# Week-01 HW- Horizons 25 (Data Science)

**1. The "Aha!" Moment: After going through the readings and slides, what was one or two of the most surprising or "aha!" ideas you encountered about data science, machine learning, or the nature of data itself? Explain what the idea was and why it stood out to you or changed your perspective**

One of the biggest "Aha!" moments I had from reading was the ubiquitous and subtle nature of bias in machine learning models and data. Prior to reading books such as Brynjolfsson and Mitchell's "What Can Machine Learning Do?" and Pedro Domingos' "A Few Useful Things to Know About Machine Learning," I knew what bias was as a theoretical concept, but I hadn't realized how far it can influence results in data systems.

The readings illustrated that machine learning models are only as accurate and fair as the data used to train them. For instance, Brynjolfsson and Mitchell point out that if historical data captures societal biases—like discrimination in employment or lending—then whatever model is constructed based on such data will also reproduce the same patterns. Domingos also points out that bias can enter not only from the data itself, but from feature selection, how the problem is defined, or even by subliminal decisions in data preprocessing.

This epiphany struck me because it shifted machine learning from a solely technical field to one that is profoundly ethical and human. It caused me to appreciate the obligation that data scientists bear—not merely to create accurate models, but to analyze data sources critically, challenge assumptions, and take active steps to counteract bias. The understanding of the centrality of bias has shifted my perspective on the whole field, keeping in mind that technology is never objective and careful, ethical management is required to create systems that benefit all equally.

**2. Data is King (or is it?): The articles "The Unreasonable Effectiveness of Data" and "The Rise of Big Data" highlight how having lots of data can be incredibly powerful, sometimes even more so than having a super-complex algorithm. • Choose one real-world application of data science or ML (either from the read- ings/slides or one you can think of). • How does the availability and use of large amounts of data (think quantity, variety, or even "messiness" as discussed in the articles) make this application possible or effective?**

An interesting real-world example of the role of data volume and diversity in driving machine learning innovation is language translation, especially in applications such as Google Translate. The first generation of translation systems used complicated, rule-based methods designed by linguists, but these fell short in dealing with the subtlety and diversity of real-world language. The turning point came with Google starting to train neural machine translation models on enormous multilingual datasets sourced from websites, documents, and user submissions. As discussed in "The Unreasonable Effectiveness of Data," this change demonstrated that raw data volume—hundreds of millions of sentence pairs—could allow even simple algorithms to learn subtle patterns, context, and meaning that hand-crafted systems couldn't. Google's neural models now employ sophisticated architectures, including sequence-to-sequence networks, and can even translate "zero-shot" between language pairs not seen during training, owing to the variety of their data. This illustration shows that, for most purposes, access to large, diverse datasets can be more revolutionary than algorithmic sophistication, empowering the development of highly capable AI systems. Language translation therefore vividly shows how data scale, and not simply brilliant engineering, has propelled progress in machine learning.

**3) Machine learning is advancing rapidly, but several read-ings (especially Domingos' "A Few Useful Things..." and Brynjolfsson & Mitchell's "What can ML do?") point out its limitations, challenges (like overfitting or the need for feature engineering), and areas where humans still excel.**
**• Discuss one challenge or limitation of current ML that you found interesting.**
**• Why do you think it's important for people learning data science to be aware**
**of this, and what role do you see humans playing in addressing or working alongside this limitation in the future?**

One of the most powerful challenges in machine learning, emphasized in Brynjolfsson and Mitchell's "What Can Machine Learning Do? ", is causal inference's difficulty. Machine learning is very good at identifying patterns and correlations within data, but it tends to struggle to establish causation—knowing whether one factor truly causes another.For instance, a model may notice that individuals who purchase sunscreen tend to purchase sunglasses as well, but it

cannot determine whether the purchase of sunscreen leads to sunglasses purchases or if both are merely an association with sunny weather.This constraint is important for data science students to grasp since many actual choices rely on causal inference instead of simple correlation. In healthcare, economics, and public policy, interventions that are ineffective or even damaging can result from acting on correlations without finding causality.It is humans who must help solve this problem by structuring experiments like randomized controlled trials to evaluate causal hypotheses, using domain knowledge to subject model output to critical interpretation, and employing tools like causal graphs or instrumental variables to steer machine learning toward more causally informed conclusions. Even as more recent research strives to bring causal inference frameworks and machine learning together, human judgment remains at the heart of question construction, verification of results, and prevention of misuse.Identifying the constraints of correlation and the value of causality sets data science students up to create more stable models and make more educated choices, highlighting that human understanding and experience will remain essential with evolving machine learning technologies.