

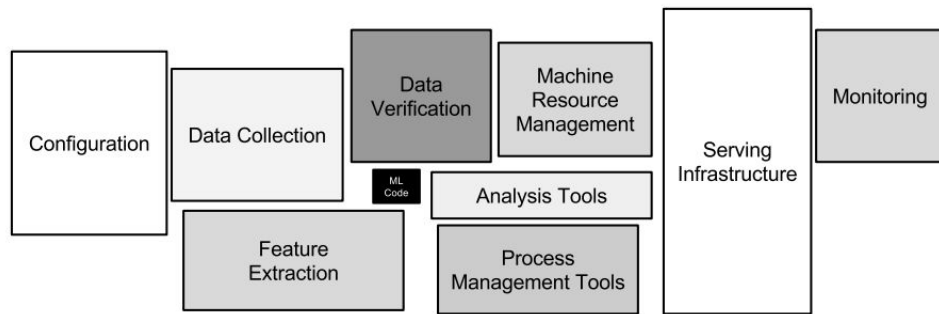
# **Machine Learning Bot for Financial Market**

**Submitted by**

**Wirachapong Suwanphibun  
Sivakorn Lerttripinyo  
Krittapasa Boontaveekul**

# Background

- Historical data in financial market is always collected.
- Collected historical prices can be used to predict a future trend.
- Machine Learning maybe able to predict by learning from data.
- Training Model is only a part of machine learning development project.



# Background: Problem

- Reliable data source
- Experiment data, such as models, hyperparameters, should be recorded
- Deploying models to the real-world application
- Changing a model version
- Underlying infrastructures to implement a system

# Background

- Objective
  - Design and implement a system going beyond the machine learning experiment.
  - Conduct an experiment on a price prediction model
    - Input
    - Output
    - Evaluation

# Outline of this Presentation

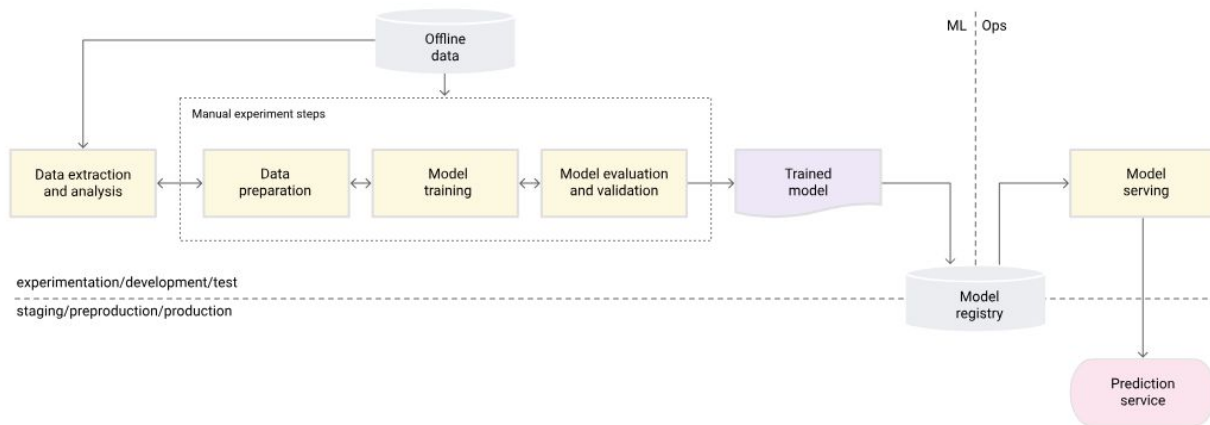
- System
  - Design
  - Experiment
  - Result and Discussion
- Model Development
  - Design
  - Experiment
  - Result and Discussion
- Conclusion
  - Assessment
  - Next Steps

# Design Requirement

- Steps for implementing ML project
  - Data Extraction
  - Data Analysis
  - Data Preparation
  - Model Training
  - Model Evaluation
  - Model Serving
  - Model Monitoring

# Design Requirement

- Applying MLOps: Use Maturity Model to evaluate
  - No standardized model, but there are some proposals by Google and Microsoft
  - Google Model: Level of Automation
  - Microsoft Model: Technical Capability
- This project starts from Google's MLOps Level 0
  - All components are run manually.



# Design Requirement

- Then, the system is improved to achieve most features from Microsoft's MLOps Level 2
  - A data pipeline automatically gathers data
  - Experiment results are tracked
  - Both the training code and the resulting models are version-controlled
  - Implementing models are heavily dependent on data scientist expertise
  - Application code has unit tests
  - Basic integration tests exist for the model



# Design Requirement

- Required Components
  - Source Code Repository
  - Model Training Infrastructure
  - Model Registry
  - ML Metadata Stores
  - Model Serving Component
  - Model Monitoring Component
  - Price Prediction Model: Be covered in the next section
    - One or more features as an input
    - A single value as an output
      - Price itself
      - Difference in price

# Infrastructure Management

- AWS and DigitalOcean
  - Free Tier for AWS
  - Free 200 USD for DigitalOcean
- Infrastructure as Code (IaC) to manage infrastructures' specification
  - Terraform
  - Backend on AWS S3 (State) and AWS DynamoDB (LockID)
    - Locking Management
    - State Management
    - Secret
  - Having additional infrastructures for managing the state is worthwhile, albeit troublesome.

```

5 resource "aws_security_group" "cap_aws_sg" {
6     name = "postgresql-rds-sg"
7
8     ingress {
9         from_port = 5432
10        to_port   = 5432
11        # from_port = 0
12        # to_port   = 65535
13        protocol   = "TCP"
14        cidr_blocks = ["0.0.0.0/0"]
15    }
16
17    egress {
18        from_port = 5432
19        to_port   = 5432
20        # from_port = 0
21        # to_port   = 65535
22        protocol   = "TCP"
23        cidr_blocks = ["0.0.0.0/0"]
24    }
25 }

```

terraform apply

-var-file="02-value.tfvars"

.tfvars is not stored in GitHub.

```

terraform > 03-main.tf
1 provider "aws" {
2     region = "ap-southeast-1"
3 }
4
5 > resource "aws_security_group" "cap_aws_sg" { ...
25 }
26
27 resource "aws_db_instance" "seniorproj_db" {
28     identifier           = "maindb-seniorproj"
29     engine               = "postgres"
30     engine_version       = "14.10"
31     instance_class       = "db.t3.micro"
32     db_name              = "seniorproj_maindb"
33     username             = var.db_username
34     password             = var.db_password
35     allocated_storage    = 20
36     storage_type         = "gp2"
37     publicly_accessible  = true
38     backup_retention_period = 2
39     skip_final_snapshot  = true
40     vpc_security_group_ids = [aws_security_group.cap_aws_sg.id]
41
42     tags = {
43         Name = "seniorproj_db"
44     }
45 }

```

# Terraform State File

```
{} terraform.json X
C: > Users > super > AppData > Local > Temp > MicrosoftEdgeDownloads > be65621f-95c9-4b6a-bbf2-2c3ec07e
7   "resources": [
76   {
81     "instances": [
82     {
84       "attributes": {
126         "license_model": "postgresql-license",
127         "listener_endpoint": [],
128         "maintenance_window": "tue:16:16-tue:16:46",
129         "manage_master_user_password": null,
130         "master_user_secret": [],
131         "master_user_secret_kms_key_id": null,
132         "max_allocated_storage": 0,
133         "monitoring_interval": 0,
134         "monitoring_role_arn": "",
135         "multi_az": false,
136         "nchar_character_set_name": "",
137         "network_type": "IPV4",
138         "option_group_name": "default:postgres-14",
139         "parameter_group_name": "default.postgres14",
140         "password": "df[REDACTED]F",
141         "performance_insights_enabled": false,
142         "performance_insights_kms_key_id": "",
143         "performance_insights_retention_period": 0,
144         "port": 5432,
```

s3/

Copy S3 URI

Objects

Properties

Objects (1) Info



Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

Show versions

< 1 > ⚙

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	terraform.tfstate	tfstate	March 21, 2024, 14:44:18 (UTC+07:00)	30.3 KB	Standard

Tables (3)

Any tag key

Any tag value

Find tables by table name

< 1 > ⚙

seniorproj-terraform-state-dynamodb

seniorproj-terraform-state-dynamodb

Autopreview

View table details

Scan or query items

Scan

Query

Select a table or index

Table - seniorproj-terraform-state-dynamodb

Select attribute projection

All attributes

Filters

Run

Reset

Items returned (1)



Actions

Create item

< 1 > ⚙

<input type="checkbox"/>	LockID (String)	Digest
<input type="checkbox"/>	seniorproj-terraform-...	2f4939c2eb86524a8e237d3183a4b480

```
terraform > 09-digitalocean.tf
15 # https://docs.digitalocean.com/reference/api/api-try-it-now/
16
17 resource "digitalocean_app" "bento_app" {
18   spec {
19     name   = "bento-app-spec-terraform"
20     region = "sgp"
21
22     env {
23       key   = "AWS_ACCESS_KEY_ID"
24       value = var.s3_access
25       type  = "SECRET"
26     }
27     env {
28       key   = "AWS_SECRET_ACCESS_KEY"
29       value = var.s3_secret_access
30       type  = "SECRET"
31     }
32   }
33 }
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

Terraform used the selected providers to generate the following execution plan. Resource actions are indicated with the following symbols:

- ~ update in-place

Terraform will perform the following actions:

```
# digitalocean_app.bento_app will be updated in-place
~ resource "digitalocean_app" "bento_app" {
  id      = "df88d3a3-1025-4eec-ab85-9860a6f21f32"
  # (7 unchanged attributes hidden)

  ~ spec {
    ~ features = [
```

## Error: Error acquiring the state lock

Error message: operation error DynamoDB: PutItem, https response error StatusCode: 400, RequestID: FS2R7PP1VTQMKC2MFU4KD0MH1VVV4KQNSO5AEMVJF66Q9ASUAAJG, ConditionalCheckFailedException: The conditional request failed

### Lock Info:










ID: aaa27854-b9ad-80ca-4c97-a1a3bbcffbfef  
Path: seniorproj-terraform-state-s3/global/s3/terraform.tfstate  
Operation: OperationTypeApply  
Who: CAQ\_DESKTOP\super@CaQ\_Desktop  
Version: 1.3.9  
Created: 2024-05-08 14:25:17.2242123 +0000 UTC  
Info:

Terraform acquires a state lock to protect the state from being written

# Database Design

- Raw Data
  - Data from BinanceAPI
  - Relational DB
    - Apply constraint easily
  - PK: (time, currency)
  - All in one table
- Derived Data
  - One currency and one set of indicators has a dedicated table.
  - Store data after performing ETL from a data pipeline
  - Relational DB
    - Apply constraint easily
  - PK: (time, currency)
- PostgreSQL on AWS RDS

# Database Design: Raw-data Database

Query   Query History   Data Output   Messages   Notifications							
        							
	currency [PK] text	time [PK] timestamp without time zone	open double precision	close double precision	low double precision	high double precision	volume double precision
1	BTCUSDT	2024-03-28 11:59:00	70637.99	70647.18	70616.41	70647.18	14.91807
2	BTCUSDT	2024-03-28 11:58:00	70673.05	70673.06	70626.33	70638	11.27007
3	BTCUSDT	2024-03-28 11:57:00	70702.01	70702.01	70673.05	70673.06	19.92578
4	BTCUSDT	2024-03-28 11:56:00	70747.01	70747.02	70700	70702	28.8656
5	BTCUSDT	2024-03-28 11:55:00	70682	70748	70682	70747.02	12.802
6	BTCUSDT	2024-03-28 11:54:00	70660	70692.06	70659.99	70682.01	22.7333
7	BTCUSDT	2024-03-28 11:53:00	70674.96	70681.99	70647.06	70660	14.35116
8	BTCUSDT	2024-03-28 11:52:00	70715.05	70715.05	70658.92	70674.96	16.10953
9	BTCUSDT	2024-03-28 11:51:00	70700	70727.93	70670.16	70715.04	28.70524
10	BTCUSDT	2024-03-28 11:50:00	70665.99	70720	70665.99	70700	28.57304
11	BTCUSDT	2024-03-28 11:49:00	70685.93	70685.94	70665.99	70666	14.85146
12	BTCUSDT	2024-03-28 11:48:00	70669.99	70701.1	70665.99	70685.94	19.95825



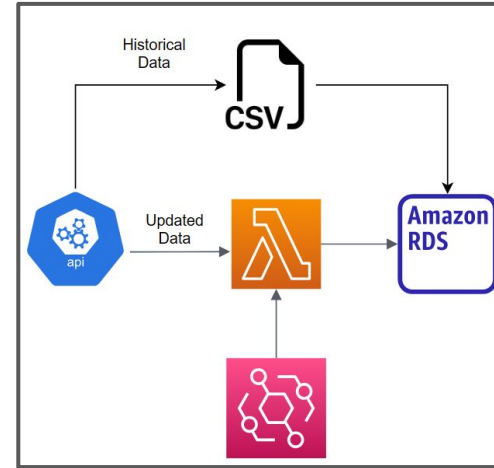
# Database Design: Derived-data Database

crypto\_ind\_adasdt  
crypto\_ind\_ethusdt  
crypto\_ind\_one

	time [PK] timestamp without time zone	currency [PK] text	close_minmax_scale double precision	close double precision	ma25_99h double precision	ma7_25h double precision	ma7_25d double precision	ma25_99h_scale double precision	ma7_25h_scale double precision	ma7_25d_scale double precision
1	2024-05-13 23:59:00	ETHUSDT	0.2227577896	2951.05	-2.862040404	6.6394285714	-122.3024571429	-0.000969838	0.002249853	-0.0414437089
2	2024-05-13 23:58:00	ETHUSDT	0.2212467244	2950.26	-3.1062505051	6.3083428571	-122.4149142857	-0.0010528735	0.0021382329	-0.0414929241
3	2024-05-13 23:57:00	ETHUSDT	0.2212467244	2950.26	-2.9413373737	6.0204571429	-122.7217142857	-0.0009969756	0.0020406531	-0.0415969149
4	2024-05-13 23:56:00	ETHUSDT	0.2192192192	2949.2	-2.8162666667	5.9367428571	-122.2489714286	-0.0009549256	0.0020130011	-0.0414515704
5	2024-05-13 23:55:00	ETHUSDT	0.2192000918	2949.19	-2.923830303	6.0408	-122.7245142857	-0.0009914011	0.0020482912	-0.0416129562
6	2024-05-13 23:54:00	ETHUSDT	0.2189705629	2949.07	-3.1964363636	5.8499428571	-122.0528571429	-0.0010838794	0.0019836568	-0.0413868973
7	2024-05-13 23:53:00	ETHUSDT	0.2176890266	2948.4	-3.1386909091	6.2338857143	-122.2893142857	-0.0010645404	0.0021143284	-0.0414765006
8	2024-05-13 23:52:00	ETHUSDT	0.2170004399	2948.04	-3.2212323232	6.0037142857	-122.7628	-0.0010926691	0.0020365105	-0.0416421758
9	2024-05-13 23:51:00	ETHUSDT	0.2175933896	2948.35	-3.5102626263	6.8582857143	-122.6173714286	-0.0011905855	0.0023261437	-0.041588472
10	2024-05-13 23:50:00	ETHUSDT	0.2169420577	2948.01	-3.0408565656	7.524	-122.5891428571	-0.0012657005	0.0025005425	-0.0416176149

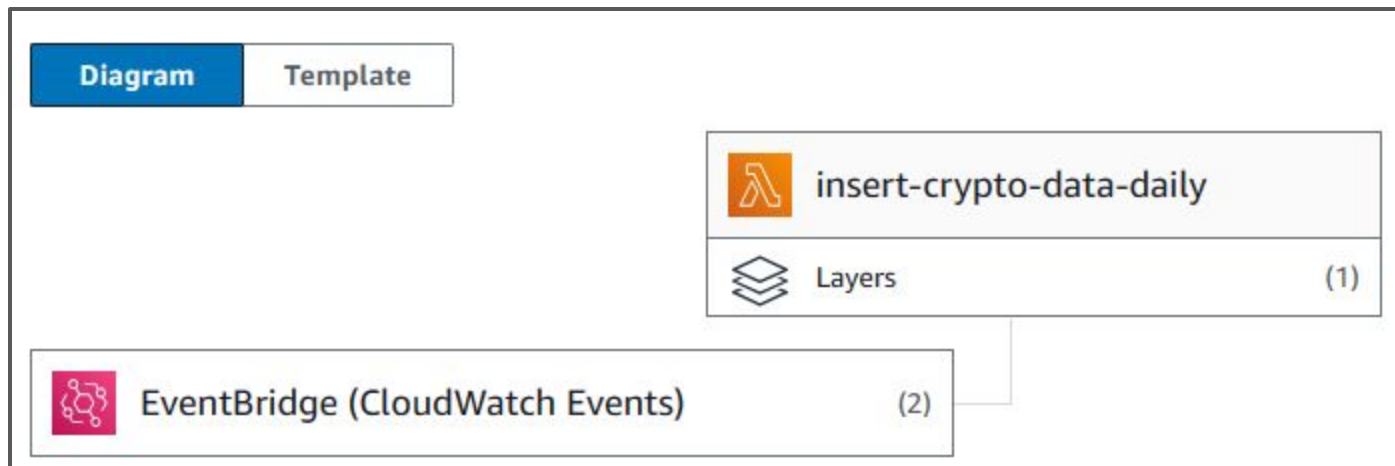
# Data Ingestion - Design

- Historical Data
  - Insert manually
    - Explore and Clean Data first
    - One-time Operation + Large data size
      - Bulk Insertion is more efficient
    - Find Pattern
      - Create Automation for further new-generated data
- New-Generated Data
  - Insert automatically
  - AWS Lambda + AWS EventBridge Rule



# Data Ingestion - Usage

- 1440 rows per day per currency (1440 minutes)
- Limitation
  - Limit 1000 rows per query for BinanceAPI
    - Solution: 2 queries per day
      - 0:00 and 12:00 UTC+0
      - 720 rows per query
      - Can increase frequency(e.g. 3 per day) to update data more real-time
  - No Available Libraries in Lambda's Python Runtime
    - Pack libraries as .zip and include as layers
    - Libraries must be downloaded in the same Python version as Lambda Runtime Env

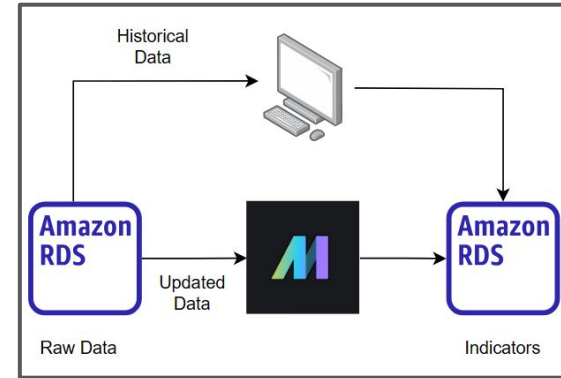


# AWS Lambda Limitation

- Tough to debug and test
  - Not suitable for implementing a complex operation.
- Tough to manage dependencies
  - Although having layers, this method is time-consuming and error-prone
  - Layer having size > 50 MB must be separated to two smaller layers (AWS Limitation)
- More appropriate tool (e.g. data pipeline tool) should be used

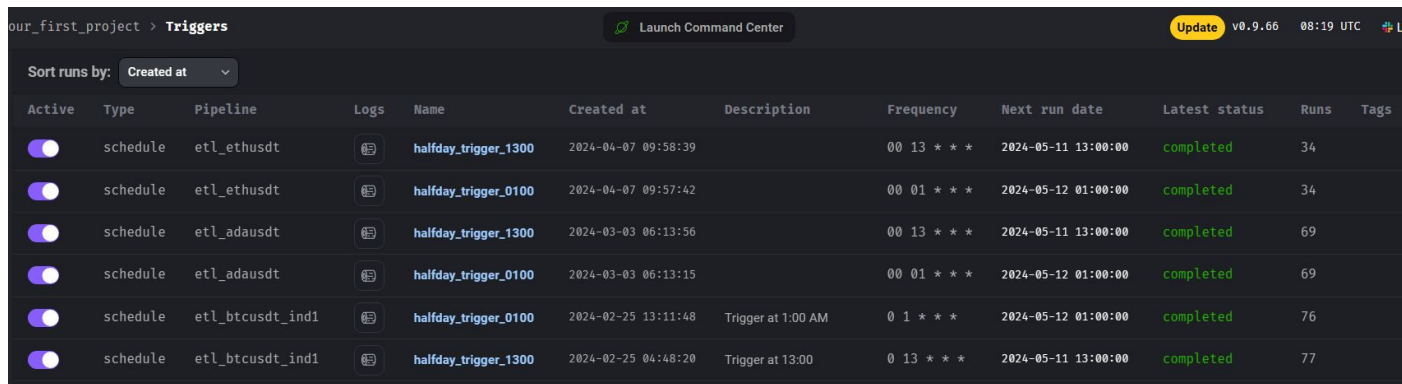
# Data Transformation - Design

- Historical Data
  - Insert Manually
  - One-time Operation + Large amount of data
    - Not worthwhile to automate this
- New-generated Data
  - Automated ETL
  - Mage.ai (Data Pipeline Tool) is applied.
    - Apache Airflow requires at least 4GB RAM
    - Has less functions but more than enough for creating data pipelines



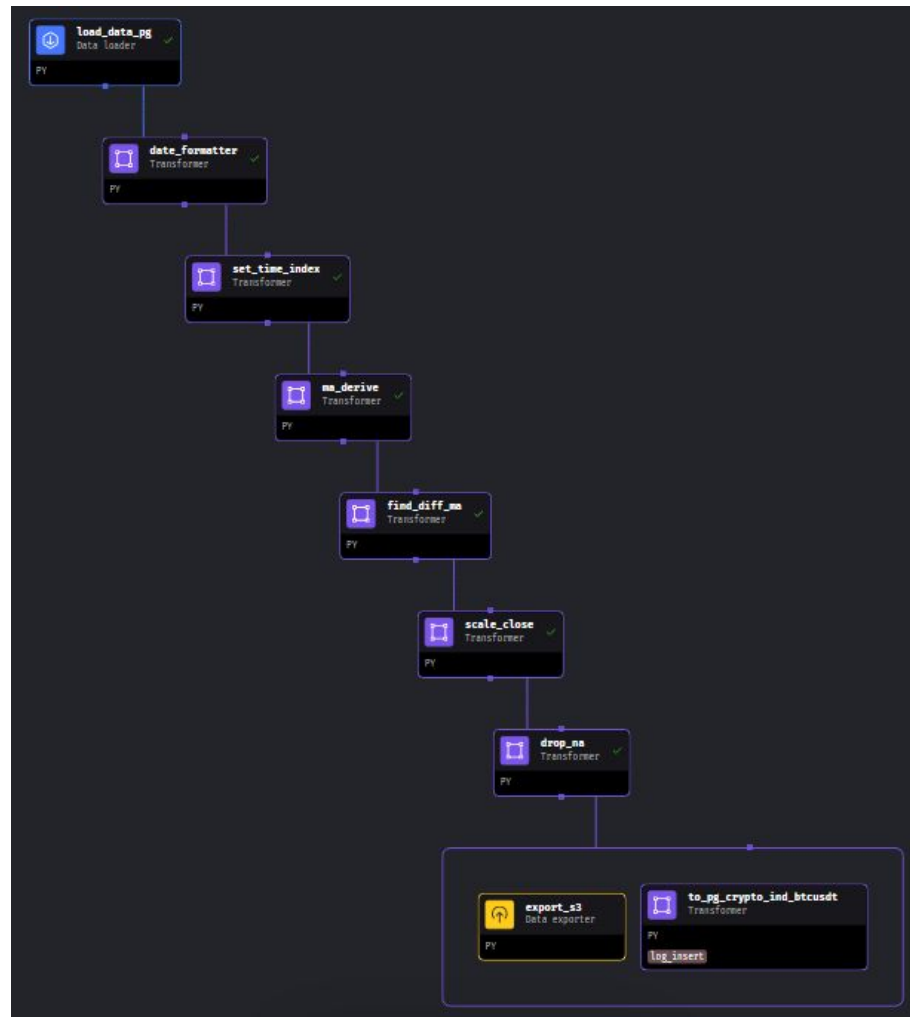
# Data Transformation - Usage

- Extract data from raw-data database
- ETL is firstly implemented in Google Colab
  - Pandas is enough for approx 3M rows
  - Test with **Historical Data**
- Then, it is adapted as a data pipeline in Mage.ai
  - Directed Acyclic Graph (DAG)
  - Use the pipelines to transform **new-generated data**
- Schedules are configured to automatically trigger the pipeline to run

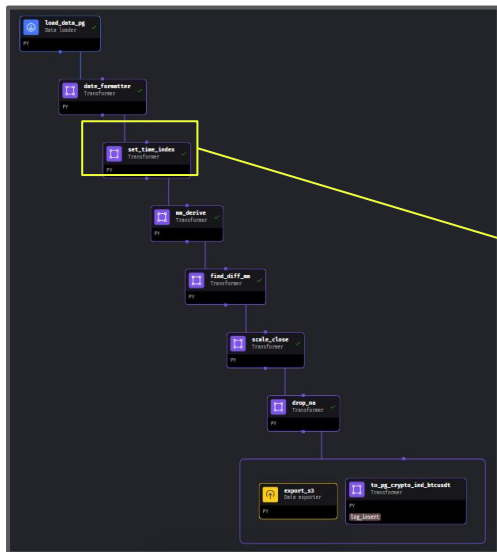


The screenshot displays the 'Triggers' section of the Mage.ai interface. At the top, there's a header with 'our\_first\_project > Triggers', a 'Launch Command Center' button, and an 'Update' button next to version 'v0.9.66' and the time '08:19 UTC'. Below the header, a 'Sort runs by:' dropdown is set to 'Created at'. The main area is a table with columns: Active, Type, Pipeline, Logs, Name, Created at, Description, Frequency, Next run date, Latest status, Runs, and Tags. There are six rows of scheduled triggers, all with 'completed' status.

Active	Type	Pipeline	Logs	Name	Created at	Description	Frequency	Next run date	Latest status	Runs	Tags
<input checked="" type="checkbox"/>	schedule	etl_ethusdt		halfday_trigger_1300	2024-04-07 09:58:39		00 13 * * *	2024-05-11 13:00:00	completed	34	
<input checked="" type="checkbox"/>	schedule	etl_ethusdt		halfday_trigger_0100	2024-04-07 09:57:42		00 01 * * *	2024-05-12 01:00:00	completed	34	
<input checked="" type="checkbox"/>	schedule	etl_adausdt		halfday_trigger_1300	2024-03-03 06:13:56		00 13 * * *	2024-05-11 13:00:00	completed	69	
<input checked="" type="checkbox"/>	schedule	etl_adausdt		halfday_trigger_0100	2024-03-03 06:13:15		00 01 * * *	2024-05-12 01:00:00	completed	69	
<input checked="" type="checkbox"/>	schedule	etl_btcsdt_ind1		halfday_trigger_0100	2024-02-25 13:11:48	Trigger at 1:00 AM	0 1 * * *	2024-05-12 01:00:00	completed	76	
<input checked="" type="checkbox"/>	schedule	etl_btcsdt_ind1		halfday_trigger_1300	2024-02-25 04:48:20	Trigger at 13:00	0 13 * * *	2024-05-11 13:00:00	completed	77	







your\_first\_project > Pipelines > etl\_btusdft\_ind1 > Edit

Launch Command Center

ALL FILES CURRENT BLOCK <

Search files

your\_first\_project

- > callbacks
- > charts
- > custom
- > data\_exporters
- > data\_loaders
- > dbt
- > extensions
- > interactions
- > pipelines
- > scratchpads
- > transformers
- > utils
- > \_\_init\_\_.py
- > io\_config.yaml
- > metadata.yaml
- > requirements.txt

File Edit Run

PY TRANSFORMER date\_formatter 1 parent 3 pipelines

PY TRANSFORMER set\_time\_index 1 parent 2 pipelines

Positional arguments for decorated function:

```
@transformer
def transform(data):
    data > date_formatter

    if 'transformer' not in globals():
        from mage_ai.data_preparation.decorators import transformer
    if 'test' not in globals():
        from mage_ai.data_preparation.decorators import test

@transformer
def transform(df, *args, **kwargs):
    """
    Template code for a transformer block.

    Add more parameters to this function if this block has multiple parent blocks.
    There should be one parameter for each output variable from each parent block.

    Args:
        data: The output from the upstream parent block
        args: The output from any additional upstream blocks (if applicable)

    Returns:
        Anything (e.g. data frame, dictionary, array, int, str, etc.)
    """
    # Specify your transformation logic here
    df = df.iloc[:: -1].reset_index(drop=True)
    df.set_index('time', inplace=True)
    df = df.sort_values(by='time')

    return df
```

# Training Management System

- Track experiments, Record relevant information, and Be a centralized source for distributing models
- MLFlow: Hosted on EC2
  - Metadata: PostgreSQL on AWS RDS
  - Artifact and Object: AWS S3 Bucket
  - Prevent data loss from the instance failure
  - IAM is required for accessing data in DB and S3
    - “aws configure” to enter access and secret access key
    - Attach “Roles” directly to EC2
- Training code is modified to utilize MLFlow to track the experiment.

## model/

Objects

Properties

### Objects (5) Info



Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly permissions. [Learn more](#)

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	conda.yaml	yaml	February 26, 2024, 14:05:05 (UTC+07:00)	239.0 B	Standard
<input type="checkbox"/>	MLmodel	-	February 26, 2024, 14:05:06 (UTC+07:00)	515.0 B	Standard
<input type="checkbox"/>	model.xgb	xgb	February 26, 2024, 14:05:06 (UTC+07:00)	11.5 KB	Standard

senioproj

- Databases (4)
- mlflow\_test
  - Casts
  - Catalogs
  - Event Triggers
  - Extensions
  - Foreign Data Wrappers
  - Languages
  - Publications
  - Schemas (1)
    - public
      - Aggregates
      - Collations
      - Domains
      - FTS Configurations
      - FTS Dictionaries
      - FTS Parsers
      - FTS Templates
      - Foreign Tables
      - Functions
      - Materialized Views
      - Operators
      - Procedures
      - Sequences
      - Tables (16)
        - alembic\_version
        - datasets
        - experiment\_tags
        - experiments
        - input\_tags
        - inputs
        - latest\_metrics

public.experiments/mlflow\_test/senioproj\_kimjongun/senioproj...

100 rows

Query Query History Data Output Messages Notifications

experiment_id	name	artifact_location	lifecycle_stage	creation_time	last_update_time
[PK] integer	character varying (256)	character varying (256)	character varying (32)	bigint	bigint
1	0 Default	s3://mlflow-bucket-test-1012-senioproj/0	active	169953685865	169953685865
2	1 nyc-taxi-experiment	s3://mlflow-bucket-test-1012-senioproj/1	active	17057447377146	17057447377146
3	2 nyc-taxi-experiment-2	s3://mlflow-bucket-test-1012-senioproj/2	active	1705746420871	1705746420871
4	3 test-crypto	s3://mlflow-bucket-test-1012-senioproj/3	deleted	1706794347828	1706793074084
5	4 Bitcoin_Price_Prediction_Minute	s3://mlflow-bucket-test-1012-senioproj/4	active	1706793272632	1706793272632
6	5 Bitcoin_Price_Prediction_GPU	s3://mlflow-bucket-test-1012-senioproj/5	active	1706858019298	1706858019298
7	6 XGBoost_Experiment	s3://mlflow-bucket-test-1012-senioproj/6	active	1708164304511	1708164304511
8	11 BTCUSD_Price_Prediction_next24hrs_Regression_HN_ind	s3://mlflow-bucket-test-1012-senioproj/11	active	1708929418074	1710228527692
9	12 ADAUSD_Price_Prediction_next24hrs_Regression_HN_ind	s3://mlflow-bucket-test-1012-senioproj/12	active	1709450735762	1710228536053
10	13 BTCUSD_Price_Prediction_Minute_Classification_MS	s3://mlflow-bucket-test-1012-senioproj/13	active	1710046000517	1710228490074
11	14 BTCUSD_Price_Prediction_Minute_Regression_MS	s3://mlflow-bucket-test-1012-senioproj/14	active	1710046072113	1710228504094
12	15 MainTradingBotXGBoostExperimentBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/15	deleted	1711618248346	1711618597798
13	18 MainTradingBotXGBoostExperimentBigDataVersion777	s3://mlflow-bucket-test-1012-senioproj/18	active	1711666366148	1711666366148
14	19 MainTradingBotLASSOBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/19	active	1711666383089	1711666383089
15	20 MainTradingBotHuberRegressorBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/20	active	1711666456956	1711666456956
16	21 MainTradingBotBayesianRidgeBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/21	active	1711666476760	1711666476760
17	22 MainTradingBotRidgeTUNEDBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/22	active	1711666489014	1711666489014
18	23 MainTradingBotLASSOARSTUNEDBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/23	active	1711666500628	1711666500628
19	24 MainTradingBotELASTICNETBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/24	active	1711666512905	1711666512905
20	25 MainTradingBotCatBoostRegressorBigDataVersion	s3://mlflow-bucket-test-1012-senioproj/25	active	1711666633272	1711666633272
21	26 MainTradingBotTradingAdditionalRound2	s3://mlflow-bucket-test-1012-senioproj/26	active	1711666833507	1711666833507
22	27 MainTradingBotTradingAdditionalRound3	s3://mlflow-bucket-test-1012-senioproj/27	active	1711667053413	1711667053413
23	28 MainTradingBotTradingAdditionalRound4	s3://mlflow-bucket-test-1012-senioproj/28	active	1711667288996	1711667288996
24	29 MainTradingBotTradingAdditionalRound5	s3://mlflow-bucket-test-1012-senioproj/29	active	1711667548555	1711667548555
25	30 MainTradingBotTradingAdditionalRound6	s3://mlflow-bucket-test-1012-senioproj/30	active	1711667768450	1711667768450

Experiment Name

```
1 TRACKING_SERVER_HOST = "..."  
2 mlflow.set_tracking_uri(f"http://{TRACKING_SERVER_HOST}<port>")  
3 mlflow.set_experiment("...")  
4 with mlflow.start_run(run_name="..."):  
5     params = {  
6         "colsample_bytree": 0.3,  
7         "learning_rate": 0.1,  
8         <more params can be added>  
9     }  
10    mlflow.log_params(params)  
11    xg_reg = xgb.XGBRegressor(**params)  
12    xg_reg.fit(X_train, y_train)  
13    y_test_pred = xg_reg.predict(X_test)  
14    test_mse = mean_squared_error(y_test, y_test_pred)  
15    mlflow.log_metric("test_mse", test_mse)  
16    mlflow.xgboost.log_model(xg_reg, "model")
```

Original  
Training Code

mlflow 2.10.2 Experiments Models

**Experiments** (+) (-)

Search Experiments

- ☐ Default
- ☐ nyc-taxi-experiment
- ☐ nyc-taxi-experiment-2
- ☐ Bitcoin\_Price\_Prediction\_Minute
- ☐ Bitcoin\_Price\_Prediction\_GPU
- ☐ XGBoost\_Experiment
- ☒ BTCUSDT\_Price\_Prediction\_next24hrs...
- ☐ ADAUSDT\_Price\_Prediction\_next24hrs...
- ☐ BTCUSDT\_Price\_Prediction\_Minute...
- ☐ BTCUSDT\_Price\_Prediction\_Minute...

**BTCUSDT\_Price\_Prediction\_next24hrs\_Regression\_HN\_ind** Provide Feedback Share

Experiment ID: 11 Artifact Location: s3://mlflow-bucket-test-1012-seniorproj/11

> Description Edit

Q metrics.rmse < 1 and params.model = "tree" Time created State: Active Datasets Sort: Created

Columns

Table Chart Evaluation **Experimental**

	Run Name	Created	Dataset	Duration	Source	Models
<input type="checkbox"/>	BTCUSDT_1	27 days ago	-	20.2s	colab_ker...	xg/3
<input type="checkbox"/>	gentle-jay-332	27 days ago	-	15.5s	colab_ker...	-
<input type="checkbox"/>	awesome-skunk-129	27 days ago	-	7.1min	colab_ker...	-

Experiment Names

Run names

mlflow 2.10.2 Experiments Models

BTCUSDT\_Price\_Prediction\_next24hrs\_Regression\_HN\_ind

**BTCUSDT\_1**

Run ID: ccc08bae3a54fe1ba002cdf3506bd25 Date: 2024-02-26 14:04:46 Source: colab\_kernel\_launcher.py User: root

Duration: 20.2s Status: FINISHED Lifecycle Stage: active

> Description Edit

predict next 24 hrs using XGBoost

> Datasets

> Parameters (11)

> Metrics (1)

> Tags

> Artifacts

Full Path: s3://mlflow-bucket-test-1012-seniorproj/11/ccc08bae3a54fe1ba002cdf3506bd25/artifacts/model

MLflow Model

The code snippets below demonstrate how to make predictions using the logged model. This model is also registered to the [model registry](#).

**Model schema**

Input and output schema for your model. [Learn more](#)

**Make Predictions**

Predict on a Spark DataFrame:

mlflow 2.10.2 Experiments Models

BTCUSDT\_Price\_Prediction\_next24hrs\_Regression\_HH\_ind >  
BTCUSDT\_1

Run ID: ccc0d8aee3a54fe1bb002cdf3506bd25 Date: 2024-02-26 14:04:46 Source: colab\_kernel\_launcher.py User: root  
Duration: 20.2s Status: FINISHED Lifecycle Stage: active

▼ Description Edit

predict next 24 hrs using XGBoost

> Datasets

> Parameters (11)

> Metrics (1)

> Tags

▼ Artifacts

model

- MLmodel
- conda.yaml
- model.xgb
- python\_env.yaml
- requirements.txt

Full Path: s3://mlflow-bucket-test-1012-seniorpro/11/ccc0d8aee3a54fe1bb002cdf3506bd25/artifacts/model

**MLflow Model**

The code snippets below demonstrate how to make predictions using the logged model. This model is also registered to the [model registry](#).

**Model schema**

Input and output schema for your model. [Learn more](#)

**Make Predictions**

Predict on a Spark Dataframe: [Link to notebook](#)

**BTCUSDT\_1**

Run ID: ccc0d8aee3a54fe1bb002cdf3506bd25 Date: 2024-02-26 14:04:46

Status: FINISHED Lifecycle Stage: active

▼ Description Edit

predict next 24 hrs using XGBoost

> Datasets

▼ Parameters (11)

Name	Value
alpha	10
colsample_bytree	0.3
currency	BTCUSDT
end-datetime	2024-02-24 23:59:00
feature	close_minmax_scale, ma7_25h_scale, ma25_99h_scale, ma7_25d_scale
learning_rate	0.1
max_depth	5
n_estimators	10
objective	reg:squarederror
prediction	result (%up/down in next 24 hrs)
start-datetime	2017-09-12 11:59:00

# Model Monitoring

- Continuously monitor ML model
- Jupyter Notebook is inconvenient for opening and sharing
- Webpage is better
  - Easier to share a link and open in any supported browser
  - Streamlit
    - Python-based: Easy to integrate with Matplotlib

Choose a table

crypto\_ind\_one

Choose a start date

2024/01/01

Choose a start time

00:00

Choose an end date

2024/01/31

Choose an end time

23:45

Choose a model:

CBR

Enter

## Visualization of Predicted Values Over Time

Start Time: 2024-01-01 00:00:00

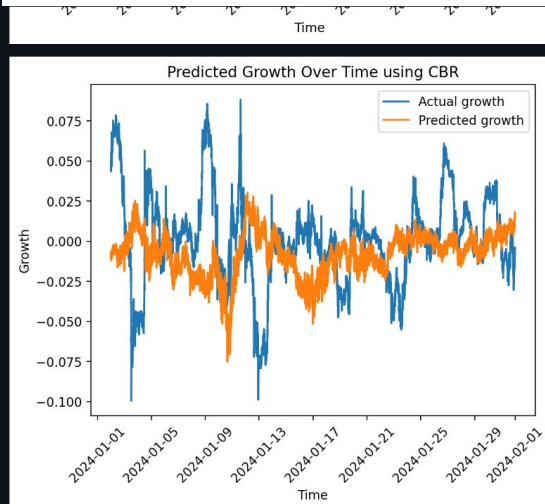
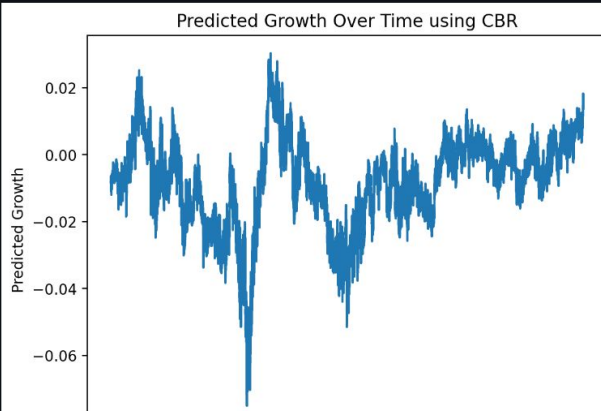
CBR

Enter

## Visualization of Predicted Values Over Time

Start Time: 2024-01-01 00:00:00

End Time: 2024-01-31 23:45:00



precision	recall	f1-score	support
-----------	--------	----------	---------

0	0.43	0.65	0.52	28444
---	------	------	------	-------

1	0.41	0.22	0.28	22742
---	------	------	------	-------

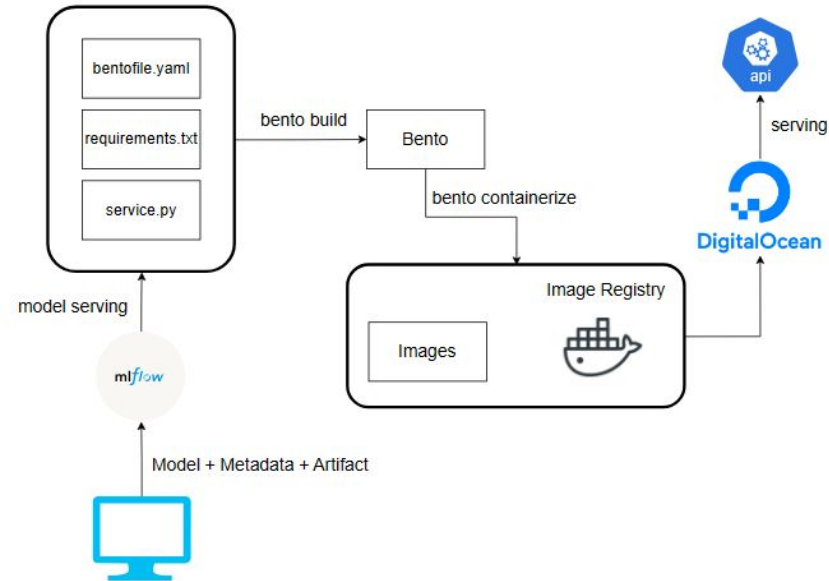
accuracy			0.42	43186
----------	--	--	------	-------

Model results will be discussed later in the next section.



# Model Building and Serving

- Bridging between the model experiment and the real-world application
- Building an API as an image and Deploying the image



InferenceAPI	Input	Output
POST <code>/predict_xgboost</code>	JSON	NumpyNdarray
POST <code>/predict</code>	JSON	NumpyNdarray

**Help**

- [Documentation](#): Learn how to use BentoML.
- [Community](#): Join the BentoML Slack community.
- [GitHub Issues](#): Report bugs and feature requests.
- Tip: you can also [customize this README](#).

[Contact BentoML Team](#)

Servers

**Service APIs** BentoML Service API endpoints for inference.

POST	<code>/predict_xgboost</code>	InferenceAPI(JSON → NumpyNdarray)
POST	<code>/predict</code>	InferenceAPI(JSON → NumpyNdarray)

Code Issues Pull requests **Actions** Projects Security Insights Settings

Debug request body #15 Re-run all jobs ...

Summary

Jobs

- build\_and\_push

Run details

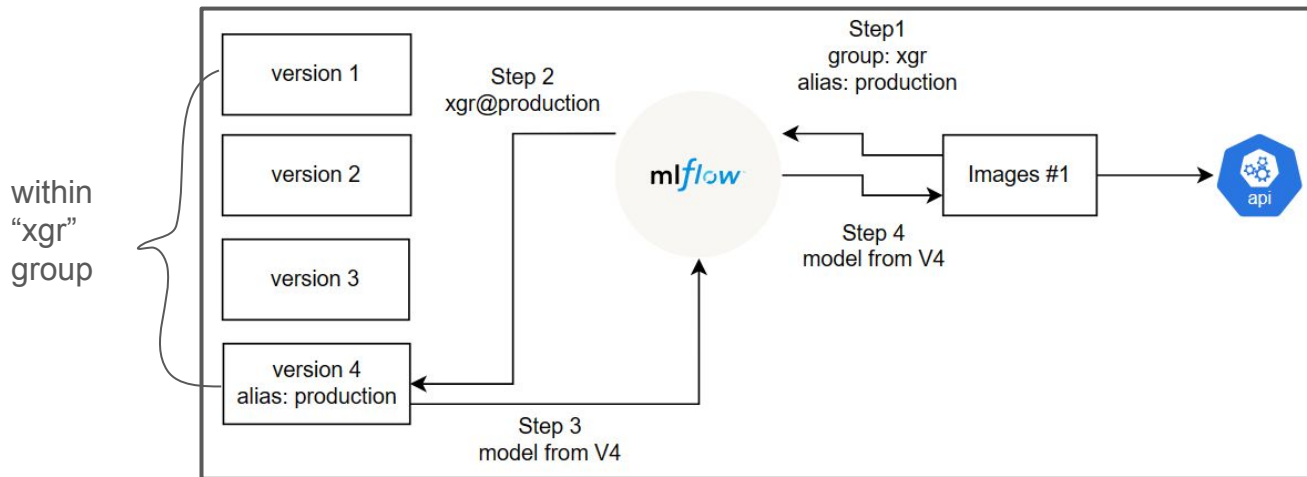
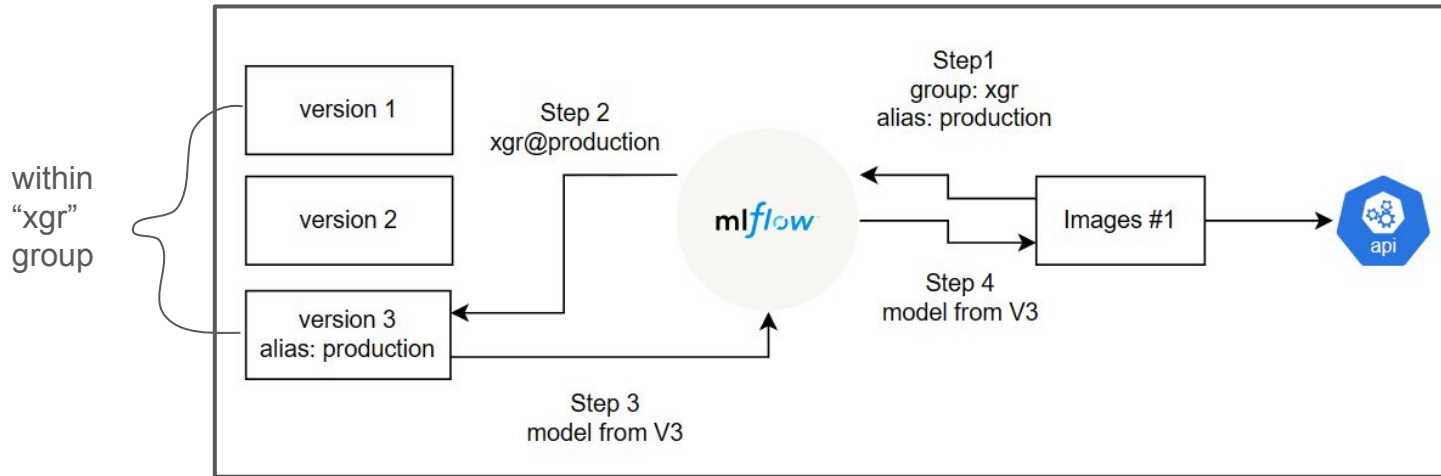
- Usage
- Workflow file

Triggered via push 2 days ago

MATT138 pushed → f70d64a main **Success** Total duration: 6m 41s Billable time: 7m Artifacts: -

cd.yaml  
on: push

build\_and\_push 6m 29s



# Experiment Setup

- A software testing concept *is adapted* to test this system.
  - A unit test tests whether a particular system functions properly.
  - A component test will test whether an interaction between two systems work properly.
  - A system test is applied to verify that all systems can work together properly.
    - An integration test is skipped as our system has a small scale, and the borderline between the system and integration test is blurred

# Experiment Setup: Infrastructure List

- AWS
  - S3 Bucket
  - Lambda (Python 3.9 as Runtime Environment) with AWS EventBridge Rules
  - EC2: t2.micro, 1vCPU, RAM 1GB
  - RDS PostgreSQL: t3.micro, 2vCPU, RAM 1 GB, gp2-20GiB Storage
  - DynamoDB: Pay per Request
- DigitalOcean
  - Droplet: RAM 4GB, 2vCPU, 80GB Storage
  - App Platform: RAM 2GB, 1vCPU
  - Container Registry: 5GB Storage

# Experimental Method: Unit Test

- IaC
  - The scripts shall provision defined infrastructures with a defined specification.
  - The state after running the script shall be updated and stored in S3 bucket.
- Data Ingestion
  - Raw-data database shall be created.
  - Historical data shall be able to be inserted manually to the raw-data database.
  - New data shall be automatically retrieved from an external API and stored in the raw-data database at the time specified by triggers.
  - No duplicated data shall exist in the raw-data database.

# Experimental Method: Unit Test

- Data Transformation
  - Derived-data database shall be created.
  - Historical data shall be able to be derived and inserted into the derived-data database manually.
  - The data pipeline shall be triggered at a specified time.
  - The data pipeline shall successfully retrieve data from a defined data source, derive new indicators, and store them in the derived-data database.
- Training Management
  - The system shall be able to create several experiments.
  - Each experiment shall be able to contain several runs.
  - Each run shall contain logged data, such as metrics and hyperparameters, and model.
  - The AWS access key and secret access key shall be required to log the experiment result from the local to the system.

# Experimental Method: Unit Test

- Model Monitoring
  - The dashboard hosted on the website shall be able to visualize the model metrics on the selected currency and date range.
  - For this test, the hard-coded model and dataset can be used.
- Model Building and Serving
  - After pushing the code into the main branch, the GitHub shall be automatically activated, and the new built image shall be automatically deployed through the API.



# Experimental Method: Component Test

- **The data ingestion system** shall always ingest new data and insert into the raw-data database. **The data transformation system** shall be able to set the raw-data database as the data source, derive new indicators, and insert them into the derived-data database.
- **The model building and serving system** shall be able to retrieve the model from **the training management system** and use it for serving. If there is a change in the model version for serving, the new version shall be able to be deployed and applied without rebuilding the image.

# Experimental Method: Component and System Test

- **The model monitoring system** shall be able to retrieve data from **derived-data database** and the model from **the training management system**. They shall be used to generate a metric and graph to show the model performance.
- System test
  - All components shall be able to work together properly.
  - The overall system shall be run and periodically check whether there is any error occur

# Non-Functional Requirements Testing

- Operational Requirements
  - The system shall provide a result of at least 99% uptime.
  - The system databases shall have 99% uptime.
  - The system shall backup all databases to prevent data loss.
- Performance Requirements
  - The system shall store daily data every day.
  - The system shall provide a response within 5 seconds after making a request.
- Security and Privacy
  - The system shall scrape data from a legal source only.
  - The system shall allow only team members to access databases.

# Test Result

Topic	Result (Pass/Partial Pass/Fail)
Infrastructure as Code	Pass
Data ingestion system	Pass
Data transformation system	Pass
Training management system	Pass
Model monitoring system	Pass
Model building and serving system	Pass

Table 2: Unit Test Result

Topic	Result (Pass/Partial Pass/Fail)
Data Ingestion + Data Transformation	Pass
Building and Serving + Management System	Pass
Model monitoring + Management System	Pass

Table 3: Component Test Result

## /aws/lambda/extract\_func

Actions ▼

View in Logs Insights

Start tailing

Search log group

## ► Log group details

Log streams

Tags

Anomaly detection

Metric filters

Subscription filters

Contributor Insights

Data protection

## Log streams (11)



Delete

Create log stream

Search all log streams



Exact match



Show expired

[Info](#)

1



Log stream



Last event time

[2024/05/06/\[\\$LATEST\]961bd54775654a8ba137f12704e418b0](#)

2024-05-06 07:00:58 (UTC+07:00)

[2024/05/05/\[\\$LATEST\]098eef322f794cfe9ef6e091647090cc](#)

2024-05-05 19:00:19 (UTC+07:00)

[2024/05/05/\[\\$LATEST\]49413c01345c4f74b7ff9b39d93f9c48](#)

2024-05-05 07:00:57 (UTC+07:00)

[2024/05/04/\[\\$LATEST\]28c32035213d4f59a3b52c7eadb039db](#)

2024-05-04 19:00:18 (UTC+07:00)

## Log events



Actions ▼

Start tailing

Create metric filter

You can use the filter bar below to search for and match terms, phrases, or values in your log events. [Learn more about filter patterns](#)

1m

1h



Local timezone ▼

Display ▼



Timestamp

Message

No older events at this moment. [Retry](#)

2024-05-06T07:00:50.075+07:00

INIT\_START Runtime Version: python:3.9.v50 Runtime Version ARN: arn:aws:lambda:ap-southeast-1::runti...



2024-05-06T07:00:51.215+07:00

START RequestId: dfb12c7c-b0d8-4fc7-9a3e-8cb09e1d61cb Version: \$LATEST



2024-05-06T07:00:58.129+07:00

END RequestId: dfb12c7c-b0d8-4fc7-9a3e-8cb09e1d61cb



2024-05-06T07:00:58.129+07:00

REPORT RequestId: dfb12c7c-b0d8-4fc7-9a3e-8cb09e1d61cb Duration: 6914.40 ms Billed Duration: 6915 ms...

No newer events at this moment. Auto retry paused. [Resume](#)



your\_first\_project &gt; Pipeline runs

[Launch Command Center](#)[Update](#)

v0.9.66

04:43 UTC

[Live help](#)

Filter



Status	Logs	ID	Block runs	Pipeline	Trigger	Pipeline tags	Execution date	Started at	Completed at	Execution time
✓ Done		331	8 / 8	etl_ethusdt	halfday_trigger_0100		2024-05-06T01:00:00	2024-05-06T01:00:04	2024-05-06T01:01:21	00:01:17.17
✓ Done		330	9 / 9	etl_btcsdt_ind1	halfday_trigger_0100		2024-05-06T01:00:00	2024-05-06T01:00:03	2024-05-06T01:01:13	00:01:09.84
✓ Done		329	8 / 8	etl_adausdt	halfday_trigger_0100		2024-05-06T01:00:00	2024-05-06T01:00:02	2024-05-06T01:01:11	00:01:09.35
✓ Done		328	8 / 8	etl_ethusdt	halfday_trigger_1300		2024-05-05T13:00:00	2024-05-05T13:00:09	2024-05-05T13:01:31	00:01:21.68
✓ Done		327	9 / 9	etl_btcsdt_ind1	halfday_trigger_1300		2024-05-05T13:00:00	2024-05-05T13:00:09	2024-05-05T13:01:30	00:01:20.95
✓ Done		326	8 / 8	etl_adausdt	halfday_trigger_1300		2024-05-05T13:00:00	2024-05-05T13:00:08	2024-05-05T13:01:28	00:01:20.33
✓ Done		325	8 / 8	etl_ethusdt	halfday_trigger_0100		2024-05-05T01:00:00	2024-05-05T01:00:08	2024-05-05T01:01:28	00:01:20.58
✓ Done		324	9 / 9	etl_btcsdt_ind1	halfday_trigger_0100		2024-05-05T01:00:00	2024-05-05T01:00:07	2024-05-05T01:01:27	00:01:19.60
✓ Done		323	8 / 8	etl_adausdt	halfday_trigger_0100		2024-05-05T01:00:00	2024-05-05T01:00:07	2024-05-05T01:01:16	00:01:09.68

# Test Result

Topic	Result (Pass/Partial Pass/Fail)
Operational	Pass
Performance	Partial Pass
Security and Privacy	Pass

Table 4: Non-functional Requirements Test

- Some operations provide a response using more than five seconds.
  - The involvement of retrieving large amount of data.

# Discussion

- The system has delivered an acceptable result.
  - It contains the essential components for creating a simple end-to-end machine learning project that includes from the data system to serving system.
- IaC
  - Handle the unstandardized infrastructure problem.
  - Define infrastructures' specification including in this project and manages the provisioning work.
- Data Ingestion and Transformation
  - Handle the data quality and availability problem
  - clean, collect, transform, and centralize data
    - The most important ingredient for working with the machine learning model.



# Discussion

- Training Management System
  - Handle the experiment information recording and the model version distribution problem
  - Store several versions of model, along with metadata
  - Centralized source for serving the model
- Model Monitoring System
  - Handle the unobserved model degradation problem
  - Dashboard on the website to display graphs and metrics
- Model Building and Serving system
  - Handle the model deployment problem
  - API for serving models
    - Ex. in this project: API is used by the model monitoring system

# Discussion

- Each component can tackle various problems of the machine learning system development.
  - When connected together as a bot, it can work without the significant error.
- This bot can demonstrate the simple end-to-end machine learning system
  - It satisfies major functional and non-functional requirements
  - This project has effectively solved the identified issues
- Our project assessment will be covered later in the final section of this presentation.

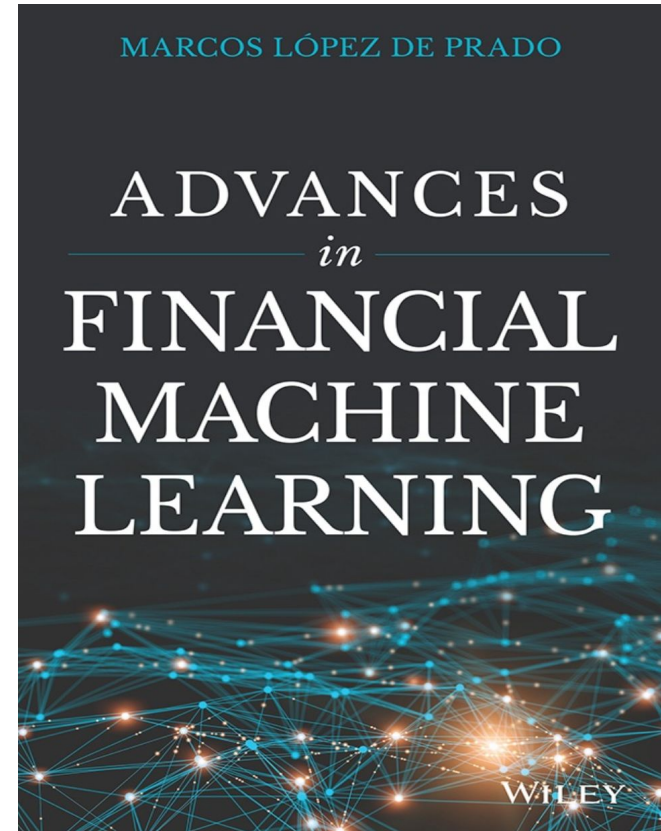
# Discussion: Design Requirement Revisited

- Required Components
  - Source Code Repository: **GitHub**
  - Model Training Infrastructure: **Local Computer and Google Colab**
  - Model Registry: **MLFlow**
  - ML Metadata Stores: **PostgreSQL and S3**
  - Model Serving Component: **API by BentoML and DigitalOcean App Platform**
  - Model Monitoring Component: **Streamlit**
  - Price Prediction Model
    - One or more features as an input
    - A single value as an output
      - Price itself
      - Difference in price

# Model Development

# Research

- Papers and Textbooks
- Other experienced teams



# Summarize

- It will be a gigantic task to actually implement a model that outperform the market by focusing only on one specific token due to our model being OVERFIT.

# Summarize

- Let's imagine that there is an actual trend that could really provide profits to users. In order to gain the highest amount of profit, users have to buy low and sell high according to those trends.
- However, there are so many competitors within the field who will be using statistical and quant approaches to detect the trends.
- Therefore, there are a significant chances that different players could end up detecting the same price trend. Once the trends are detected, some users will frontrun others to sell first. Resulting in the trend to always be changing.

# Summarize

- This means the trend of cryptocurrency prices might still be close to before. However, the change in trend could always be big enough to make the trained model perform worse than just simply holding tokens.



# Our decision

- That's why we shifted our attention toward building models to gain the most gains from shifting between tokens instead of focusing only on one token.
- For example, allocating our fund in day1 to token A and shift our fund in the next day to other tokens once the token A is predicted to have its value reduced.

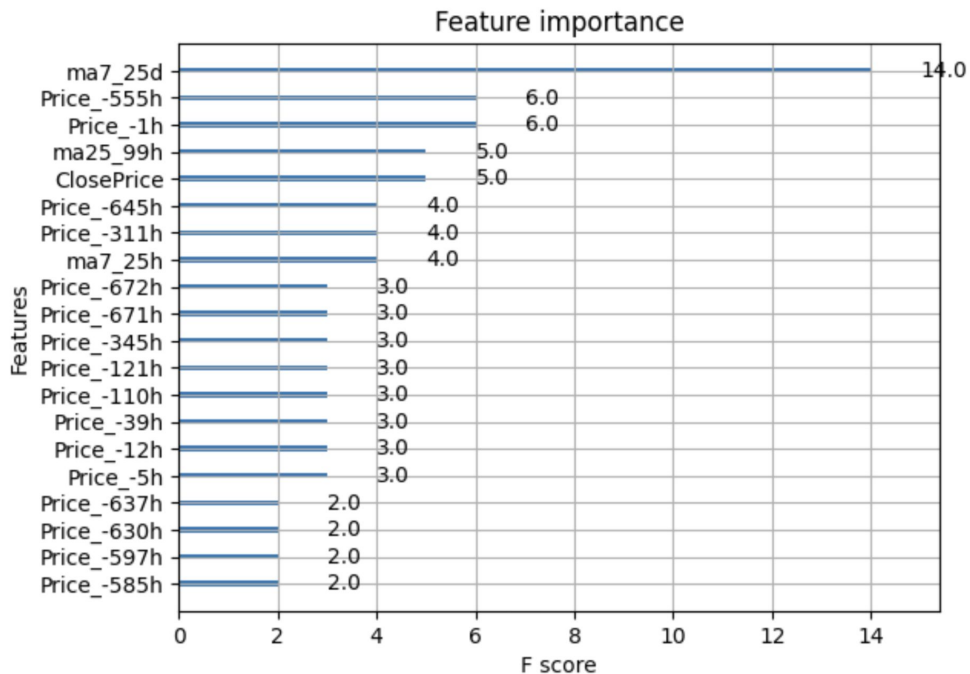
# Preprocess

The indicators we experimented on include but are not limited to the followings:

- Close price (normalized)
- ma7h-ma25h
- ma25h-99h
- ma7d-ma25d
- ma25d-ma99d
- RSI
- MACD
- Every Hour of Historical Price (normalized)
  - Using data from the previous 1,2,3,4,..., 672 hours
- Bollinger Bands

# Preprocess

Then, we plot out the feature significance of those indicators by training them with some of our models.



# Preprocess

We summarize that the indicators that contribute the most to the model performances are the followings

- Close price (normalized)
- ma7h-ma25h
- ma25h-99h
- ma7d-ma25d
- Every 24 Hour Data of Historical Price (normalized)
  - Using data from the previous 1,25,49,73,..., 635, 672 hours

# Model Selection

- We listed out some of the models that are convenient to be implemented. With more time provided, our team could experiment on more models and might receive better results

XGBoost
LassoLars
Bayesian Ridge
Ridge
Lasso
Elastic Net
Huber
CatBoost

# Hypertuning models

- We adopted Optuna in hypertuning all our models.
- We set mean square error (MSE) as the objective function for training.
  - The fewer, the better

```
def objective_lars(trial):  
    alpha = trial.suggest_float('alpha', 1e-6, 1.0)  
    model = LassoLars(alpha=alpha)  
    model.fit(X_train, y_train)  
    y_pred = model.predict(X_test)  
    return mean_squared_error(y_test, y_pred)
```

# Model Training

- We trained our model by having the training process recorded on MLFlow.

```
with mlflow.start_run(experiment_id=experiment_id) as run:  
    best_model = Ridge(**study.best_params)  
    best_model.fit(X_train, y_train)  
    y_pred = best_model.predict(X_test)
```

# Model Evaluation



# Model Evaluation

- Evaluation Datasets
  - 20% of Historical price data of 194 currencies
  - Recent price data of BTCUSDT spanning from January 1<sup>st</sup> to January 31<sup>st</sup> 2024
- Evaluation Metrics
  - Pool Value
  - Classification Report
  - Actual vs Predicted Value Plot

# Pool Value

- Uniquely designed metric that we designed ourselves
- Capable of telling whether the model can outperform the market.
- The logic underlying it is provided as following:
  - We divide our money into 24 parts.
  - Every start of an hour, we will predict which tokens will perform the best on the next 24 hours and we will allocate our money to it.
  - Total pool value represent the percent gains across all the pools.

```
pools = np.array([1.0] * 24)
pool_index = 0
for value1,value2 in zip(filtered_y_pred_data,filtered_y_test_data):
    # print(index,value)
    if value1 > 0.001:
        pools[pool_index] *= (1+value2)
        pool_index = (pool_index + 1) % 24 # Increment the pool index and loop back after the 24th

# Sum the values of all pools

total_pool_value = np.sum(pools)/24
```

# Pool Value

- Higher pool value indicates better performance
  - Pool value is not ROI due to how some data from the test cases are from the same time and we can't use our fund in different places at the same time.
- Pool value of a model surpasses that of the control groups(the pool value from just holding the tokens in all scenarios), indicating that some of our models are able to outperform the market.

Models	Pool Value
XGBoost	$1.226 \times 10^{49}$
LassoLars	$1.216 \times 10^{-6}$
Bayesian Ridge	$1.661 \times 10^{34}$
Ridge	$2.693 \times 10^{36}$
Lasso	$1.516 \times 10^{-6}$
Elastic Net	$1.516 \times 10^{-6}$
Huber	$1.216 \times 10^{30}$
CatBoost	$4.046 \times 10^{125}$
Control group	$1.144 \times 10^{-3}$

# Classification Report

- Evaluated using recent price data of BTCUSDT
- Positive values are classified as Class 1.
- Negative values are classified as Class 0.
- Class 1 precision indicates model accuracy in buying decisions.
- Class 1 recall indicates missed buying opportunities.
- Models can be grouped into 2 groups based on their performance.

# Classification Report

- Group 1
  - Class 1 recall of 1 and class 0 recall of 0.
  - Indicates that the model only predict upward trends; therefore, these models should not be used.
  - The model in this group, including Lasso, LassoLars, and Elastic Net, displays similar performance.

Classification report for ElasticNet model

	precision	recall	f1-score	support
0	0.00	0.00	0.00	20458
1	0.53	1.00	0.69	22742
accuracy			0.53	43200
macro avg	0.26	0.50	0.34	43200
weighted avg	0.28	0.53	0.36	43200

Classification report for Lasso model

	precision	recall	f1-score	support
0	0.00	0.00	0.00	20458
1	0.53	1.00	0.69	22742
accuracy			0.53	43200
macro avg	0.26	0.50	0.34	43200
weighted avg	0.28	0.53	0.36	43200

# Classification Report

- Group 2
  - The model in this group includes XGBoost, Bayesian Ridge, Ridge, Huber, and CatBoost.
  - Huber is the only model showing higher than 50% class 1 precision during the evaluation period.
  - Class 0 recall is significantly higher than the class 1 recall, which may implies the tendency to predict negative values for the model in this group.

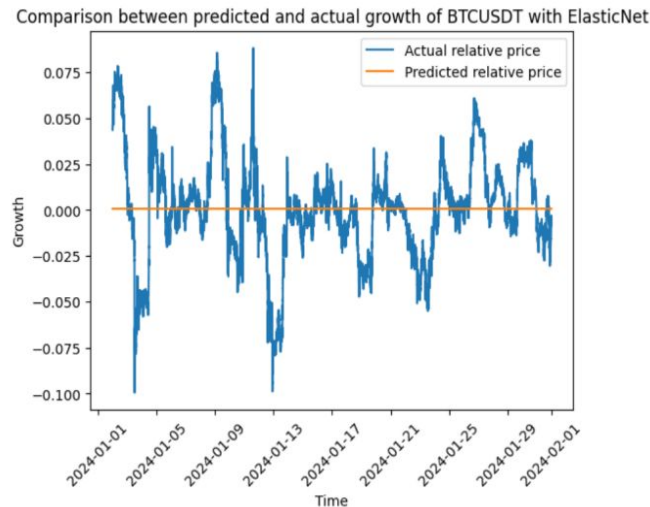
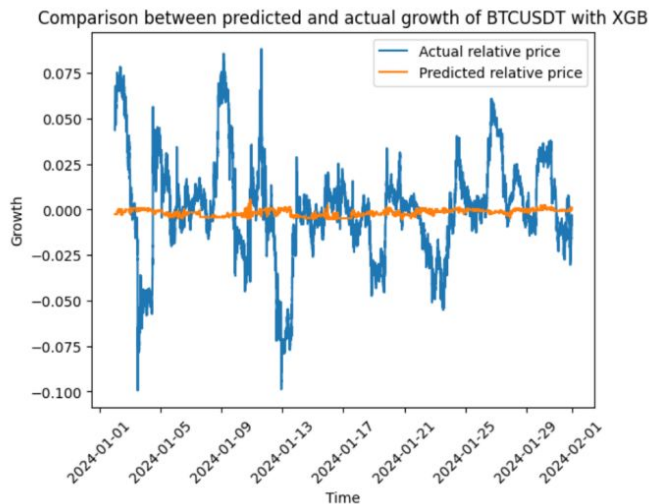
Classification report for Huber model				
	precision	recall	f1-score	support
0	0.47	0.70	0.57	20458
1	0.53	0.30	0.38	22742
accuracy			0.49	43200
macro avg	0.50	0.50	0.47	43200
weighted avg	0.50	0.49	0.47	43200

# Actual vs Predicted Value Plot

- Evaluated using recent price data of BTCUSDT.
- The actual relative price and the predicted value are plotted in the same axis.
- Visualize if the model can capture the price trend.
- Requires human observation and judgement.
- Models can be grouped into 2 groups.

# Actual vs Predicated Value Plot

- Group 1
  - The model in this group includes Lasso, LassoLars, Elastic Net, and XGBoost.
  - Grouping is similar to that of the classification report metric.
  - Fail to capture the price trend.
  - Most predicted values hovering around zero.





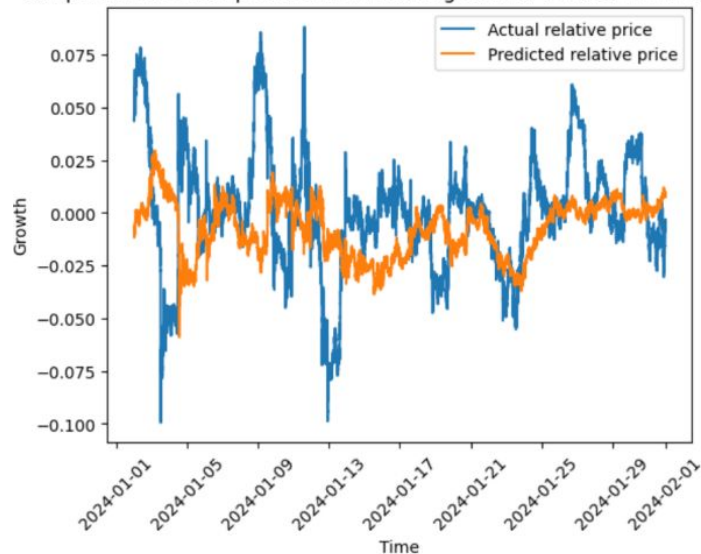
# Actual vs Predecited Value Plot

- Group 2
  - The model in this group includes Bayesian Ridge, Ridge, Huber, and CatBoost.
  - Able to partially capture the price trend
  - With some limitations
    - Tendency to predict negative values aligning with speculation made when observing classification reports
    - Narrower ranges of predicted outputs
    - Some delay in the prediction trend compared to the actual trend

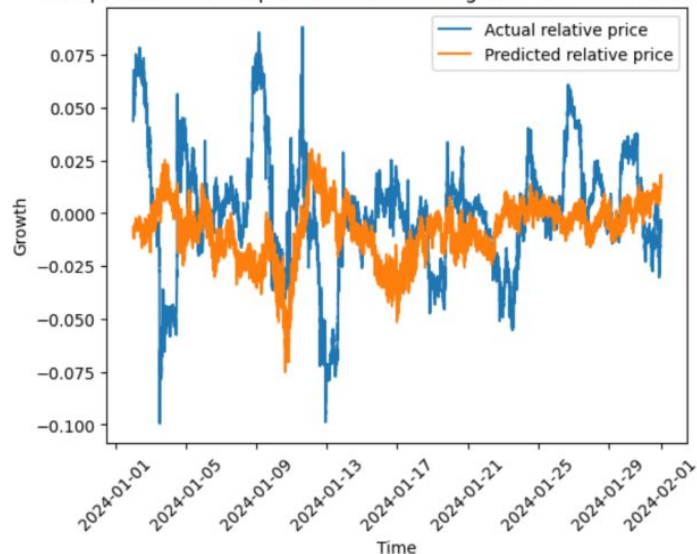
# Actual vs Predicted Value Plot

- Group 2
  - Bayesian Ridge, Ridge, and Huber predict very similar trends, while CatBoost predict a slightly different trend.

Comparison between predicted and actual growth of BTCUSDT with Bayesian



Comparison between predicted and actual growth of BTCUSDT with CBR



# Discussion

- When market conditions resemble those of the train dataset
  - CatBoost emerges as the best-performing model.
- With current market conditions
  - Huber emerges as the best-performing model.
- There is a performance degradation problem on most of the models.
  - Class 1 precision falls below 50% for many models.
  - Hard to gain profits with this precision.
- The Huber model is either more tolerant to performance degradation or coincidentally performs well in the current market condition.

# Discussion

- Overall model performance can be further improved.
  - While the models are able to partially capture the underlying trend, they struggle to predict accurate values.
  - There is a tendency to predict values in a certain class.
- The selected model should be able to perform acceptably in the current market condition.
  - Class 1 precision is at 53%.
  - The model can make accurate buying decisions when the model predicts positive outcomes.

# Discussion

- Reasons that factors to the underperformance of the model are speculated as follows
  - Using several currencies: Although the approach is taken to reduce the chance of overfitting, it may come with trade off as each currency may have different natures.
  - Constantly changing of market conditions: This is evident by the degradation in performance of the models.
  - Cryptocurrency is a highly volatile market: Simple models may not be able to fully capture this volatile nature.
  - MSE might not be the most appropriate metric to be used for the optimization: Other metrics should also be considered.

# Conclusion

# Assessment: System

- Successfully implement a minimum viable product(mvp) of the E2E ML system
- Weakness
  - Data Ingestion and Transformation can be indeed implemented on Mage.ai
    - However, also want to demonstrate serverless method
  - Use multiple cloud providers
    - Cannot fully integrate AWS features to DigitalOcean resources
  - Alias, allowing immediate model changing, requires the model to be downloaded every time.
    - Consume a large amount of memory
    - Still be a problem of our APIs
    - Maybe because we still don't fully understand how BentoML works
      - Can't properly optimized

## Next Step: System

- Implement an automated model training pipeline and refine other components to fully fulfill Microsoft's MLOps Level 2
- Find a method to improve the manual method to handle the historical data.
- Store timezone information along with datetime
  - May have a problem with daylight saving time
- A method of always downloading the model should be refined
- Grafana for model and API monitoring instead of using Streamlit



# Assessment: Model

- There are room for further improvement as evident by the evaluation result using the recent data
- The best performing model should outperform the market in the current condition.
- In the lab environment, several models have successfully outperformed the market.

## Next Step: Model

- Several potential methods for improving model performance and robustness are proposed as follows:
  - Probabilistic forecasting models, including AR, ARMA, and ARIMA, can be alternative solutions.
  - Deep Neural Networks (DNNs) are renowned for their ability to capture highly complex relationships.
  - The problem can be framed as a binary classification problem, opening up other objective functions for model optimization.

Q & A