

Prompt written by: Cody York

LLM: OpenAI ChatGPT(May 2023)

Version: GPT-4

Date: 05/18/2023

LLM Prompt: Write a one-page paper on how companies use data science and machine learning in production to perform Data Engineering. Data Engineering is referring to managing large changing datasets. Include citations for 3 sources used in the paper.

****Title:** Data Engineering in the Age of Data Science and Machine Learning: An Insight into Modern Companies**

Data Engineering has emerged as a critical domain in the modern technological landscape. Big Data and new data analytics enable new approaches to data generation and analyses to be implemented that make it possible to ask and answer questions in new ways. (Zikopoulos & Eaton, 2011).

Companies are increasingly using Data Science and Machine Learning (ML) to optimize data engineering processes and to derive actionable insights from the voluminous data that they generate daily. Data Science provides a systematic framework to extract, process, and visualize data, enabling companies to recognize patterns and trends more efficiently (Provost & Fawcett, 2013). On the other hand, ML automates the predictive analysis of data, freeing up data engineers to focus on more complex tasks, such as architecting scalable data pipelines.

To manage large, changing datasets in production, companies are integrating ML algorithms with traditional data engineering techniques. The outcome is a dynamic data processing pipeline that can adapt to the evolving data landscape. For example, companies leverage ML to automate the identification and correction of anomalies in data streams, thereby enhancing data quality (Kandel et al., 2012). This ML-assisted anomaly detection is particularly useful in managing streaming data, where the volume, velocity, and variety of data pose significant challenges.

Companies are also utilizing ML for data imputation, the process of filling missing values in a dataset. Traditional methods, which relied on statistical techniques, were often insufficient when faced with complex datasets with high dimensionality. ML algorithms, in contrast, can handle multiple variables simultaneously and more accurately predict missing values, thereby improving the overall data quality.

In the realm of data storage and retrieval, ML is enabling companies to structure and classify unstructured data, which comprises a significant portion of Big Data. Techniques such as Natural Language Processing (NLP) and image recognition algorithms are used to extract meaningful features from unstructured data, which can then be organized and stored in a structured format, making it accessible for future analysis.

In conclusion, Data Science and Machine Learning are transforming the face of Data Engineering in production environments. They allow companies to manage large, changing datasets more efficiently, accurately, and scalably, leading to improved business intelligence and decision-making.

****References****

1. Zikopoulos, P., & Eaton, C. (2011). Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data. McGraw-Hill Osborne Media.
2. Provost, F., & Fawcett, T. (2013). Data Science for Business: What You Need to Know about Data Mining and Data-analytic Thinking. O'Reilly Media.
3. Kandel, S., Paepcke, A., Hellerstein, J., & Heer, J. (2012). Enterprise Data Analysis and Visualization: An Interview Study. IEEE Transactions on Visualization and Computer Graphics, 18(12), 2917–2926.