

## Data Drift: The Challenge of Changing Data Over Time

As the amount of data being collected continues to grow, so too does the challenge of dealing with data drift. Data drift refers to the phenomenon of changes in data over time, which can negatively impact the accuracy and reliability of machine learning models. In this article, we will explore what data drift is, its advantages, techniques to measure data drift, examples where data drift issues need to be addressed, and how to use KL divergence to calculate data drift.

**Data drift** is a common problem in machine learning where the statistical properties of data used to train a model change over time. These changes can happen due to a wide range of factors, such as seasonality, changes in user behavior, new data sources, and many more. These changes can negatively impact the accuracy and reliability of a model. Therefore, detecting and addressing data drift is crucial to maintaining the performance of machine learning models.

There are **several advantages** to dealing with data drift. By detecting and addressing data drift, we can:

- Improve the accuracy and reliability of machine learning models
- Avoid costly errors and lost opportunities
- Reduce the time and resources required to maintain and update models
- Increase transparency and accountability in decision-making processes

There are several techniques to measure data drift, such as:

- **Statistical tests:** These tests compare the distribution of data from different time periods to detect changes in statistical properties.
- **Model-based approaches:** These approaches involve comparing the performance of a model trained on data from different time periods to detect changes in model performance.

**Data drift can occur** in a wide range of scenarios, such as:

**E-commerce websites:** Changes in user behavior, such as a shift in the types of products being purchased or changes in purchasing patterns based on the time of year, can impact the accuracy of models used to recommend products to users.

**Financial services:** Changes in economic conditions or regulatory environments can impact the accuracy of models used to detect fraud or assess credit risk.

**Healthcare:** Changes in patient populations or treatment options can impact the accuracy of models used to predict health outcomes.

Data drift can occur in spam email detection ML models due to various reasons, some of which are:

- **Changes in the email content:** Spammers often change the content of their emails to bypass spam filters. This can result in changes in the statistical properties of the data used to train the model.
- **Changes in the email structure:** Spammers can also change the structure of their emails, such as adding or removing certain elements, to evade spam filters. This can result in changes in the features used to train the model.

- **Changes in the email source:** Spammers can change the source of their emails, such as using different email providers or IP addresses, to avoid being detected. This can result in changes in the distribution of the data used to train the model.

All of these changes can lead to data drift in the spam email detection ML models, causing a decrease in their performance over time. It is important to monitor the data used to train and test the model regularly to detect and address any data drift issues.

**KL divergence is a popular method for measuring data drift.** KL divergence measures the difference between two probability distributions. By calculating the KL divergence between the distribution of data at two different time periods, we can quantify the amount of data drift. A high KL divergence value indicates a significant difference between the two distributions, while a low KL divergence value indicates that the distributions are similar. Here are the steps we can follow:

1. Collect a historical data that you want to monitor – this data will be used as a baseline to compare against the new data
2. Calculate the probability distribution of each input feature for each time period (baseline and new data).
3. Calculate the KL divergence between each pair of probability distributions – this will give you a measure of the difference between the two distributions.
4. Set a threshold value for the KL divergence that indicates a significant enough difference between the distributions to signal data drift.
5. Compare the KL divergence values against the threshold. If any value exceeds the threshold, the model may be experiencing data drift and requires retraining.

Finally, data drift is a common challenge faced by data scientists and machine learning practitioners. Detecting and addressing data drift is crucial to maintaining the accuracy and reliability of machine learning models. By using techniques such as KL divergence, we can measure data drift and take proactive steps to address it.