

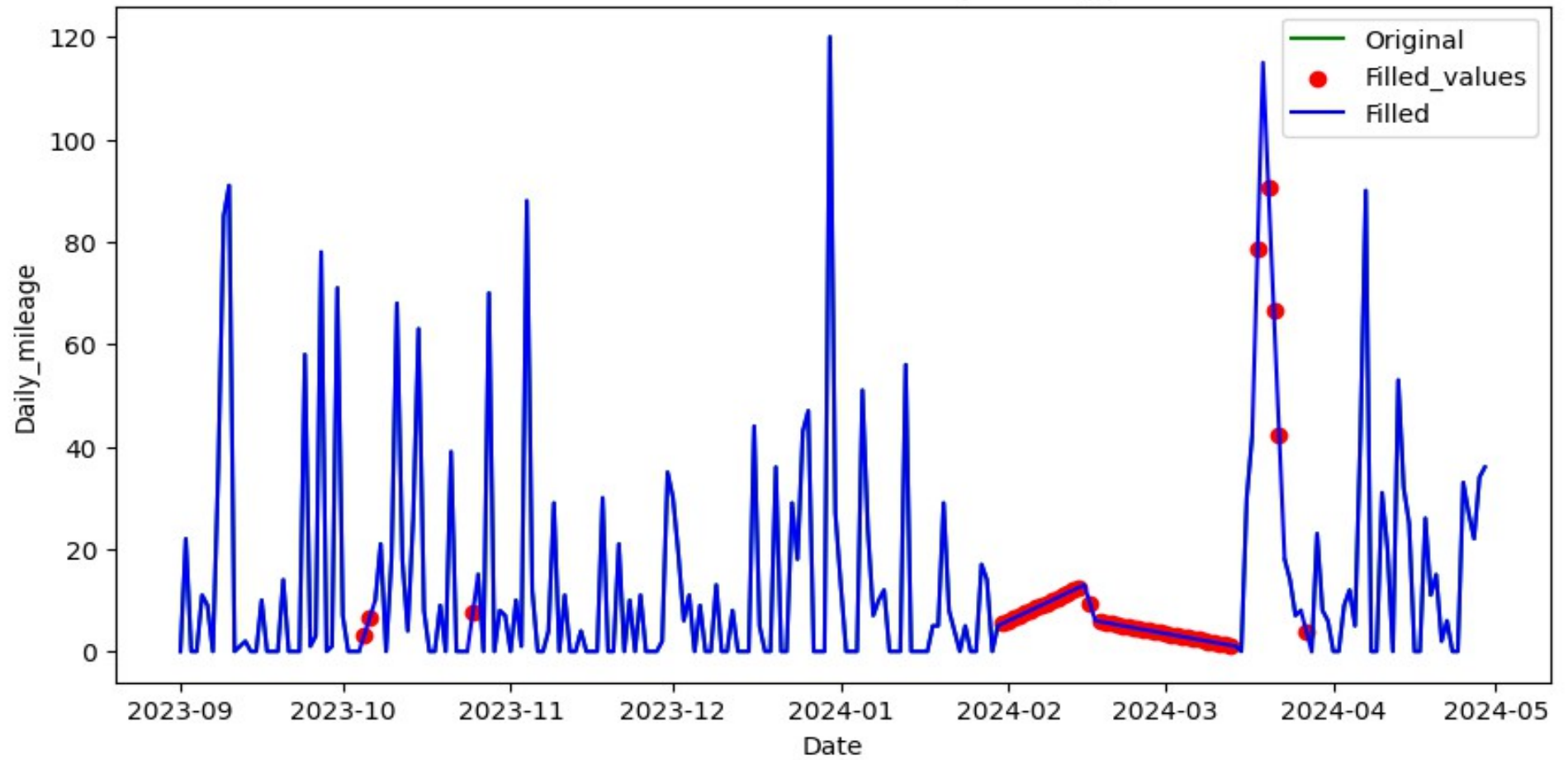
Mileage Forecasting

- I. Preprocessing and Model Training
- II. Model Deployment
- III. Testing

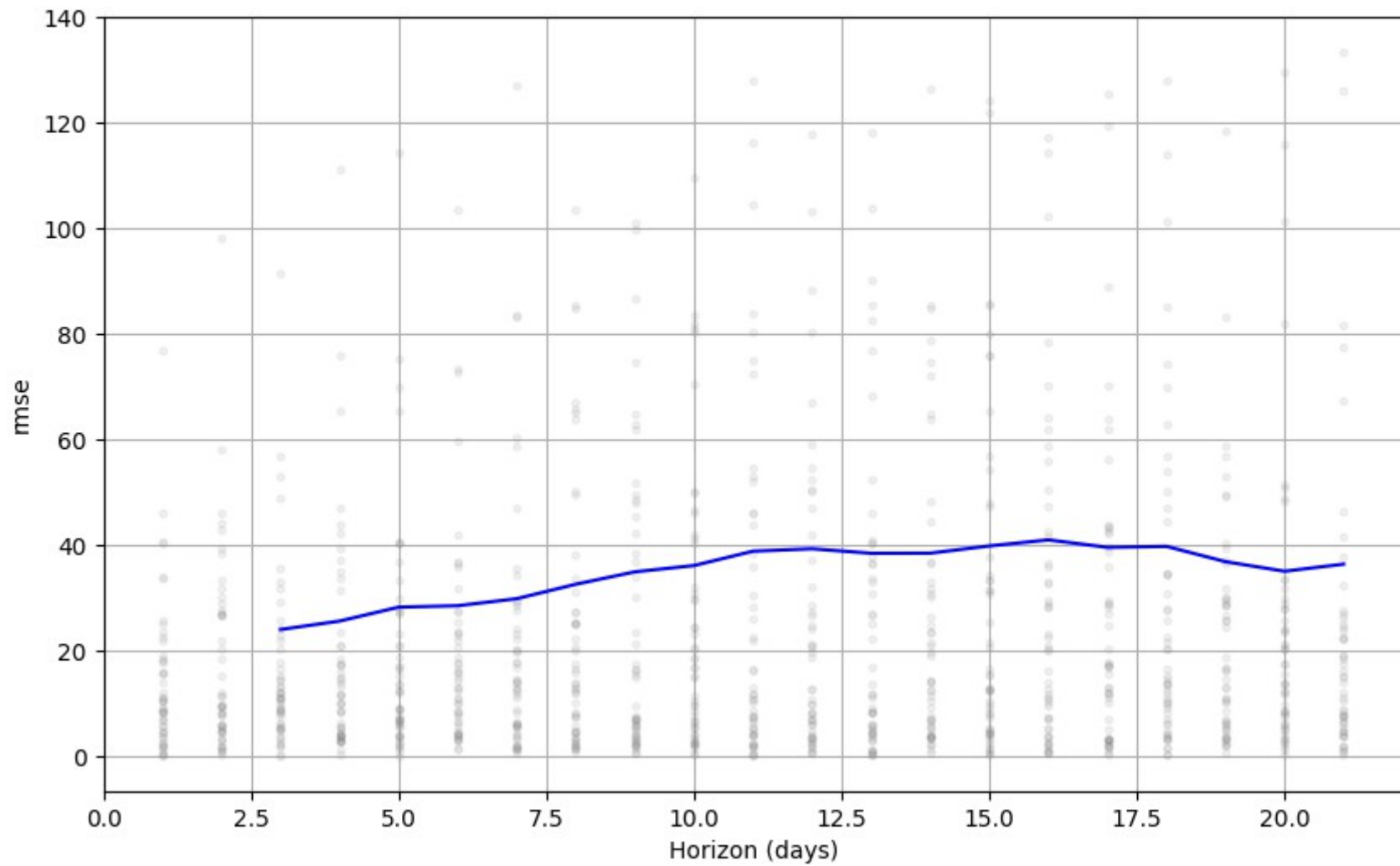
Preprocessing and Model Training

- First step: Exploratory Data Analysis
 - Exploring missing values
 - Identifying outliers
 - Exploring model (prophet) performance with different methods of handling missing values
- Handling of missing values
 - Best results achieved by setting missing values to constants like zero, mean or median
 - Assumption: Missing rows might indicate no vehicle activity
 - Choice of method: setting to zero

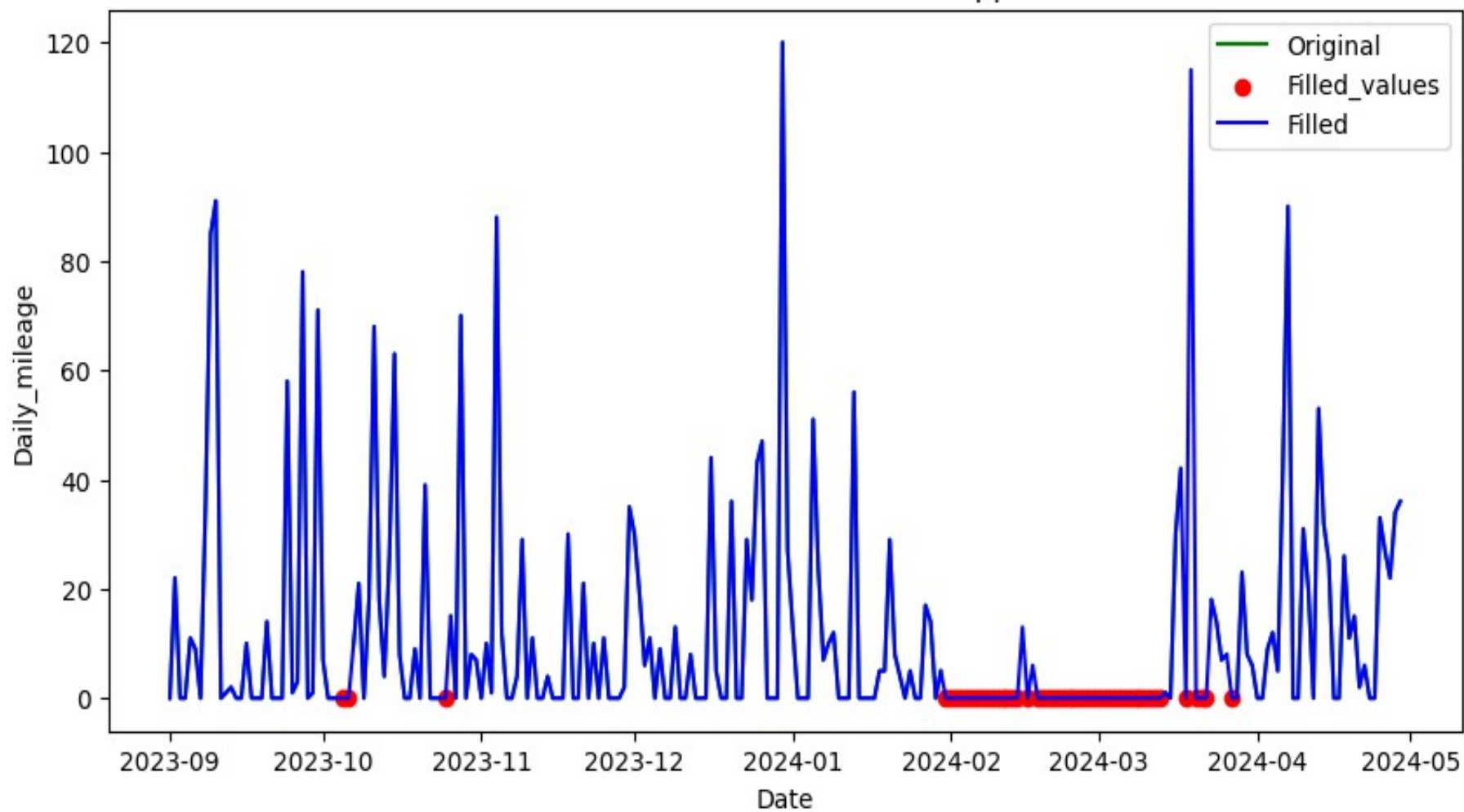
Time Series Data filled with interpolate approach



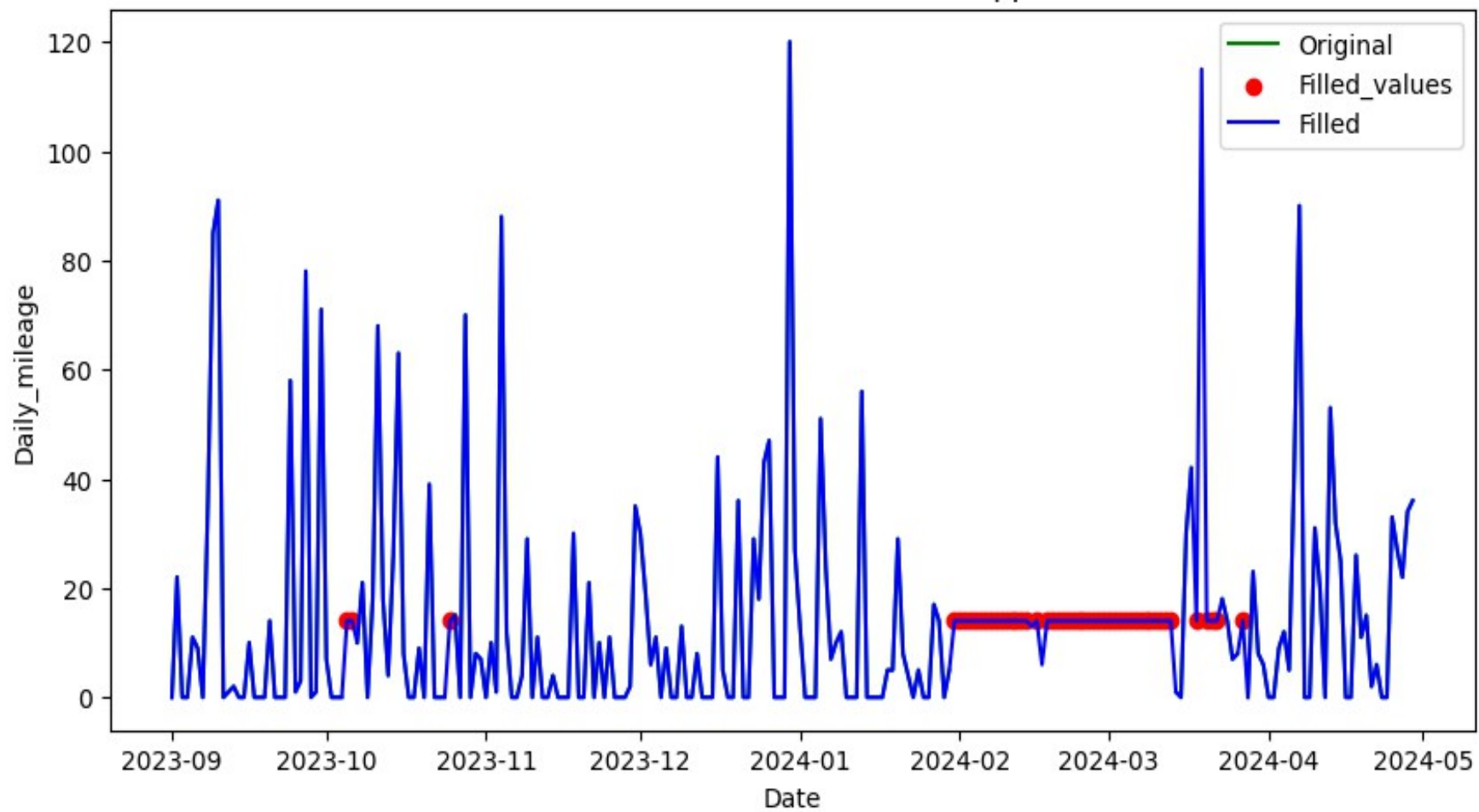
model with MV set to mean



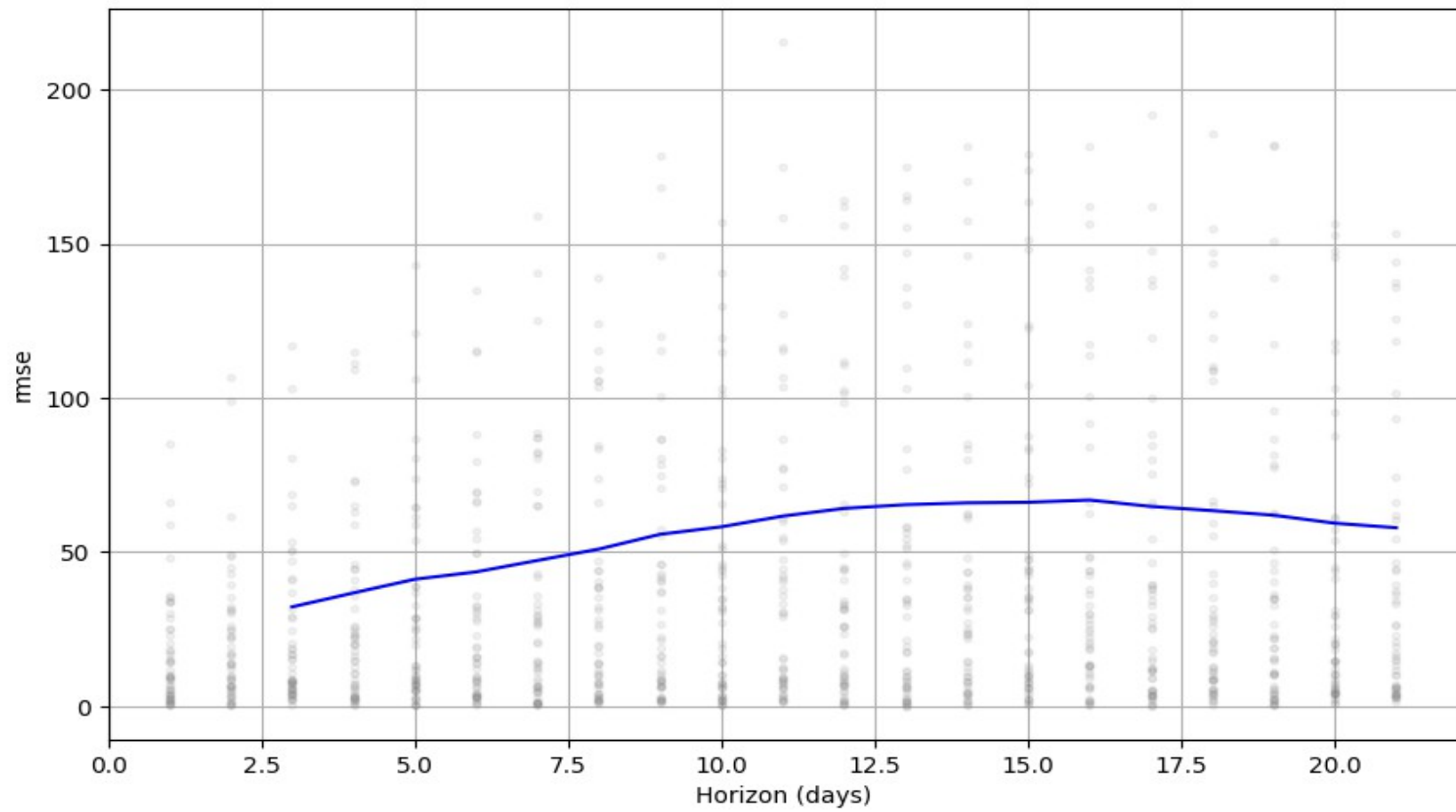
Time Series Data filled with zero approach



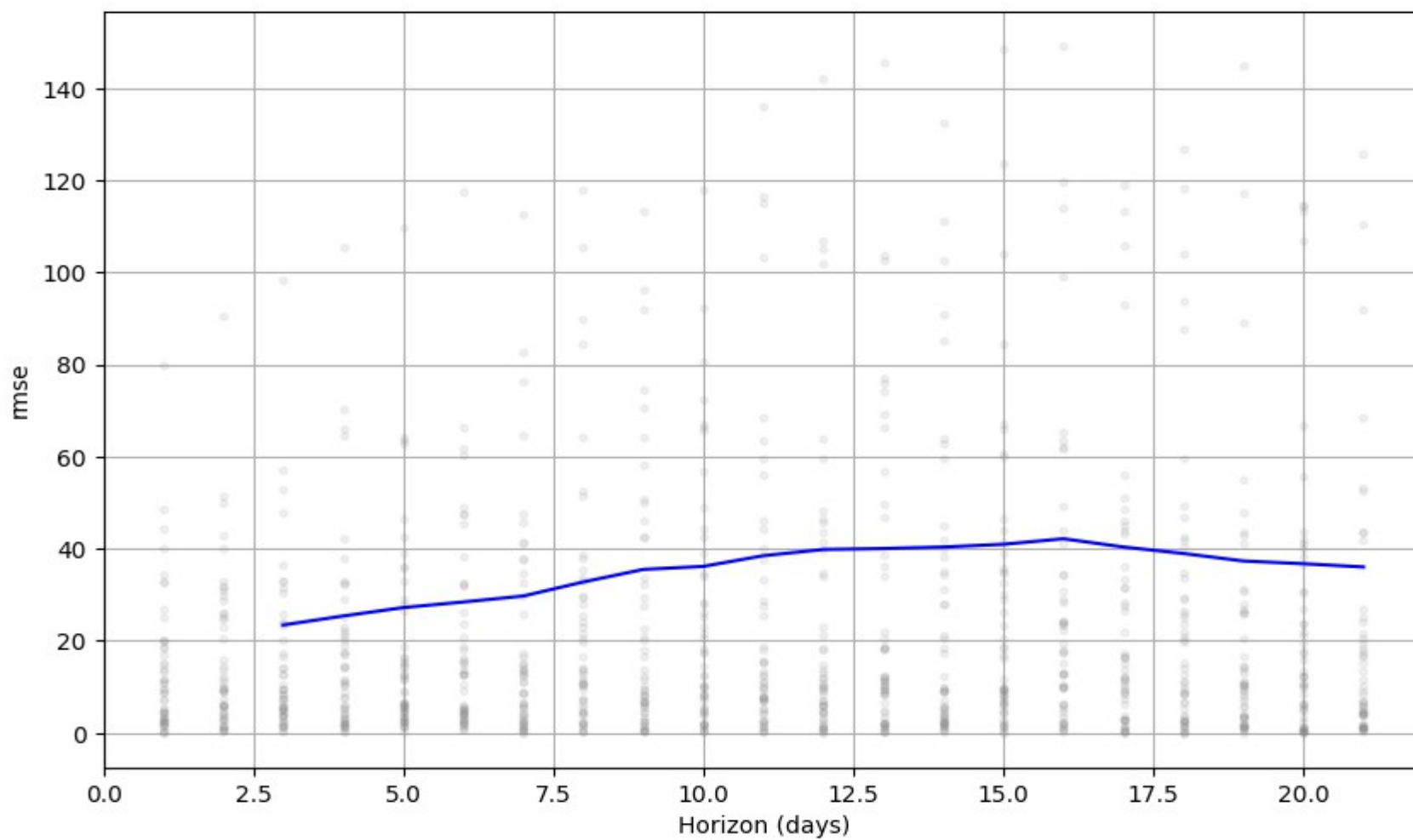
Time Series Data filled with mean approach



model with MV interpolated



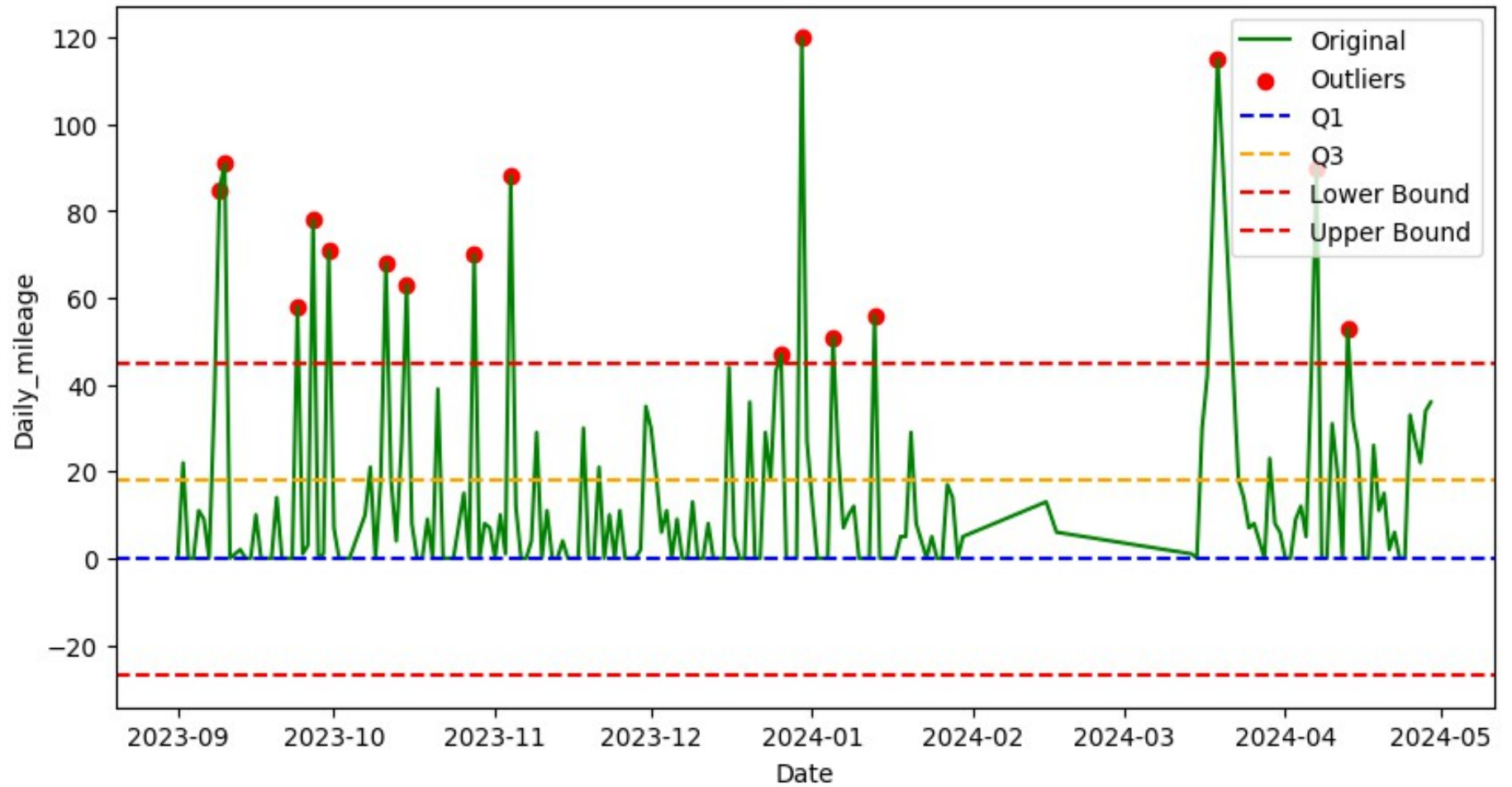
model with MV set to zero



Preprocessing and Model Training

- Identify outliers by Interquartile Range Method
 - Q1: 25th percentile: value below which 25% data falls
 - Q3: 75th percentile: value below 75% of data falls
 - Interquartil: Range between Q1 and Q3
 - Outliers: everything that is below lower bound = $1,5 * Q1$ and above $1,5 * Q3$
- Chose method to handle outliers:
 - Set them to NaN as recommended in the Prophet documentation

Time Series Data Outliers



Model Evaluation

- Instead of splitting to train and testing sets:
Cross Validation
 - Done by rolling time window for time series
 - Advantage:
 - More data for training
 - More than just a single evaluation
 - More reliability on performance estimates
 - Better detection of overfitting
- Metric : RMSE Root mean square error
 - Intuitive interpretation as it is in the same unit as data
 - Penalize larger errors more than smaller errors

Hyperparameter Tuning

- Use grid search and cross-validation to identify best hyperparameters for
 - Changepoint:
 - Detects when trend changes and adapts predictions
 - Specify potential change points
 - Seasonality: model can try to identify seasonal patterns

• Model Deployment

- Core Idea

- Deploy model automatically with Kubeflow Pipeline:
 - Takes care of containerization of components to have reusable components
 - Portability : can be run locally or in cloud
 - Track and monitor models
 - Efficiency: Optimize resource usage
 - Runs on top of Kubernetes

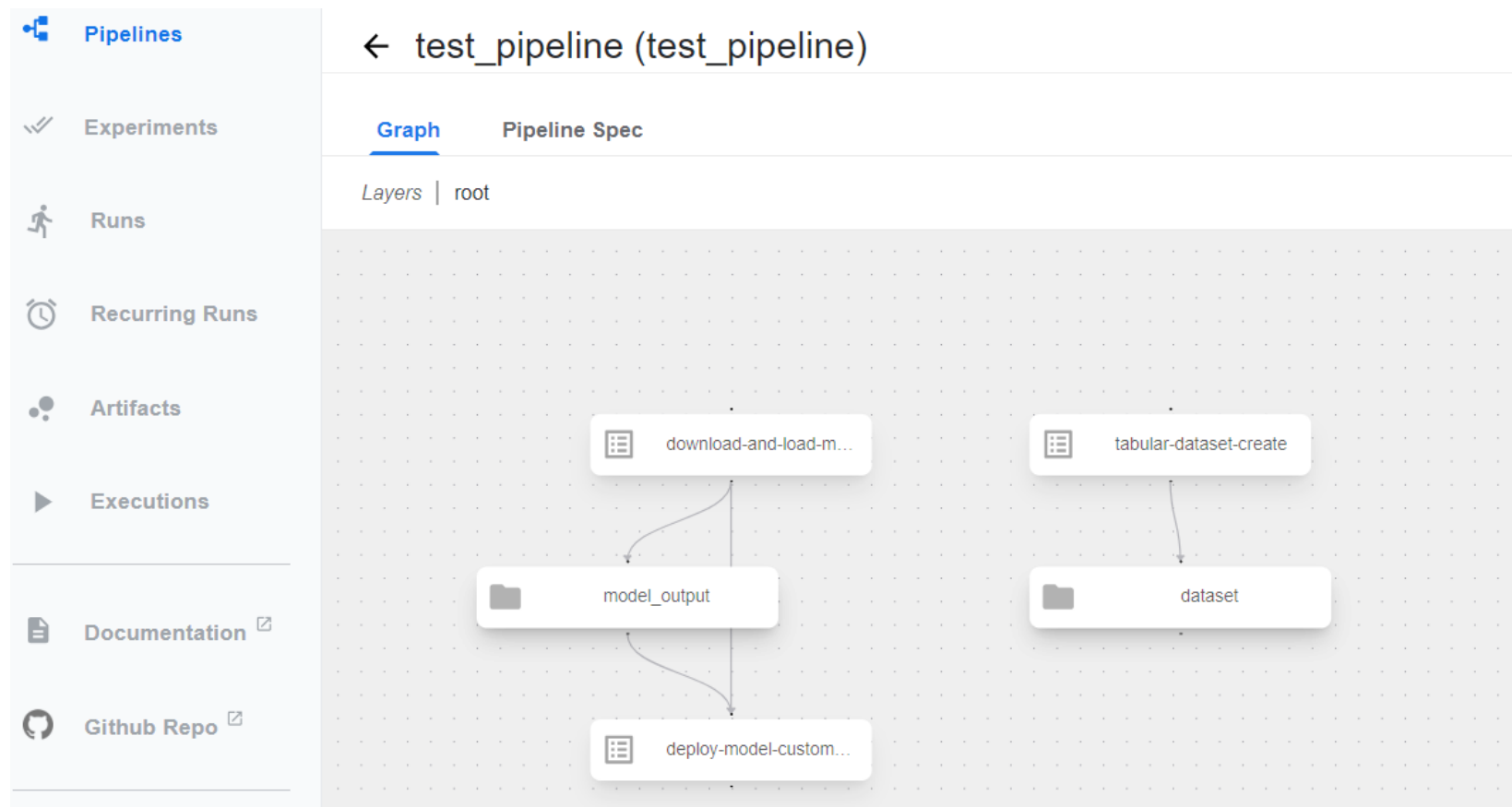
- Ways I tried to do this:
 - Initial idea: Use Google Pipeline components to deploy model: (<https://cloud.google.com/vertex-ai/docs/pipelines/build-pipeline>)
 - Problem it does not support Prophet library
 - Containerize model, load it to Vertex Ai Model registry and run Kubeflow pipeline from there:
 - Problem: Vertex Ai model registry does not support Prophet library models
 - Load model files from Github by writing custom component and deploy it with a custom component
 - Problem: very difficult to set up Kubeflow locally on windows, lost a lot of time fixing this

- Ways i tried to solve this issue with Kubeflow
 - Install kubernetes minikube locally and install Kubeflow on the minikube container:
 - 1.Problem: containersize by default too small
 - Solved by changing default container size in WSL
 - 2. Problem: kubeflow installation on minikube failes
 - Install kubeflow using kind following this tutorial :
 - <https://www.kubeflow.org/docs/components/pipelines/v1/installation/localcluster-deployment/>

- Further problems: Compling the pipeline the Client could not find the host
 - Solved by forwarding port to 3000 :
mlpipeline)
 - Then kubeflow could be opend in the browser
(localhost:3000)
- ```
C:\Users\aisa\MLPipeline\MLPipeline\src> kubectl
port-forward svc/ml-pipeline-ui 3000:80
--namespace kubeflow
```



- Then I ran out of time to properly debug my pipeline



- Lesson learned:
  - When you want to use Google Best Practise better use a Tensorflow model

# Testing

- First step: create a python project from the notebook
- Set up unit test for modules, preprocessing, training
- Kubeflow pipeline can also be covered with unit test

# Testing

- Automated testing with github action



2 workflow runs



Merge branch 'main' of <https://github.com/AishaLichtner/MLPipeline>

main

Python application #18: Commit [1f28173](#) pushed by AishaLichtner



Create python-app.yml

main

Python application #17: Commit [edfc711](#) pushed by AishaLichtner

- Problem setting up the automated testing
  - Python path wasnt set correctly, so my modules could not be found and therefore unit tests could not run
  - Solved by adding a job to pipeline to set the python path

```
- name: Set Python Path
 run: |
 echo "PYTHONPATH=$(pwd)/src" >> $GITHUB_ENV
```

## build

succeeded 7 hours ago in 56s

- > ☒ Set up job
- > ☒ Run actions/checkout@v4
- > ☒ Set up Python 3.7.16
- > ☒ Install dependencies
- > ☒ Set Python Path

### ✓ Test with pytest

```
1 ▶ Run pytest
8 ===== test session starts =====
9 platform linux -- Python 3.7.16, pytest-7.4.4, pluggy-1.2.0
10 rootdir: /home/runner/work/MLPipeline/MLPipeline
11 collected 6 items
12
13 tests/test_preprocess.py [100%]
14
15 ===== 6 passed in 1.35s =====
```

- > ☒ Post Set up Python 3.7.16
- > ☒ Post Run actions/checkout@v4
- > ☒ Complete job

# Things I would have done with more time

- Modulize code better, eg.
  - Add config file to avoid hard coding
  - Make it more flexible for different structured input data
- Try github actions for model deployment
- Try combination of github action for containerization and loading to GCP and Kubeflow pipeline.