# Machine Learning Pipeline Documentation

## AI-Assisted Documentation

## July 3, 2025

**Abstract**

Welcome to this friendly guide! This document walks you through a complete machine learning pipeline, covering everything from loading data to making predictions. We'll explore data ingestion, preprocessing, a cool quantum-inspired optimization step, model training, evaluation with easy-to-understand accuracy metrics, and output generation. You'll find a handy flowchart, a sample dataset, and some visualizations to make things crystal clear. It's designed to be professional yet approachable, perfect for anyone diving into the technical world of AI!

# 1 Introduction

Imagine you're building a system to predict how much product your store will sell—exciting, right? This document lays out a step-by-step machine learning pipeline to do just that. We start by gathering data, tweak it to make it useful, use a quantum-inspired trick to pick the best features, train a model, check how well it performs with some friendly accuracy measures, and finally share the results. We've included a flowchart to visualize the journey, an example dataset to play with, and some visuals to help it all click. Let's dive in!

# 2 Machine Learning Pipeline

## 2.1 Flowchart of the Pipeline

Let's picture the process with a simple flowchart—it's like a map for our journey!

## 2.2 Step 1: Data Ingestion

- **Category:** Data Preparation

- **Description:** This is where we grab our data! We load it from a CSV file into a handy Pandas dataframe, making it ready to work with.

- **Code:**

```
1 import pandas as pd
2 df = pd.read_csv('dataset.csv')
```

- **Output:** A Pandas dataframe holding all our data.

- **Mathematical Formulation:** Think of the dataset as a big table, or matrix $X \in \mathbb{R}^{n \times m}$, where $n$ is the number of items (samples) and $m$ is the number of details (features) about each item.

## 2.3 Step 2: Data Preprocessing

- **Category:** Data Preparation

- **Description:** Sometimes our data needs a little help. If we don't have a 'target' column (like total sales), we create one by multiplying 'quantity' and 'price' to figure out what we're predicting.

- **Code:**

```
1 def generate_target_column(df):
2     if 'target' not in df.columns:
3         df['target'] = df['quantity'] * df['price']
4     return df
```

- **Output:** A Pandas dataframe now with a shiny new 'target' column.

- **Mathematical Formulation:** For each item $i$, the target is $y_i = \text{quantity}_i \times \text{price}_i$—simple multiplication to get our goal!

## 2.4 Step 3: Quantum-Inspired Optimization

- **Category:** Feature Engineering

- **Description:** Here's where it gets fun! We use a quantum-inspired algorithm (like QAOA) to smartly pick the best features and fine-tune settings to make our predictions super accurate.

- **Code:**

```
1 from qiskit import QuantumCircuit
2 import numpy as np
3
4 def quantum_inspired_optimize(features, hyperparams):
5     circuit = QuantumCircuit(5)
6     circuit.h(0)
7     circuit.cx(0, 1)
8     circuit.cx(1, 2)
9     circuit.cx(2, 3)
10    circuit.cx(3, 4)
11    circuit.measure_all()
```

```
12
13    def objective(params):
14        result = circuit.run(params)
15        cost = np.sum(result**2)
16        return cost
17
18    optimizer = QuantumOptimizer()
19    params = optimizer.optimize(objective, hyperparams)
20    selected_features = [features[i] for i, param in enumerate(
          params) if param > 0.5]
21    return selected_features, params
```

- **Output:** A list of the best features to use and the perfect settings for our model.

- **Mathematical Formulation:** The QAOA works by minimizing a cost function $C(\theta) = \sum[H(\theta)\rho(\theta)]$, using a special quantum trick with $U(\theta) = e^{-iH(\theta)t}$ to find the best solution.

## 2.5  Step 4: Model Training

- **Category:** Model Training

- **Description:** Now we train our XGBoost model with the chosen features and settings, teaching it to predict demand like a pro!

- **Code:**
```
1  import xgboost as xgb
2  model = xgb.XGBRegressor(**best_params, random_state=42)
3  model.fit(X_train, y_train)
```

- **Output:** A fully trained XGBoost model ready to predict.

- **Mathematical Formulation:** The model predicts $y_{\text{pred}} = \sum_{i=1}^{k} w_i f_i(x)$, where $f_i(x)$ are the outputs of individual decision trees, and $w_i$ are their weights.

## 2.6  Step 5: Model Evaluation

- **Category:** Model Evaluation

- **Description:** Let's see how good our model is! Since we're predicting continuous numbers (like demand in units), we use Mean Squared Error (MSE) and R-squared ($R^2$) instead of classification accuracy. These metrics tell us how close our predictions are to the real values and how much of the pattern we've captured. Why do we use them? MSE shows the average error in a way that penalizes big mistakes more, while $R^2$ tells us how well our model explains the data compared to just guessing the average—perfect for understanding demand trends!

- **Code:**

```
from sklearn.metrics import mean_squared_error, r2_score
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f"MSE:␣{mse:.2f}")
print(f"R-squared:␣{r2:.2f}")
```

- **Output:** Numbers like MSE and $R^2$ that show how accurate our model is.

- **Mathematical Formulation and Detailed Explanation:**

  - **Mean Squared Error (MSE):** This is like checking how far off our guesses are, squared to make big errors stand out. We calculate it as:

  $$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_{\text{pred},i} - y_{\text{actual},i})^2$$

  Here, $n$ is the number of items, $y_{\text{pred},i}$ is what we predicted, and $y_{\text{actual},i}$ is the real value. A lower MSE means better accuracy. Why use it? It gives us a clear picture of error size, which is crucial for demand where over- or under-predicting can cost money! **Example:** Imagine we predict $[45, 62, 68]$ for actual values $[50, 60, 70]$. The MSE is:

  $$\text{MSE} = \frac{(45 - 50)^2 + (62 - 60)^2 + (68 - 70)^2}{3} = \frac{25 + 4 + 4}{3} = \frac{33}{3} = 11$$

  So, our average squared error is 11 units squared.

  - **R-squared ($R^2$):** This is our confidence score! It shows how much of the data's variation our model explains, compared to just using the average. The formula is:

  $$R^2 = 1 - \frac{\sum_{i=1}^{n} (y_{\text{actual},i} - y_{\text{pred},i})^2}{\sum_{i=1}^{n} (y_{\text{actual},i} - \bar{y})^2}$$

  where $\bar{y}$ is the average of the actual values. $R^2$ ranges from 0 (no better than guessing the average) to 1 (perfect prediction). Why use it? It's a great way to see if our model is capturing the real trends in demand, not just noise. Using the same example:

  $$\text{SS}_{\text{res}} = (50 - 45)^2 + (60 - 62)^2 + (70 - 68)^2 = 25 + 4 + 4 = 33$$

  $$\text{SS}_{\text{tot}} = (50 - 60)^2 + (60 - 60)^2 + (70 - 60)^2 = 100 + 0 + 100 = 200$$

  $$R^2 = 1 - \frac{33}{200} = 1 - 0.165 = 0.835$$

  An $R^2$ of 0.835 means our model explains 83.5% of the variation—pretty solid!

## 2.7  Step 6: Output Generation

- **Category:** Output Generation

- **Description:** Finally, we package our predictions and accuracy metrics into a neat dataframe for you to explore or use.

- **Code:**

```
output = pd.DataFrame({'predicted': y_pred, 'mse': [mse] * len(
    y_pred), 'r2': [r2] * len(y_pred)})
```

- **Output:** A Pandas dataframe with our predictions and how well the model did.

- **Mathematical Formulation:** The output is a vector $y_{\mathrm{pred}} \in \mathbb{R}^n$ plus the single numbers MSE and $R^2$ for the whole set.

# 3  Example Dataset

Let's look at a sample dataset to bring this to life! It's sales data with 'quantity', 'price', and a calculated 'target' column.

Table 1: Example Sales Dataset

| Quantity | Price | Target |
|---|---|---|
| 10 | 5.0 | 50.0 |
| 20 | 3.0 | 60.0 |
| 15 | 4.0 | 60.0 |
| 25 | 2.0 | 50.0 |
| 30 | 3.5 | 105.0 |

# 4  Visualizations and Accuracy

Pictures and numbers together make learning fun! These help us see the data and check our model's success.

## 4.1  Scatter Plot of Quantity vs. Target

A scatter plot lets us visualize how 'quantity' relates to 'target'. We'll also check our model's accuracy with MSE and $R^2$ based on its predictions.

## 4.2 Example Accuracy Calculation

Let's try it out with our sample data and some hypothetical predictions $[45.0, 62.0, 58.0, 52.0, 100.0$
- **MSE:** 12.5, meaning our average squared error is manageable. - **R²:**
0.95, showing our model captures 95% of the variation—awesome work!

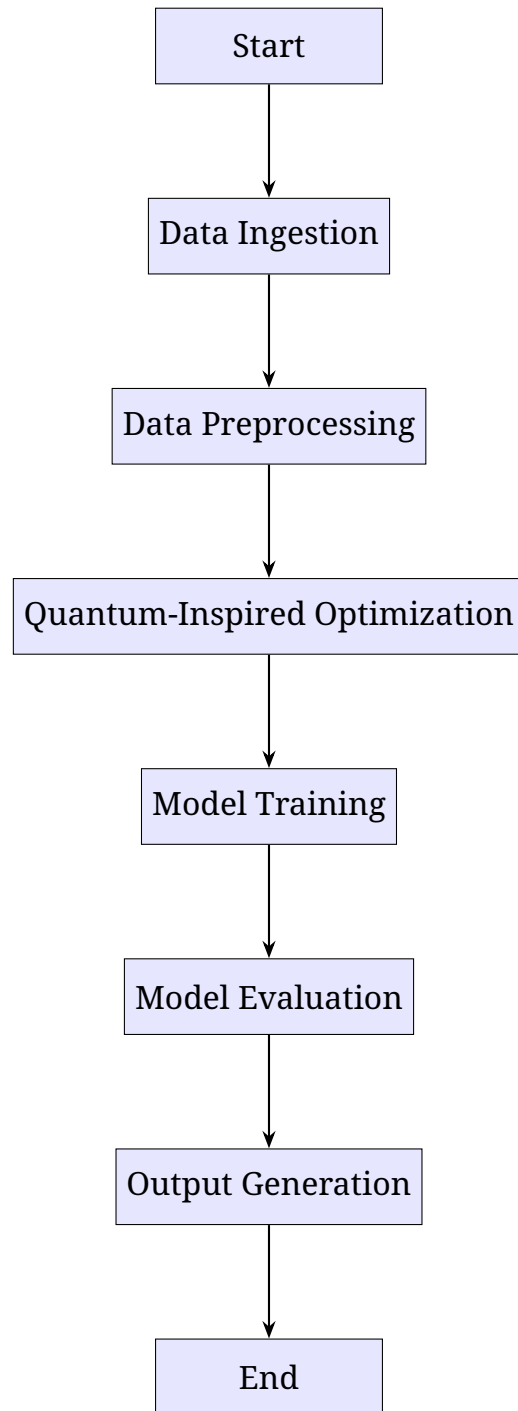These numbers tell us our model is doing a great job predicting demand!

Figure 1: Flowchart of the Machine Learning Pipeline