# THE real-world MACHINE LEARNING TUTORIAL

BY PAU LABARTA BAJO

# Why is this course different?

#### In this course you will...

- Solve a **real-world** prediction problem.

#### In this course you will...

- Solve a **real-world** prediction problem.

 Build a complete ML service, as opposed to just training an ML model.

#### In this course you will...

- Solve a **real-world** prediction problem.

 Build a complete ML service, as opposed to just training an ML model.

Learn all the tools you need to solve the problems YOU
 CARE ABOUT.

# How are we going to get there?

#### 1. Business problem 💼

A ride-sharing company wants to maximize revenue

1. Business problem 💼

A ride-sharing company wants to maximize revenue

2. Data preparation 🔣

From raw data to tabular data for Supervised ML

1. Business problem 💼

A ride-sharing company wants to maximize revenue

2. Data preparation ...
From raw data to tabular data for Supervised ML

3. Model training Tra

- 1. Business problem A ride-sharing company wants to maximize revenue
- 2. Data preparation From raw data to tabular data for Supervised ML
- 3. Model training T First without ML, then we add ML.
- 4. Model operationalization (MLOps) 
  Solution of the second of

- 1. Business problem A ride-sharing company wants to maximize revenue
- 2. Data preparation From raw data to tabular data for Supervised ML



- 3. Model training The First without ML, then we add ML.
- 4. Model operationalization (MLOps) 

  Without this our model has 0 business value

1. Business problem A ride-sharing company wants to maximize revenue

- 2. Data preparation From raw data to tabular data for Supervised ML
- 3. Model training T First without ML, then we add ML.
- 4. Model operationalization (MLOps) 
  Without this our model has 0 business value





# 1. Business problem

#### **Business Problem**

#### **Business Problem**

- You work as a data scientist and in a ride-sharing app company (e.g. Uber)

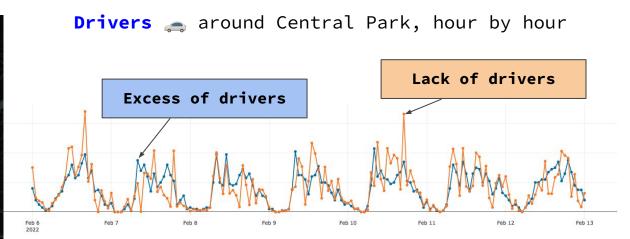
 Your job is to help the operations team keep the fleet as busy as possible.

# Supply and demand 🩋

**Drivers** around Central Park, hour by hour North Bergen Feb 13 NEW YORK Users looking for a ride 🙋 around Central Park

# Supply and demand 🥷

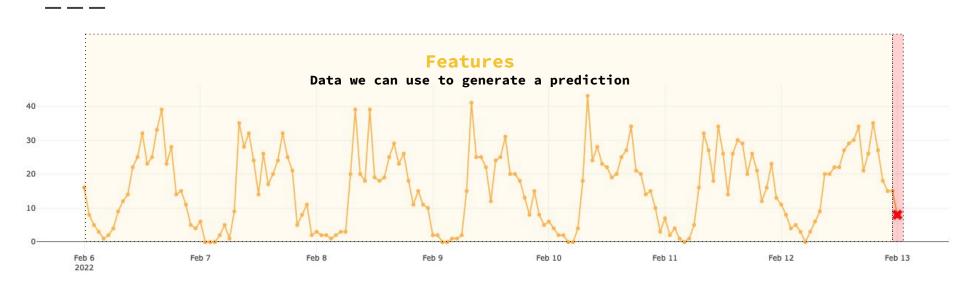
NEW YORK



Users looking for a ride 🙋 around Central Park



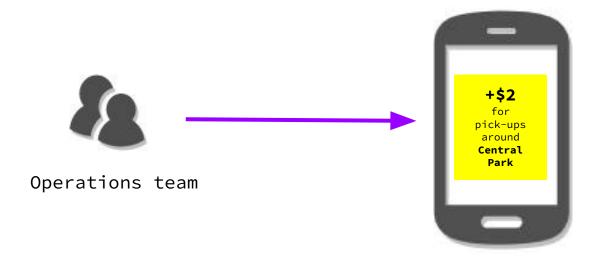
#### Let's predict user demand in the next hour



#### Target

This is what we want to predict #users who will request a ride in the next hour

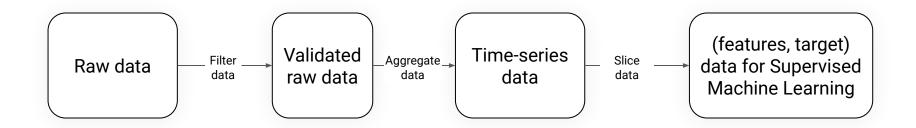
#### Ops team can adjust the distribution of the drivers

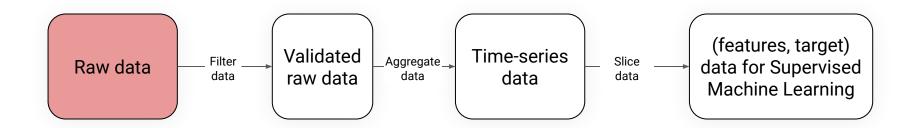


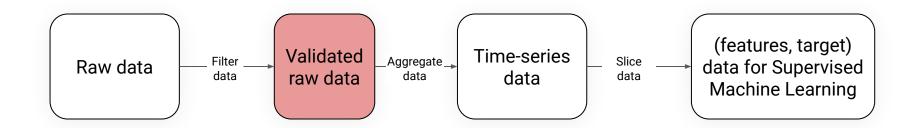
# Development environment **\***

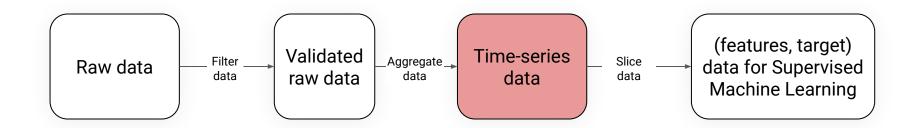
- **VSCode** for Python development
  - Other options: JupyterLab, PyCharm, Sublime...
- Poetry for virtual environments and packaging
  - Other options: virtualenv, conda, pipenv
- Git and GitHub for version control

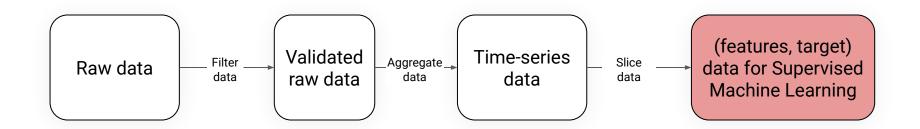
# 2. Data preparation



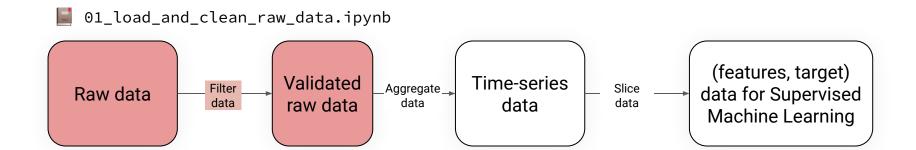




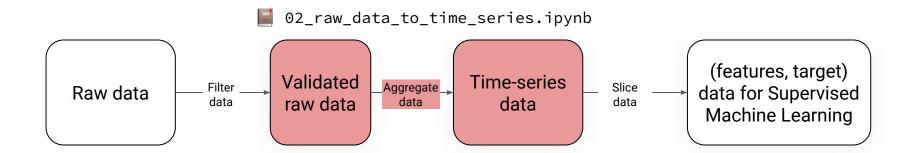




#### Step 1 - Data validation



#### Step 2 - Raw data into time-series data



# Step 3 - Time-series data into (features, target) data

Raw data

Filter data

Validated raw data

Validated raw data

Validated raw data

Filter data

Validated raw data

Validated raw data

Validated raw data

Validated raw data

Aggregate data

Aggregate data

Validated raw data

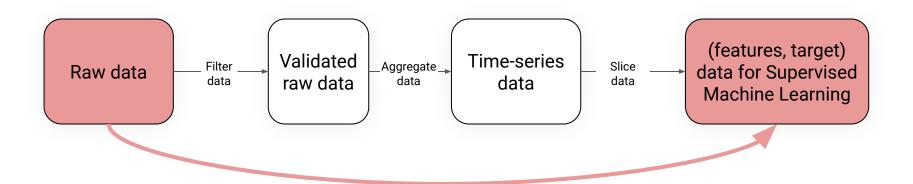
Aggregate data

Aggregate data

Validated raw data

### Step 4 - From raw data to training data

\_\_\_\_



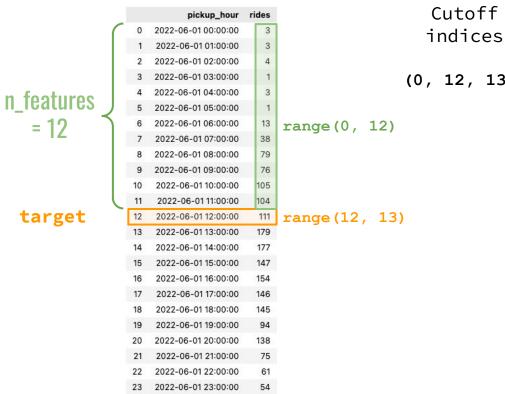
04\_raw\_data\_to\_features\_and\_target.ipynb

Slice & Slide

#### Supervised ML Tabular data

location id = 43

2022-06-02 00:00:00



27

features

target

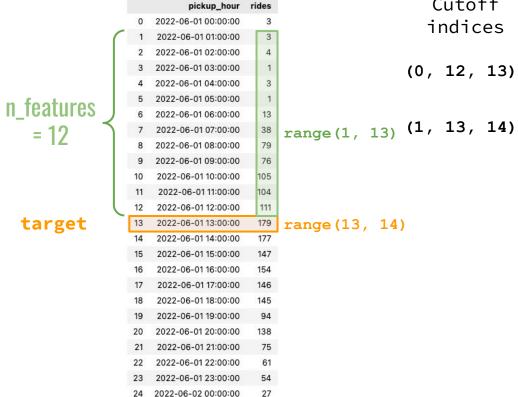
(0, 12, 13)

[ 3., 3., 4., ..., 76., 105., 104.]

Slice & Slide

#### Supervised ML Tabular data

location id = 43



Cutoff

(0, 12, 13)

features

111.

target

[ 3., 4., 1., ..., 105., 104., 111.]

[ 3., 3., 4., ..., 76., 105., 104.]

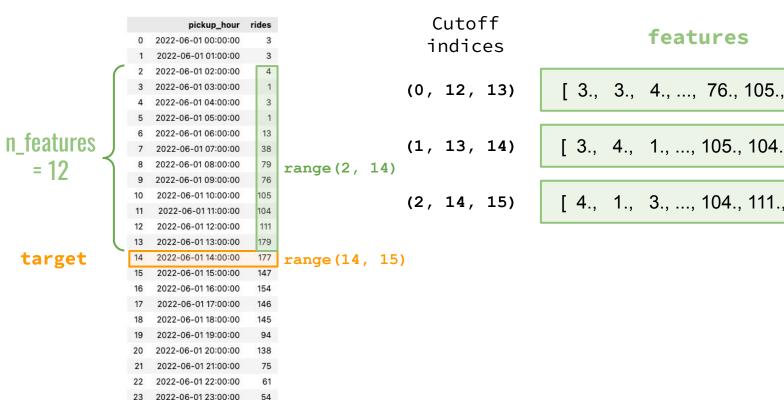
#### Slice & Slide

#### Supervised ML Tabular data

location id = 43

2022-06-02 00:00:00

27



target

111.

[ 3., 3., 4., ..., 76., 105., 104.]

[ 3., 4., 1., ..., 105., 104., 111.]

179.

[ 4., 1., 3., ..., 104., 111., 179.]

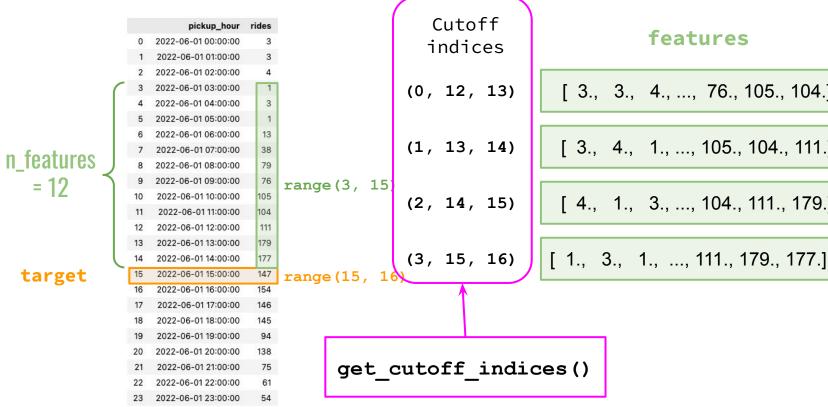
location id = 43

2022-06-02 00:00:00

27

#### Slice & Slide

#### Supervised ML Tabular data



target

[ 3., 3., 4., ..., 76., 105., 104.]

[ 3., 4., 1., ..., 105., 104., 111.]

[ 4., 1., 3., ..., 104., 111., 179.]

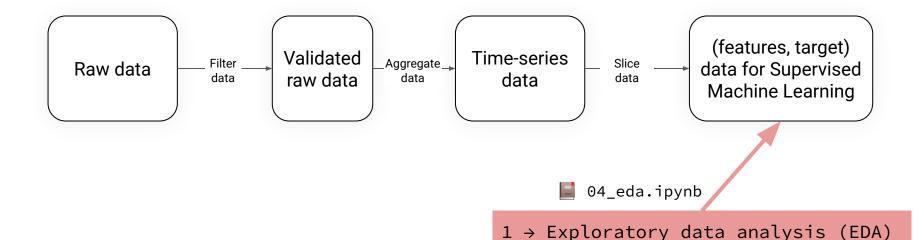
147.

111.

179.

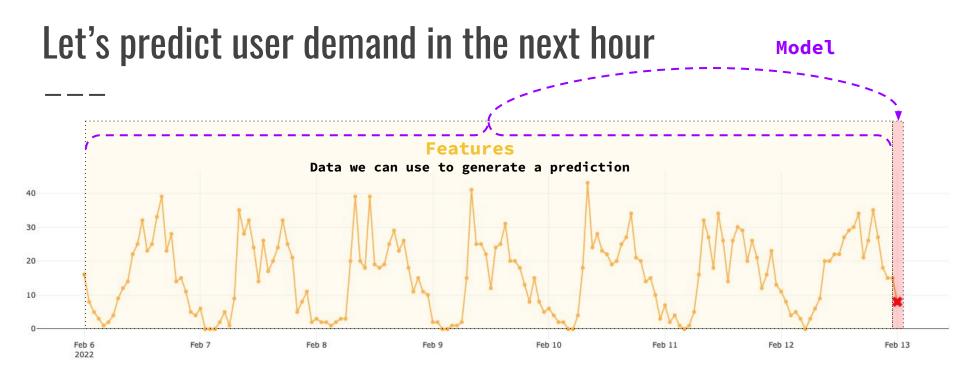
## Step 5 - Explore and visualize the final dataset

\_\_ \_ \_



2 → Plot individual examples

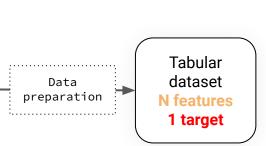
# 3. Model training



#### Target

This is what we want to predict #users who will request a ride in the next hour

# Split the data

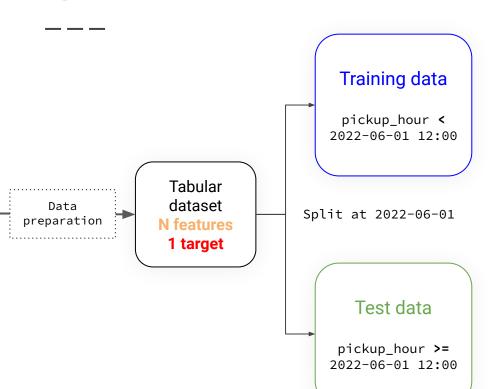


#### features

### target

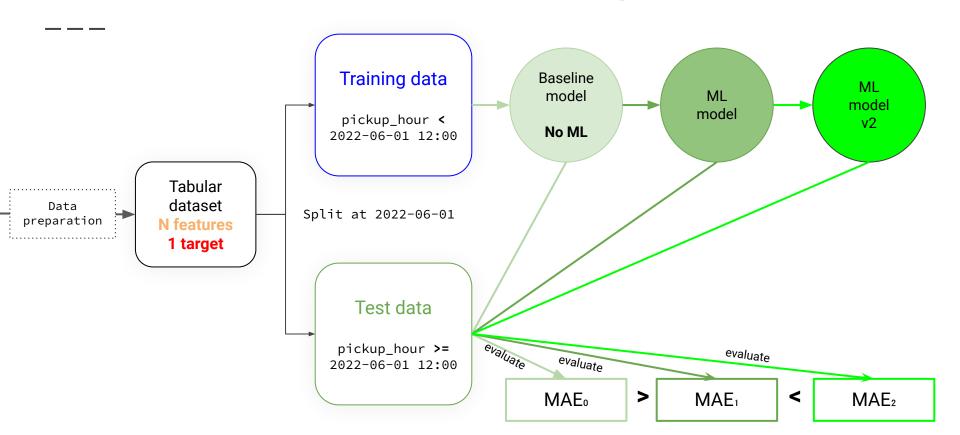
rides_24_ho urs_ago	rides_23_ho urs_ago	 rides_1_h ur_ago	pickup_ hour	pickup_lo cation_id	target_ride s_next_hou r
341	452	 232	2022-0 1-01 12:00	43	453
452	213	 566	2022-0 6-01 13:00	211	322
452	213	 453	2022-0 6-01 12:00	43	212
29	34	 54	2022-1 0-01 12:00	256	122

# Split the data

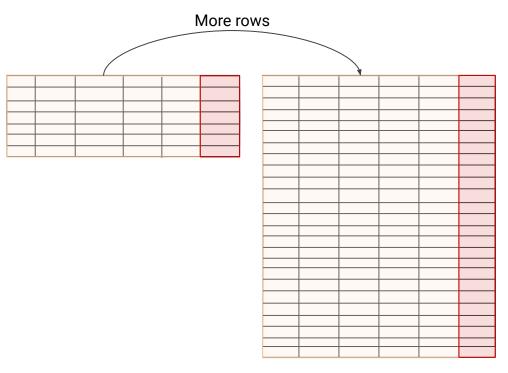


rides_24_ho urs_ago	rides_23_ho urs_ago	 rides_1_ho ur_ago	pickup_ hour	pickup_lo cation_id	target_ride s_next_hou r
341	452	 232	2022-0 1-01 12:00	43	453
452	213	 566	2022-0 6-01 13:00	211	322
452	213	 453	2022- 06-01 12:00	43	212
29	34	 54	2022-1 0-01 12:00	256	122

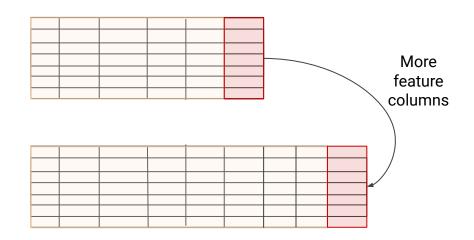
### How do you build a good ML model? Steps



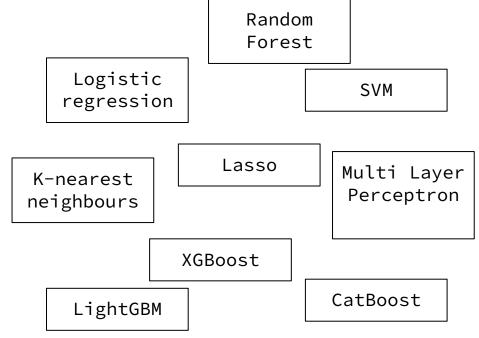
1. Increase training data size



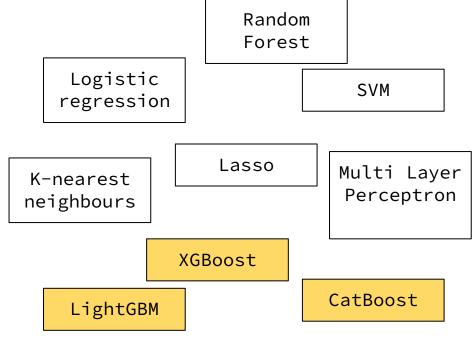
- Increase training data size
- 2. Add more features to the training data



- Increase training data size
- 2. Add more features to the training data
- 3. Try another algorithm



- Increase training data size
- 2. Add more features to the training data
- 3. Try another algorithm



**Boosting algorithms** 

- Increase training data size
- 2. Add more features to the training data
- 3. Try another algorithm
- 4. Tune algorithm hyper-parameters

#### A few of XGBoost hyper-parameters

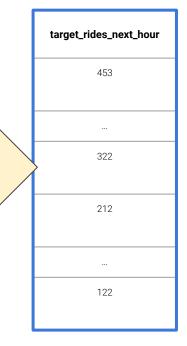
```
params = {
    'max_depth': trial.suggest_int('max_depth', 1, 9),
    'learning_rate': trial.suggest_loguniform('learning_rate', 0.01, 1.0),
    'n_estimators': trial.suggest_int('n_estimators', 50, 500),
    'min_child_weight': trial.suggest_int('min_child_weight', 1, 10),
    'gamma': trial.suggest_loguniform('gamma', 1e-8, 1.0),
    'subsample': trial.suggest_loguniform('subsample', 0.01, 1.0),
    'colsample_bytree': trial.suggest_loguniform('colsample_bytree', 0.01, 1.0),
    'reg_alpha': trial.suggest_loguniform('reg_alpha', 1e-8, 1.0),
    'reg_lambda': trial.suggest_loguniform('reg_lambda', 1e-8, 1.0),
    'eval_metric': 'mlogloss',
    'use_label_encoder': False
}
```

# Feature engineering

#### features

rides_24_ hours_ago	rides_23_h ours_ago	 rides_1_h our_ago
341	452	 232
452	213	 566
452	213	 453
29	34	 54

### target



LightGBM

# **Idea #1** $\rightarrow$ Add average rides 7, 14, 21 and 28 days ago

	featu	res		New feature		targ
rides_24_ hours_ago	rides_23_h ours_ago		rides_1_h our_ago	avg_rides_last_4 _weeks		target_rides_n
341	452		232	342		453
452	213		566	432	LightGBM	322
452	213		453	287		212
29	34		54	45		122

### **Idea #2** → Extract cyclical features from 'pick\_hour'

	featu	res		New feature	datetime		target
rides_24_ hours_ago	rides_23_h ours_ago		rides_1_h our_ago	avg_rides_last_4 _weeks	pickup_hour		target_rides_next_hour
341	452		232	342	2022-01-01 12:00		453
452	213		566	432	2022-06-01 13:00	LightGBM	322
452	213		453	287	2022-06-01 12:00		212
29	34		54	45	2022-10-01 12:00		122

# **Idea #2** → Extract cyclical features from 'pick\_hour'

	featu	res		New feature	New fe	eatures		target
rides_24_ hours_ago	rides_23_h ours_ago		rides_1_h our_ago	avg_rides_last_4 _weeks	hour	week_day		target_rides_next_hour
341	452		232	342	12	5		453
							<u> </u>	
452	213		566	432	13	2	LightGBM	322
452	213		453	287	12	2		212
29	34		54	45	12	5		122

### **Idea #3** → Encode categorical feature 'pickup\_location\_id'

	featu	res		New feature	New fe	eatures	categorical		target
rides_24_ hours_ago	rides_23_h ours_ago		rides_1_h our_ago	avg_rides_last_4 _weeks	hour	week_day	pickup_location_id		target_rides_next_ hour
341	452		232	342	12	5	23		453
452	213		566	432	13	2	23	LightGBM	322
452	213		453	287	12	2	256		212
29	34		54	45	12	5	256		122

# **Idea #3** → Encode categorical feature 'pickup\_location\_id'

	featu	res		New feature	New fe	eatures	New fe	atures		target
rides_24_ hours_ago	rides_23_h ours_ago		rides_1_h our_ago	avg_rides_last_4 _weeks	hour	week_day	latitude	longitude		target_rides_next_ hour
341	452		232	342	12	5	40.13	74.3		453
452	213		566	432	13	2	40.13	74.3	LightGBM	322
452	213		453	287	12	2	40.32	74.5		212
29	34		54	45	12	5	40.32	74.5		122

rides_24_ hours_ago	rides_23_h ours_ago	 rides_1_h our_ago
341	452	 232
29	34	 54

### Scikit-learn pipeline

**Transformation**add\_average\_rides
\_last\_4\_weeks

**Transformation**add\_temporal\_fe
atures

Estimator
LightGBM
Regressor

target_rides_next_hou	ır
453	
122	

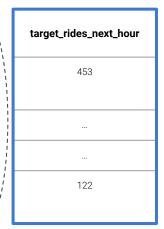
rides_24_ hours_ago	rides_23_h ours_ago	 rides_1_h our_ago
341	452	 232
29	34	 54

### Scikit-learn pipeline

**Transformation**add\_average\_rides
\_last\_4\_weeks

**Transformation**add\_temporal\_fe
atures

Estimator
LightGBM
Regressor



#### 1. Create scikit-learn pipeline

```
from sklearn.pipeline import make_pipeline
pipeline = make_pipeline(
   add_average_rides_last_4_weeks,
   add_temporal_features(),
   LGBMRegressor()
)
```

### X\_train

rides_24_ hours_ago	rides_23_h ours_ago	 rides_1_h our_ago
341	452	 232
29	34	 54

### Scikit-learn pipeline

Transformation
add\_averages\_ride
s\_last\_4\_weeks
.fit()

Transformation
add\_temporal\_fe
atures
.fit()

Estimator
LightGBM
Regressor
.fit()

y\_train

	target_rides_next_hour
	453
1	
	122

#### 1. Create scikit-learn pipeline

```
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    add_average_rides_last_4_weeks,
    add_temporal_features(),
    LGBMRegressor()
)
```

#### 2. Train

pipeline.fit(X\_train, y\_train)

#### **X\_test**

				•	
n	re	rb	ct	า ด	ns
r		•		. •	

rides_24_ hours_ago	rides_23_h ours_ago	 rides_1_h our_ago
341	452	 232
29	34	 54

### Scikit-learn pipeline

Transformation
add\_averages\_ride
s\_last\_4\_weeks
.transform()

Transformation
add\_temporal\_fe
atures
.transform()

Estimator
LightGBM
Regressor
.predict()

target\_rides\_next\_hour

453

...

122

#### 1. Create scikit-learn pipeline

LGBMRegressor()

```
from sklearn.pipeline import make_pipeline
pipeline = make_pipeline(
   add_average_rides_last_4_weeks,
   add_temporal_features(),
```

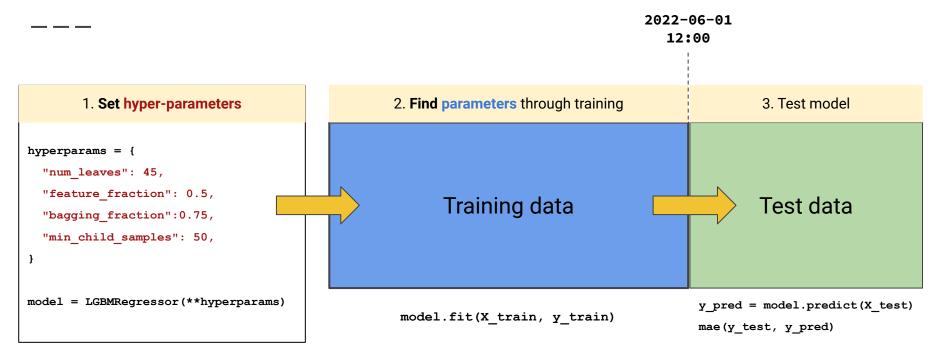
```
2. Train
```

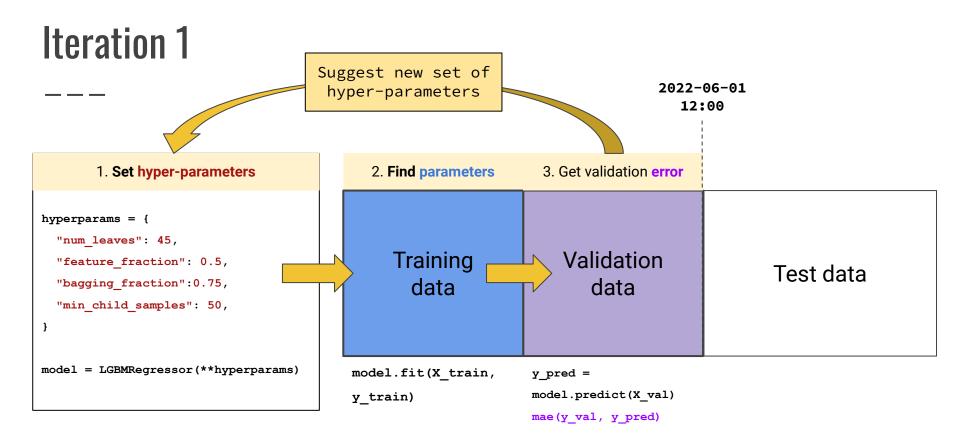
```
pipeline.fit(X_train, y_train)
```

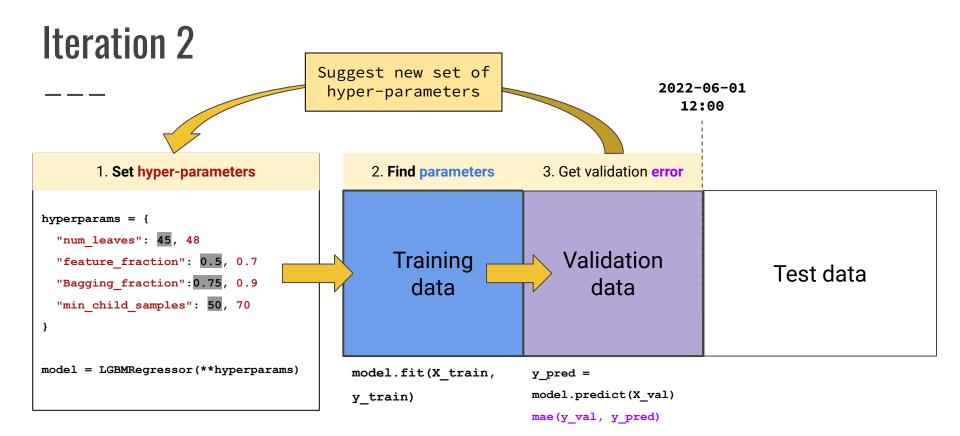
#### 3. Predict

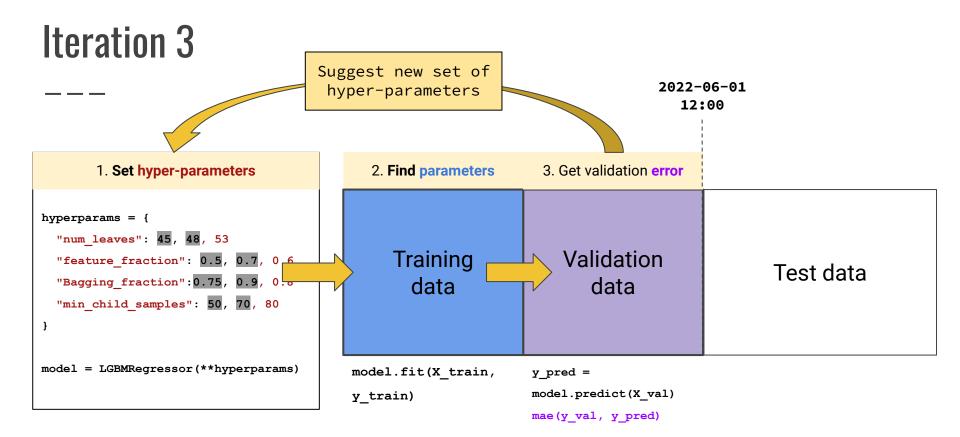
pipeline.predict(X\_test)

### Hyper-parameter tuning

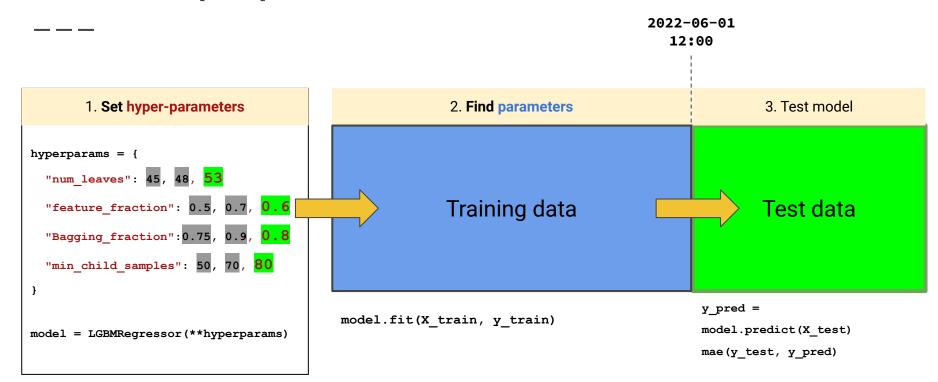








### Use best hyper-parameters, train and test



1. Set hyper-parameters

```
hyperparams = {
   "num_leaves": 45,
   "feature_fraction": 0.5,
   "bagging_fraction":0.75,
   "min_child_samples": 50,
}
model = LGBMRegressor(**hyperparams)
```

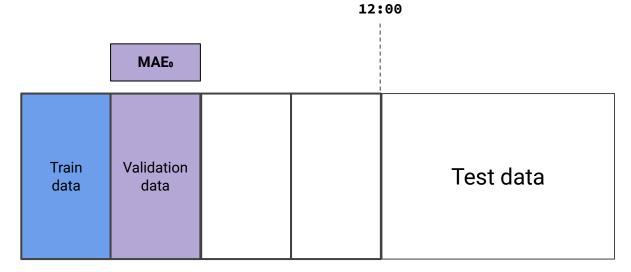
Split 1	Split 2	Split 3	Split 4	Test data

2022-06-01 12:00

```
1. Set hyper-parameters

hyperparams = {
    "num_leaves": 45,
    "feature_fraction": 0.5,
    "bagging_fraction":0.75,
    "min_child_samples": 50,
}

model = LGBMRegressor(**hyperparams)
```

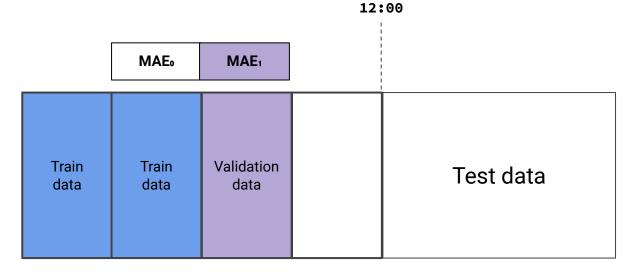


2022-06-01

1. Set hyper-parameters

hyperparams = {
 "num\_leaves": 45,
 "feature\_fraction": 0.5,
 "bagging\_fraction":0.75,
 "min\_child\_samples": 50,
}

model = LGBMRegressor(\*\*hyperparams)

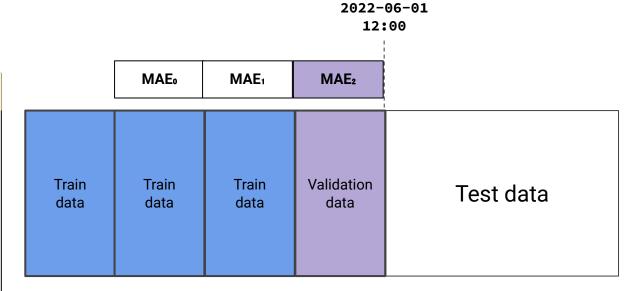


2022-06-01

```
1. Set hyper-parameters

hyperparams = {
    "num_leaves": 45,
    "feature_fraction": 0.5,
    "bagging_fraction":0.75,
    "min_child_samples": 50,
}

model = LGBMRegressor(**hyperparams)
```



Suggest new set of hyper-parameters

2022-06-01 12:00

#### 1. Set hyper-parameters

hyperparams = {

"num leaves": 45,

```
"feature_fraction": 0.5,
   "bagging_fraction":0.75,
   "min_child_samples": 50,
}
model = LGBMRegressor(**hyperparams)
```

### Average MAE<sub>0</sub>, MAE<sub>1</sub> MAE<sub>2</sub>

Split 1	Split 2	Split 3	Split 4	Test data

### Hyperparameter search with Optuna

\_\_\_

#### 1. Set hyper-parameters

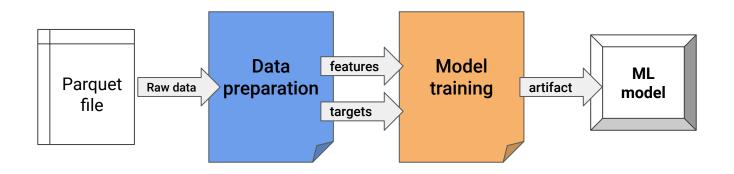
```
hyperparams = {
   "num_leaves": 45,
   "feature_fraction": 0.5,
   "bagging_fraction":0.75,
   "min_child_samples": 50,
}
model = LGBMRegressor(**hyperparams)
```

objective()

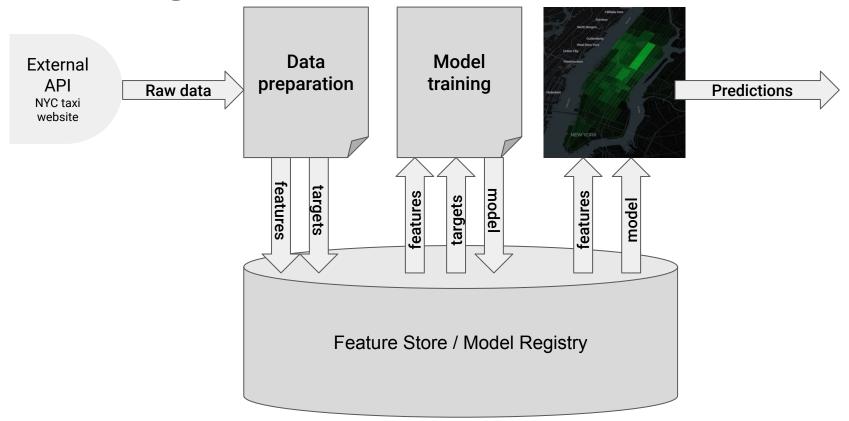
**Validation MAE** 

# 4. MLOps

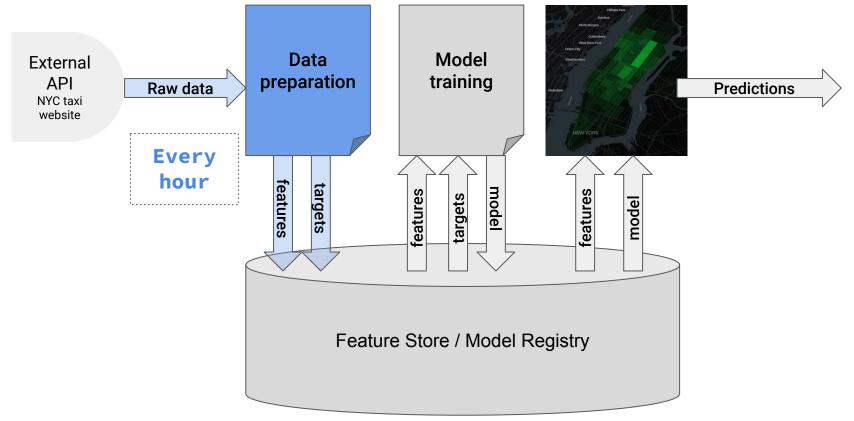
### So far we did this...



### **Batch-scoring system**



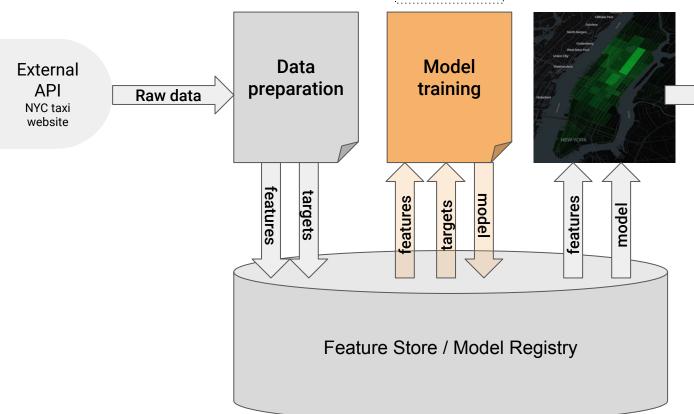
### 1. Prepare data



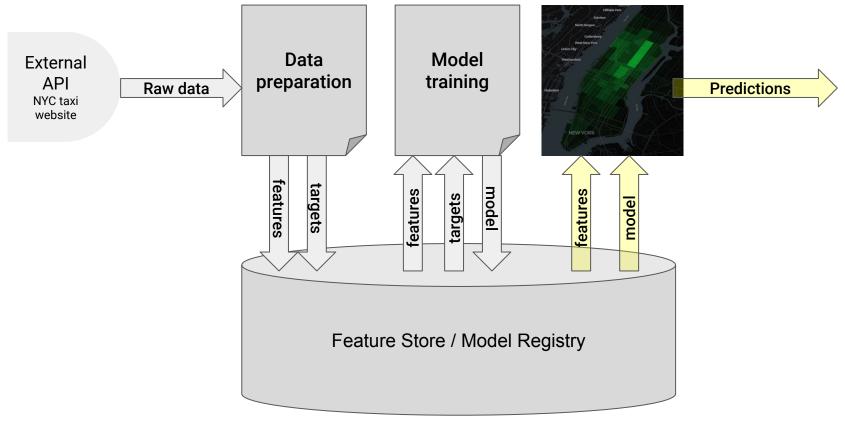
### 2. Train ML model

On demand

**Predictions** 



### 3. Generate predictions on recent data



# Serverless MLOps tools \*\*

- **Hopsworks** as our Feature Store
- **GitHub Actions** to schedule and run jobs

# What's a Feature Store?

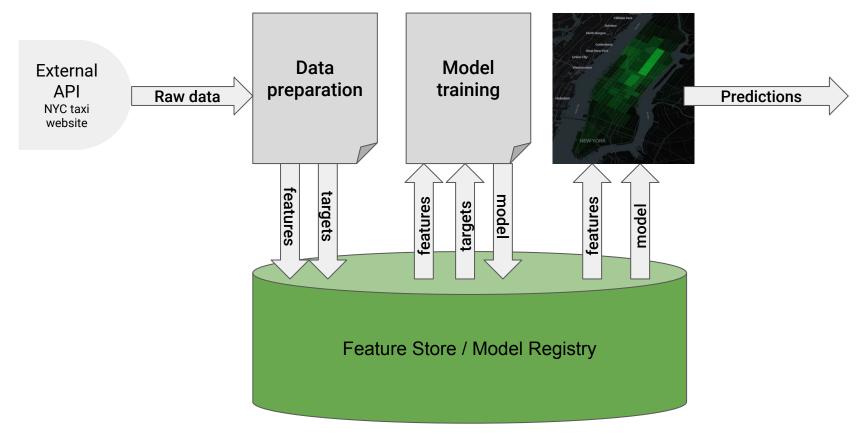
HOPSWORKS Data sources ML apps API Online A Pipelines **Applications, Services Online Apps Feature Groups Feature Views Event Bus** API WRITE API **READ API Batch Apps** Offline, Data warehouse Model Development Feature Store

HOPSWORKS Data sources ML apps API Online A Pipelines **Applications, Services Online Apps Feature Groups Feature Views** API **Event Bus** WRITE API **READ API Batch Apps** Offline, Data warehouse Model Development Feature Store

HOPSWORKS Data sources ML apps API Online A Pipelines **Applications, Services Online Apps Feature Groups Feature Views** API **Event Bus** Feature F WRITE API **READ API Batch Apps** Offline, Data warehouse Model Development Feature Store

HOPSWORKS Online API Feature Pipelines **Applications, Services Online Apps Feature Views Feature Groups** API **Event Bus WRITE API READ API Batch Apps** Offline, Data warehouse Model Development Feature Store

HOPSWORKS API Online A Pipelines **Applications, Services Online Apps Feature Views Feature Groups** API **Event Bus** Feature F **WRITE API READ API Batch Apps** Data warehouse Model Development Feature Store



### **Backfill the Feature Store**

