

Assignment 1:

1.)The Aha Moment:

While reviewing the study materials, I had two key “aha” moments:

i.) Biases formed by machine learning:

In the pre-course Microsoft Form, I noted that machine learning can introduce **bias** and **overfitting** when analyzing datasets like crime distribution in towns. In **Pedro Domingos’** *“A Few Useful Things to Know About Machine Learning,”* bias is defined as **“a learner’s tendency to consistently learn the same wrong thing.”**

To address this, Domingos suggests **cross-validation**, writing:

“The most common way to estimate generalization performance is cross-validation. In k-fold cross-validation, the data is split into k subsets, and the learner is trained and tested k times...”

This technique helps models generalize better and avoid misleading performance from overfitting. A simple analogy: imagine a student preparing for an exam by dividing a textbook into five parts. Each time, they study four parts and test themselves on the fifth—rotating through all sections. This ensures genuine understanding, not memorization, just as cross-validation helps models perform well on unseen data.

ii.) Semantic Web versus Semantic Interpretation:

In *“The Unreasonable Effectiveness of Data,”* **Halevy, Norvig, and Pereira** argue that **“simple models and a lot of data trump more elaborate models based on less data.”** They demonstrate that

massive datasets allow machines to perform tasks like translation or intent detection through **statistical pattern recognition**, without needing hand-coded semantic rules. They note:

“Because of a huge shared cognitive and cultural context, linguistic expression can be highly ambiguous and still often be understood correctly.”

This insight shows how **quantity and variety of data** empower models to **memorize and generalize patterns**, making large datasets more valuable than complex model architecture.

Let me know if you'd like it formatted for a specific platform (e.g., Word doc, Google Doc, printed page).

2.) Data is King (or is it?):

i.) Applications of data analysis in the real world:

In this segment I would like to mention two examples of data analysis being used in the real world, one of which is based on my interests whereas the other was one that was talked about and explained during the lecture:

a.) Character selection rates in online games: One clear application of **data analysis** is in **online gaming**, where developers track **character selection rates** to maintain gameplay balance. Using tools like **2D arrays**, they map player choices to experience levels, helping adjust characters that are over- or underused. For example, experienced players might select more complex

characters, while newcomers prefer simpler ones. These insights lead to **balance patches** that not only adjust gameplay but also enhance the **player experience**. As **Helldivers 2** director **Pilestedt** stated, “**You may have a balanced game but is it fun, probably not...**” — highlighting that **enjoyment**, not just fairness, is the goal. **League of Legends**, a prominent example, has implemented over **160 such patches**, relying heavily on real-time data.

b.) Genomes and phenotypes: genomics and phenomics, where scientists study the relationship between an organism’s **genotype** and **phenotype**. Using **high-throughput sequencing** and **machine learning**, researchers link genetic variations to traits like disease risk or drought resistance. This helps in **personalized medicine**, **crop breeding**, and **understanding evolution**. By analyzing massive datasets, researchers can predict which gene combinations result in beneficial traits — a form of biological balance patching.

ii.) Availability and use of large amounts of data types:

As discussed in **Pedro Domingos’** work, “**More data usually beats a cleverer algorithm...**” Similarly, in “**The Unreasonable Effectiveness of Data**,” the authors argue that **data volume and diversity** often outweigh model complexity. Even “**messy**” or **incomplete data** can reveal meaningful patterns when large enough. Text, images, behavioral logs—this **variety** boosts machine learning capabilities across disciplines. In essence, having **large, diverse datasets** is often more powerful than building intricate models alone.

3.) Humanity in the loop:

i.)Overfitting and Its Causes

Since I've already explained overfitting before, **I'll now dive deeper into its root causes**, which are **bias** and **variance**. In simple terms, overfitting happens when models are **too flexible** and try to fit **every detail** in the training data, including random noise. As Pedro Domingos states, the key challenge in machine learning is **generalization**—how well a model performs on **unseen data**, not just the training set. He notes that **“the more complex the hypothesis class, the more prone the learner is to overfitting.”**

- **Bias** is the error introduced by using a **simplified model** to approximate a complex real-world problem. A model with high bias makes **strong assumptions**, often resulting in **underfitting**, where it misses important patterns (e.g., using a linear model for nonlinear data).
- **Variance** refers to the model's sensitivity to **minor fluctuations** in the training data. High variance causes the model to **memorize noise**, leading to **overfitting**. This is common in models like deep neural networks or high-degree polynomials, which fit the training data perfectly but **fail to generalize**.

ii.)Importance of Humanity in the Machine-Human Relationship

While techniques like **cross-validation** help combat overfitting, **humans play an essential role**. In *The Unreasonable Effectiveness of Data*, Halevy, Norvig, and Pereira argue that **more data** reduces overfitting by helping the model focus on **general patterns** rather than noise. **Humans are key providers of this data**—through labeling, collection, and curation. Their involvement ensures that models are trained on **relevant, diverse, and high-quality data**. As

machine learning continues to evolve, **human collaboration remains vital** to building systems that are robust, fair, and grounded in real-world understanding.

4.) Understanding LLMs

From what I understand about **AI models**, especially **Large Language Models (LLMs)** like ChatGPT, they rely heavily on being trained with **enormous amounts of data**. These models are built using **complex learning algorithms**, often based on deep learning architectures like **transformers**, which allow them to process, understand, and generate language that resembles **human communication**.

At the core of LLMs is a process called **training**, where the model is exposed to **billions of words** from books, websites, articles, conversations, and other publicly available text sources. The goal is for the model to **learn patterns** in how humans use language—grammar, meaning, tone, context, and even subtle things like irony or metaphor. This learning isn't based on understanding in the human sense but on recognizing **statistical relationships** between words and phrases. The model essentially **predicts the next word** in a sentence based on what came before, using probabilities learned from its training data.

The **learning algorithm** is what allows the model to update its internal parameters during training to improve its predictions. Over time, and with enough data, the model becomes remarkably good at generating text that seems coherent, informative, and contextually appropriate.

Ultimately, ChatGPT and other LLMs reflect the language they are trained on. Their ability to produce **human-like responses** comes not from true understanding, but from the power of **pattern recognition** on a massive scale.