-equipped all-seeing AI are clearly not engaging in any *Creative Decision Making*.

1: See the syllabus note from Dr. Mollick at the beginning of this chapter.

### 3.4 Garbage In, New Perspective Out?

For creative tasks, it generally doesn't matter that a deep learning model is unscientific or trained on a lopsided dataset. An informed user of AI knows this and can account for that in their decision making, especially when engaging in creative decision making. The situation becomes problematic once we allow deep learning models to engage in critical decision making by themselves.

Generative models will come up with amazing (but derivative) works of art, But, they will never "change the game". A model of sculpture trained on past sculptures in 1917 will never come up with Marcel Duchamp's "Fountain". There is a term that machine learning engineers use for this, when the rules of the game are changed, it's called "Concept Drift".



**Figure 3.4:** Marcel Duchamp's "Fountain". A urinal that blew peoples minds <https://www.tate.org.uk/art/artworks/duchamp-fountain-t07573>.

## 3.5 Concept Drift and the End of Usefulness

Concept drift refers to the phenomenon where the distribution of data changes over time, causing the performance of deep learning models built on historical data to degrade. The model becomes less useful because it is trained to make predictions based on the relationship between the inputs and outputs in the data it was trained on, and if this relationship changes, the model may start making incorrect predictions. This is particularly problematic in real-world applications, where the data is constantly evolving and the relationship between inputs and outputs is subject to change. To mitigate the effects of concept drift, it is often necessary to continually retrain deep learning models on updated data.

The frequency with which a deep learning model needs to be retrained to mitigate the effects of concept drift depends on several factors, including the rate at which the data distribution changes, the size and complexity of the model, and the availability of computational resources.

For some applications with relatively stable data distributions, retraining the model once every few months or even once per year may be sufficient. However, in other applications where the data is changing rapidly, it may be necessary to retrain the model more frequently, such as once per week or even once per day.

Ultimately, the frequency of retraining will depend on the specific requirements of the application, and the trade-off between the cost of retraining and the potential cost of incorrect predictions. In general, it's recommended to monitor the performance of the model over time and to retrain it as needed to ensure that it remains accurate and relevant.

If a deployed deep learning model is not monitored, there are several risks that can emerge:

- ▶ Accuracy degradation: As the data distribution changes over time, the model may become less accurate, leading to incorrect predictions. This can result in financial losses, reduced customer satisfaction, or even harm to individuals.
- ▶ Bias amplification: deep learning models can be biased, and if this bias is not monitored and addressed, it can be amplified over time as the model continues to make incorrect predictions. This can result in discriminatory outcomes, such as unequal treatment of individuals based on protected characteristics such as race, gender, or age.
- ▶ Legal liability: In some cases, incorrect predictions made by deep learning models can result in legal liability, particularly if the model is being used to make decisions that have significant consequences, such as in the criminal justice system or in medical diagnosis.
- ▶ Reputational damage: If a deep learning model is making incorrect predictions, it can damage the reputation of the organization deploying the model, potentially leading to a loss of customers or investors.

It is important to monitor deep learning models once they are deployed, and to take action to address any issues that arise, such as retraining the



**Figure 3.5:** "Plastic surgery gone wrong" made with Stable Diffusion. Imagine a model that classifies images as "human face" or "not human face", and imagine that model was trained on images of human faces before 1900, maybe you would not be surprised if you gave it a picture of a human face that had a lot of cosmetic surgery done to it, and that model might say "this is not a human face", the idea of what a human face is has changed over time this is called "concept drift".

model or adjusting its parameters, in order to mitigate these risks and ensure that the model continues to perform well over time.

## 3.6 The Impossibility of Fairness

Achieving fairness in deep learning models can be challenging, and to some extent it may be impossible to completely eliminate all forms of bias. This is because deep learning models are trained on historical data, which may contain biases and disparities that are perpetuated in the model's predictions. One can attempt to "fix the training set" but in practice models will continue to bias their prediction to past data, or their creators careful curation of the past.<sup>[7]</sup>

However, it is possible to reduce the impact of bias in deep learning models through careful design and monitoring of the model's performance. This may include techniques such as fairness constraints, algorithmic transparency, and regular auditing of the model's predictions to identify and address any issues of bias.

It's important to note that fairness is a complex and multi-faceted concept, and different definitions of fairness may be appropriate for different applications. For example, some definitions of fairness may prioritize equal treatment of all individuals, while others may prioritize proportional representation or equal opportunities.

Ultimately, the extent to which fairness is achievable in deep learning models will depend on the specific requirements of the application and the level of effort that is put into designing and monitoring the model to ensure that it is making fair and unbiased predictions.

[7]: Christian (2020), *The Alignment Problem: Machine Learning and Human Values*

## 3.7 Transfer Learning Everywhere

Transfer learning is a machine learning technique that involves transferring knowledge from one model trained on a task to another model trained on a related task. The idea behind transfer learning is that a model that has been trained on one task can be fine-tuned for another task, reducing the amount of labeled data required to train the new model.

For example, imagine that you have a deep learning model that has been trained to recognize objects in natural images (cats, dogs, balls, cars). You can use the knowledge learned by this model as a starting point to train a new model that recognizes objects in medical images, such as X-rays or MRI scans. The new model can start with the weights of the pre-trained model and fine-tune them on the new task, using a much smaller amount of labeled data than would be required to train the model from scratch.

Transfer learning can be useful in many different applications, particularly when labeled data is scarce or expensive to obtain. By leveraging knowledge

from a pre-trained model, transfer learning can help to improve the performance of new models, reduce the amount of data required for training, and accelerate the development of new machine learning applications.

I asked ChatGPT to "tell me a story of transfer learning going horribly wrong in the context of online dating, with a model trained in the Middle East being deployed by an American dating app." and here's what it said:

*"Once upon a time, an American online dating company decided to use transfer learning to improve their matchmaking algorithms. They trained a state-of-the-art machine learning model on a large dataset of profiles and dating preferences from the Middle East. The model was able to suggest matches with remarkable accuracy and the company was confident that it would perform just as well in the American market.*

*However, things went terribly wrong when they deployed the model on their American dating app. The model was not properly trained to consider the cultural and social differences between the Middle East and America. As a result, it suggested matches that were inappropriate and offensive, causing widespread backlash and criticism from users.*

*For example, the model suggested matches between people of different religious backgrounds, which was not acceptable to many users. It also suggested matches between people with significant age gaps, which was not socially acceptable in the American market.*

*The company was forced to remove the model and start from scratch, causing significant harm to their reputation and financial losses. They learned the hard way that transfer learning can only be effective if the model is properly trained on relevant and culturally appropriate data.*

*This story serves as a cautionary tale about the importance of considering cultural differences and ethical implications when deploying machine learning models in different regions and markets."*

Not bad I think... picking good training data and doing good "data janitoring" is the main job of a machine learning engineer<sup>2</sup>.

If the past is not like the future, you are also doing transfer learning. Most models steal data from other sources so are doing transfer learning too. This is fine, but do we know that we are doing this? A model that is deployed in a domain experiencing concept drift and continues to make predictions without being retrained can be considered to be doing a harmful form of transfer learning. This is because the model has been trained on a different distribution of data (the past), and is being applied to a new domain with a different distribution (the present and future).

In traditional transfer learning, the goal is to transfer knowledge from one domain to another related domain, where the data distributions are similar enough to enable the model to make accurate predictions. However, in the case of concept drift, the data distributions are changing over time, and the model is becoming less accurate as a result.<sup>3</sup>

By continuing to make predictions without being retrained, the model is essentially "transferring" its knowledge from a historical data distribution

2: The better your machine learning engineer understands your business problem, the better job they'll do at creating the dataset.

3: Think about your life and if you "transferred" your model of thinking from when you were 12 years old, to "today" when you are 35 years old, that is a bad model to be operating on! That model needs to be updated. We'll explore many examples of where this is and is not a problem in the next chapter."

to a new, changed data distribution, which may not be a valid assumption. This can result in incorrect predictions and other negative outcomes, such as harm to individuals or organizations.

### 3.8 Industrial-Scale Plagiarism

Aside from regurgitating the past and predictions from other domains, deep learning models can also enable plagiarism on an industrial scale. Early text generation models could be made to write entire sections of *Harry Potter* when fed the opening lines of a chapter. Even as models grow large and more sophisticated, users, researchers and lawyers are still able to "extract the training data" from large models[8], causing headaches for their creators and adding to the work of intellectual property lawyers.

To avoid this creators of deep learning models have the following tools at their disposal:

- ▶ Best practices in machine learning: data preprocessing, augmentation, regularization, model architecture choices, early stopping and validation are all things that we are taught to do to prevent overfitting.
- ▶ Contractual agreements: Microsoft (owner of github and cocreator of the [Github Copilot](#) code generation model) has a special contract for anyone submitting code to github that essentially says that "we are allowed to use this code to train our models, and if our models regenerate your copyrighted code, you can go fuck yourself."
- ▶ Release the Kraken: StabilityAI's Stable Diffusion model was released as open source software, they still managed to get sued [9], but they basically said, "this model was trained on all copyrighted and copylefted images on the web, sometimes it'll generate stuff that violates copyright law, but we are in Germany and will give this thing away for free to the public, and see how Getty Images, photographers and graphic designers of the world handle it... OKbye!".

These techniques can help to reduce the risk of having deep learning models reproduce the training data and cause intellectual property disputes, but there is no way to completely eliminate these risks, they are simply side-effects of the method of machine learning that we are doing nowadays. If you are an IP lawyer and need an expert witness, I'm your man [brad@bradflaugh.com](mailto:brad@bradflaugh.com).

[8]: Carlini et al. (2023), *Extracting Training Data from Diffusion Models*



**Figure 3.6:** "man 30s talking to a doctor" Getty Images and Stable Diffusion comparison. Is Stable Diffusion creating something new here or spitting out its training data? Is training a model like this fair use? I don't know, have fun lawyers! Note rights to each Getty photo runs around \$450, and Stable Diffusion is free to use, and they give users the rights to the model output. The Kraken has been released!

### 3.9 Working Together is Best

It turns out sometimes the best combination is a reasonably smart (but not too cocky) human that gets to get recommendations from a few deep learning models.<sup>4</sup>

Garry Kasparov, the legendary chess player we heard about in chapter 1, witnessed a tournament where humans and AI could compete together. The tournament was designed to showcase the strengths of both humans

4: Funnily enough at Wharton I recently took [a class](#) where Dr. Cade Massey said basically "If you want to bring people along and get them to use your model, let them play with the output a bit, allowing for even the slightest bit (2%) adjustment in the outputs will get users to adopt your model much faster." Thanks Cade!

and AI, and how they could work together to achieve incredible results here is what Garry says about the tournament:

*"Once again, the chess world offers a useful test case for how this collaboration can play out. In 2005 the online chess playing site Playchess.com hosted what it called a "freestyle" chess tournament in which anyone could compete in teams with other players or computers. What made this competition interesting is that several groups of grandmasters working with computers also participated in this tournament. Predictably, most people expected that one of these grandmasters in combination with a supercomputer would dominate this competition ; but that's not what happened. The tournament was won by a pair of amateur American chess players using three computers. It was their ability to coordinate and coach effectively their computers that defeated the combination of a smart grandmaster and a PC with great computational power."* [10]

The tournament was a huge success, and showed that when humans and AI work together, they can achieve amazing things. The tournament participants learned that AI is not just a tool, but a valuable partner, and that by combining their strengths, they could achieve results that neither could have accomplished on their own. The tournament inspired many people to explore the potential of human-AI collaboration, and showed that by embracing technology, we can create a brighter future for all.

AI is an amazing partner, but we need to think for ourselves too. We can't blindly trust AI, but we can use it to inspire and challenge us. [11]

[10]: Kasparov (2021), *Ai should augment human intelligence, not replace it*: Harvard Business Review: March 18, 2021

[11]: Mansharamani (2020), *Think for yourself: Restoring common sense in an age of experts and artificial intelligence*

## 3.10 Key Takeaways

- **When you are using deep learning models, you are almost always engaging in some kind of transfer learning.** You are assuming the past will be like the future, or that training set reflects reality, which isn't always the case.<sup>5</sup>
- **Models can puke out seemingly creative things, but they are still a product of their training data.** deep learning models are deterministic systems, they slap together data based on rules they guessed from their training data. This process can still be very useful, but deep learning models are not sentient creative creatures from science fiction.
- **Even if you think you "own the rights" to the output of a model, models can sometimes puke out their training data** modelers try their best to prevent this but models can spit out their training data, which may cause headaches for everyone except the intellectual property lawyers.
- **Users should consider the quality and appropriateness of the data the model they are using is trained on before deciding to use a model.** Remember, machine learning engineers spend most of their time cleaning up and gathering data, if you ask questions about how and when the data was collected and cleaned up you can get ahead of potential problems of concept drift and blind transfer learning (engaging in transfer learning without your consent).

5: It's as simple as the common financial disclaimer *Past performance is no guarantee of future results*. Even though it is everywhere, the general public and even big swingers in finance don't seem to get this one. In AI there is an additional disclaimer, *Training data may not reflect the the "real world"*.

- **Work with the technology** don't make your model do everything, sometimes it'll come up with stupid answers. Also don't assume you are always right, try and remember the story of the freestyle chess tournament and the grandmasters who got beat by the amateurs getting reasonable advice from a few chess engines. Working with your model, as if it were a coworker that you don't fully trust, but still think is smart is probably the best way to use deep learning models.

# 4

## Model Cards and Case Studies

I like huggingface's model cards this chapter is inspired by their outline, but adds some of my commentary <https://huggingface.co/blog/model-cards>

### 4.1 Fluid Dynamics

- ▶ **Description** This AI model is designed to count sperm in a video of semen under a microscope.
- ▶ **Training Data** The model was trained on a large dataset of semen videos taken under a microscope. The training data was annotated with the number of sperm in each video, allowing the model to learn the patterns and features of sperm in different types of semen samples.
- ▶ **Evaluation Metrics** The model was evaluated using precision, recall, and F1 score, which are commonly used metrics for object counting tasks. The model achieved high scores on all metrics, indicating that it is effective at accurately counting sperm in the semen videos.
- ▶ **Use Cases** This model can be used in clinical or research settings to quickly and accurately count sperm in semen samples. This can be useful for diagnosing and monitoring male infertility, as well as for understanding the impact of various factors on sperm count.
- ▶ **Limits and Risks** Concept Drift and Transfer Learning! The model is only trained to count sperm in semen videos taken under a microscope and may not perform well on videos taken under different conditions or with different types of microscopes. It is important to ensure that the semen samples are of high quality and that the video recording conditions are consistent in order to obtain accurate results.
- ▶ **Disclaimer** This model is not intended to replace human expert judgment and should be used as a tool to support decision making. The results generated by this model are not a substitute for professional medical advice, diagnosis, or treatment.
- ▶ **The Myth** Talk about how folks infer causality, even though we are not doing science here.

4.1	Fluid Dynamics . . . . .	28
4.2	Shakespearean Text Generation . . . . .	29
4.3	Lithium Mining Classifier . . . . .	29
4.4	Chess Playng Model . . . . .	29
4.5	Go Playing Model . . . . .	29
4.6	Large Stock Order Detector . . . . .	29
4.7	Stock Trading Bot . . . . .	29
4.8	Sperm Counter . . . . .	29
4.9	Handwriting Recognizer . . . . .	29
4.10	Liver Drug Discovery . . . . .	29
4.11	Online Dating Matcher . . . . .	29
4.12	Online Ad Server Classifier . . . . .	29
4.13	Tennis Playing Bot . . . . .	29
4.14	Hate Speech Classifier . . . . .	30
4.15	Contract Review Bot . . . . .	30
4.16	Facial Recognition . . . . .	30
4.17	Smartwatch Danger Classifier . . . . .	30
4.18	CCTV Threat Classifier . . . . .	30
4.19	War Robot . . . . .	30
4.20	Harvard Acceptance Bot . . . . .	30
4.21	Simple Credit Score . . . . .	30
4.22	Social Credit Score . . . . .	30
4.23	Artificial General Intelligence Chatbot . . . . .	30

## 4.2 Shakespearean Text Generation

### 4.3 Lithium Mining Classifier

### 4.4 Chess Playng Model

### 4.5 Go Playing Model

<https://arstechnica.com/information-technology/2022/11/new-go-playing-trick-defeats-world-class-go-ai-but>

### 4.6 Large Stock Order Detector

### 4.7 Stock Trading Bot

### 4.8 Sperm Counter

### 4.9 Handwriting Recognizer

### 4.10 Liver Drug Discovery

<https://www.psychologytoday.com/us/blog/the-future-brain/202301/ai-finds-drug-candidate-for-liver-cancer-in-30-days>

### 4.11 Online Dating Matcher

### 4.12 Online Ad Server Classifier

### 4.13 Tennis Playing Bot

[https://en.wikipedia.org/wiki/Moravec%27s\\_paradox](https://en.wikipedia.org/wiki/Moravec%27s_paradox)

#### **4.14 Hate Speech Classifier**

#### **4.15 Contract Review Bot**

#### **4.16 Facial Recognition**

#### **4.17 Smartwatch Danger Classifier**

<https://www.nytimes.com/2023/02/03/health/apple-watch-911-emergency-call.html>

#### **4.18 CCTV Threat Classifier**

#### **4.19 War Robot**

<https://www.extremetech.com/extreme/342413-us-marines-defeat-darpa-robot-by-hiding-under-a-cardboard-box>

#### **4.20 Harvard Acceptance Bot**

#### **4.21 Simple Credit Score**

#### **4.22 Social Credit Score**

#### **4.23 Artificial General Intelligence Chatbot**

<https://openai.com/blog/forecasting-misuse/>

# 5

## Self-Driving With Statistics

- ▶ *First Law: A robot may not injure a human being or, through inaction, allow a human being to come to harm.*
- ▶ *Second Law: A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.*
- ▶ *Third Law: A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.*

*Asimov's Three Laws of Robotics*

5.1 Self-Driving Horses . . . . .	31
5.2 Semiautonomy is Stupid . . . . .	31
5.3 Trolley Problems . . . . .	31
5.4 Concept Drift (Reprise) . . . . .	31
5.5 Multicollinearity (Reprise) . . . . .	31
5.6 Explaining The Unexplainable in Court . . . . .	31
5.7 A Train is a Self-Driving Car, Right? . . . . .	32

### 5.1 Self-Driving Horses

### 5.2 Semiautonomy is Stupid

Talk about the levels 0-4, deep learning is a subsystem of these levels.

Semiautonomy is dumb because it feels autonomous, but the user is expected to take over at any time. This is the worst of both worlds. [12]

and this article [Tesla Crash Illustrates Problem With Semi-Automated Driving](#) and comments

### 5.3 Trolley Problems

### 5.4 Concept Drift (Reprise)

### 5.5 Multicollinearity (Reprise)

1

1: You can't get your self-driving car dirty either, it'll mess up those sensors.

### 5.6 Explaining The Unexplainable in Court

*"Limitations and Pitfalls of Explainable and Interpretable Methods*

*Before diving into the exact methods for interpreting and explaining models, let's take a look at some of the pitfalls of these methods.*

*First off, if you need to make high-stakes decisions, make sure to use inherently interpretable models<sup>2</sup>. These are models such as decision trees that are more readily converted to output explanations.*

2: This sentence is very important, they are saying DO NOT USE DEEP LEARNING TECHNIQUES FOR CRITICAL DECISION MAKING. This is wild, we are already using these models everywhere, it's like asking people not to use computers.... TODO EXPAND ON THIS.

*Before choosing a method, you need to be absolutely clear about what you want out of it. Are you trying to understand the nature of the data procurement process? How a decision was made? How the model works on a fundamental level? Some tools might be appropriate for some of these goals but not others.*

*If your goal is to make sense of the data generation process, this is only possible if you know that your model already generalizes well to unseen data.*

*Decision interpretability can be misleading. It highlights things like correlations, but doesn't go into the level of causal detail that causal inference does. Remember that correlation does not (always) imply causation.*

**WARNING** Spurious correlations can result in inaccurate interpretations even with advanced interpretability methods like saliency methods and attention-based methods.

*Tools such as feature importance usually estimate mean values, but one should beware the error bars on those means and take stock of the confidence intervals.*

*A lot of machine learning involves working with extremely high-dimensional spaces. There's no way around it: high-dimensional data and feature spaces are hard to make sense of without grouping the data or features together first.*

*Even if you do find important features in those matrices, remember that this does not imply causality (we've said this before and we'll say it again)."* [13]

## 5.7 A Train is a Self-Driving Car, Right?

Discuss how the problem space has changed in warehousing, we don't actually have self-driving forklifts, we have moving shelves.

Maybe talk about elon musk and manufacturing.

# 6

## Unplugging Skynet

*"The second requirement of goal-misalignment risk is that an intelligent machine can commandeer the Earth's resources to pursue its goals, or in other ways prevent us from stopping it... We have similar concerns with humans. This is why no single person or entity can control the entire internet and why we require multiple people to launch a nuclear missile. Intelligent machines will not develop misaligned goals unless we go to great lengths to endow them with that ability. Even if they did, no machine can commandeer the world's resources unless we let it. We don't let a single human, or even a small number of humans, control the world's resources. We need to be similarly careful with machines." Jeff Hawkins, 2022 [4]*

6.1 AI and Human Safety . . . . .	33
6.2 Useful Incompatibility . . . . .	33
6.3 Checks and Balances . . . . .	33

### 6.1 AI and Human Safety

### 6.2 Useful Incompatibility

It's a feature not a bug that no single computer can control all others, AKA Security by Obscurity

### 6.3 Checks and Balances

<https://www.acm.org/media-center/2023/january/techbrief-safer-algorithmic-systems>

# 7

## Revolutionary for Whom?

*"The inhabitant of London could order by telephone, sipping his morning tea in bed, the various products of the whole earth – he could at the same time and by the same means adventure his wealth in the natural resources and new enterprise of any quarter of the world – he could secure forthwith, if he wished, cheap and comfortable means of transit to any country or climate without passport or other formality."* John Maynard Keynes, 1920 [14]

### 7.1 The Battle of the Assistants

Butler vs Indian Virtual Assistant vs Siri

7.1 The Battle of the Assistants . . . . .	34
7.2 Employees That Are Better Than You . . . . .	34
7.3 Slow on the Uptake . . . . .	34
7.4 Free-Rider Problems . . . . .	34
7.5 When You Can't Tell The Difference . . . . .	34
7.6 Sucker Traps . . . . .	35
7.7 Dead Inside . . . . .	35
7.8 This Book is a Case Study . . . . .	35

### 7.2 Employees That Are Better Than You

Respond directly to Jon Krohn's TED talk about monkeys being dumber than us... what about construction equipment that's stronger than us physically, or racism/eugenics people that are dumber than us [15]

[15]: (2022), *Jon Krohn*

### 7.3 Slow on the Uptake

<https://www.economist.com/finance-and-economics/2023/02/02/the-ai-boom-lessons-from-history>  
talk about historical perspective of technological change. Counter the narrative of "fastest tech to 1 million users" for ChatGPT.

### 7.4 Free-Rider Problems

I can steal your AI quite easily from outputs

### 7.5 When You Can't Tell The Difference

Talk about Taleb's aphorisms "Another definition of modernity: conversations can be more and more completely reconstructed with clips from other conversations taking place at the same time on the planet.", "You are alive in inverse proportion to the density of cliches in your writing."

## 7.6 Sucker Traps

*"(Traditional) search engines are databases, organized collections of data that can be stored, updated, and retrieved at will. (Traditional) search engines are indexes. a form of database, that connect things like keywords to URLs; they can be swiftly updated, incrementally, bit by bit (as when you update a phone number in the database that holds your contacts).*

*Large language models do something very different: they are not databases; they are text predictors, turbocharged versions of autocomplete. Fundamentally, what they learn are relationships between bits of text, like words, phrases, even whole sentences. And they use those relationships to predict other bits of text. And then they do something almost magical: they paraphrase those bits of texts, almost like a thesaurus but much much better. But as they do so, as they glom stuff together, something often gets lost in translation: which bits of text do and do not truly belong together." Gary Marcus, 2023 [16]*

"The sucker's trap is when you focus on what you know and what others don't know, rather than the reverse."

## 7.7 Dead Inside

"If you know, in the morning, what your day looks like with any precision you are a little bit dead - the more precision the more dead you are."

## 7.8 This Book is a Case Study

Talk about "organic content" market.

## **Errors and Omissions**

**8**

# Bibliography

Here are the references in citation order.

- [1] John von Neumann and Ray Kurzweil. *The Computer and the Brain (The Silliman Memorial Lectures Series)*. New Haven, CT, USA: Yale University Press, Aug. 2012 (cited on page ii).
- [2] Matt Welsh. *The end of programming*. Jan. 2023. URL: <https://cacm.acm.org/magazines/2023/1/267976-the-end-of-programming/fulltext> (cited on page 1).
- [3] Elaine Rich, Kevin Knight, and Shivashankar B. Nair. *Artificial Intelligence*. Tata McGraw-Hill, 2009 (cited on page 1).
- [4] Jeff Hawkins. *A thousand brains: A new theory of intelligence*. Basic Books, 2022 (cited on pages 7, 33).
- [5] Caglar Aytekin. 'Neural Networks are Decision Trees'. In: (2022). doi: [10.48550/ARXIV.2210.05189](https://doi.org/10.48550/ARXIV.2210.05189) (cited on page 8).
- [6] Francois Chollet. *Deep learning with python, second edition*. Manning Publications, 2022 (cited on pages 14, 16).
- [7] Brian Christian. *The Alignment Problem: Machine Learning and Human Values*. New York, NY, USA: W. W. Norton & Company, Oct. 2020 (cited on page 23).
- [8] Nicholas Carlini et al. *Extracting Training Data from Diffusion Models*. 2023. doi: [10.48550/ARXIV.2301.13188](https://doi.org/10.48550/ARXIV.2301.13188). URL: <https://arxiv.org/abs/2301.13188> (cited on page 25).
- [9] Kyle Barr. *Getty Images Claims Stable Diffusion's Creator 'Copied' 12 Million Copyrighted Images*. <https://news.yahoo.com/getty-images-claims-stable-diffusions-220000016.html>. (Accessed on 02/08/2023). Feb. 2023 (cited on page 25).
- [10] Garry Kasparov. *Ai should augment human intelligence, not replace it: Harvard Business Review: March 18, 2021*. Mar. 2021. URL: <https://www.kasparov.com/ai-should-augment-human-intelligence-not-replace-it-harvard-business-review-march-18-2021/> (cited on page 26).
- [11] Vikram Mansaramani. *Think for yourself: Restoring common sense in an age of experts and artificial intelligence*. Harvard Business Review Press, 2020 (cited on page 26).
- [12] Jason Torchinsky and Beau Boeckmann. *Robot, take the wheel: The road to autonomous cars and the lost art of driving*. Apollo Publishers, 2019 (cited on page 31).
- [13] Matthew McAteer Yada Pruksachatkun and Subhabrata (Subho) Majumda. *Practicing Trustworthy Machine Learning*. <https://learning.oreilly.com/library/view/practicing-trustworthy-machine/9781098120269/ch03.html>. (Accessed on 02/08/2023). Feb. 2023 (cited on page 32).
- [14] John Maynard Keynes, Elizabeth Johnson, and Donald Moggridge. *The Collected Writings of John Maynard Keynes (Volume 5)*. Cambridge, England, UK: Cambridge University Press, Dec. 2012 (cited on page 34).
- [15] Jon Krohn. [Online; accessed 18. Oct. 2022]. Oct. 2022. URL: <https://www.jonkrohn.com/posts/2022/10/7/tedx-talk-how-neuroscience-inspires-ai-breakthroughs-that-will-change-the-world> (cited on page 34).
- [16] Gary Marcus. *Is chatgpt really a 'code red' for google search?* Jan. 2023. URL: <https://cacm.acm.org/blogs/blog-cacm/268376-is-chatgpt-really-a-code-red-for-google-search/fulltext> (cited on page 35).