

Machine Learning Task for TDT4173

(Modern Machine Learning in Practice)

Task Overview

This task is part of a consultancy project conducted by Append on behalf of Hydro.

Hydro is a Norwegian industrial company with operations in energy, aluminum, and recycling. Append is a consultancy firm that works with data science, artificial intelligence, and system development to help companies make better use of data and technology in their operations.

The aim of the project is to support improved production planning at one of Hydro's refineries by developing more accurate forecasts of incoming raw materials.

Task Description

You are provided with historical data on raw material deliveries and orders through the end of 2024. Each raw material is identified by a unique `rm_id`. The goal is to develop a model that forecasts the cumulative incoming weight of each raw material from January 1, 2025, up to any specified end date between January 1 and May 31, 2025.

For any end date within this range, the model should predict the total weight of a raw material delivered between January 1, 2025, and that date.

Dataset Overview

The datasets are organized as follows:

- **data/kernel/receivals.csv**: The primary dataset containing historical records of material receivals. Each entry includes a timestamp, the quantity received, and the corresponding `rm_id`.
- **data/kernel/sales_orders.csv**: Contains information on ordered quantities and expected deliveries.
- **data/extended/materials.csv** (Optional): Metadata on various raw materials, including categories and classifications.
- **data/extended/suppliers.csv** (Optional): Information about suppliers, potentially useful for identifying patterns in delivery reliability or frequency.

- **data/extended/transportation.csv** (Optional): Transportation-related data that could affect delivery times and consistency.

Evaluation

Quantile Error at 0.8 (Asymmetric Loss)

Let there be N raw materials indexed by $i = 1, \dots, N$. Over a forecasting window of h days, define

$$A_i = \sum_{t=1}^h y_{i,T+t}, \quad F_i = \sum_{t=1}^h \hat{y}_{i,T+t},$$

as the actual and forecast total deliveries for material i , respectively.

To evaluate performance, we compute the quantile loss at the 0.8 level:

$$\text{QuantileLoss}_{0.8}(F_i, A_i) = \max(0.8 \cdot (A_i - F_i), 0.2 \cdot (F_i - A_i)).$$

The overall metric is the average quantile loss across all materials:

$$\text{QuantileError}_{0.8} = \frac{1}{N} \sum_{i=1}^N \text{QuantileLoss}_{0.8}(F_i, A_i).$$

This metric penalizes overestimation more than underestimation, which aligns with the practical needs of smelting. If we underestimate the available materials, the smelt can usually continue with what is on hand. However, if we overestimate, we risk planning a smelt that cannot be completed due to missing resources. Therefore, it's better for the model to be slightly cautious and predict too little rather than too much.