# How Not to Let Your Model and Data Drift Away Silently

Chengyin Eng

# About

## Chengyin Eng
### Data Scientist @ Databricks

- Machine Learning Practice Team

- Experience
  - Life Insurance
  - Teaching ML in Production, Deep Learning, NLP, etc.

- MS in Computer Science at University of Massachusetts, Amherst

- BA in Statistics & Environmental Studies at Mount Holyoke College, Massachusetts
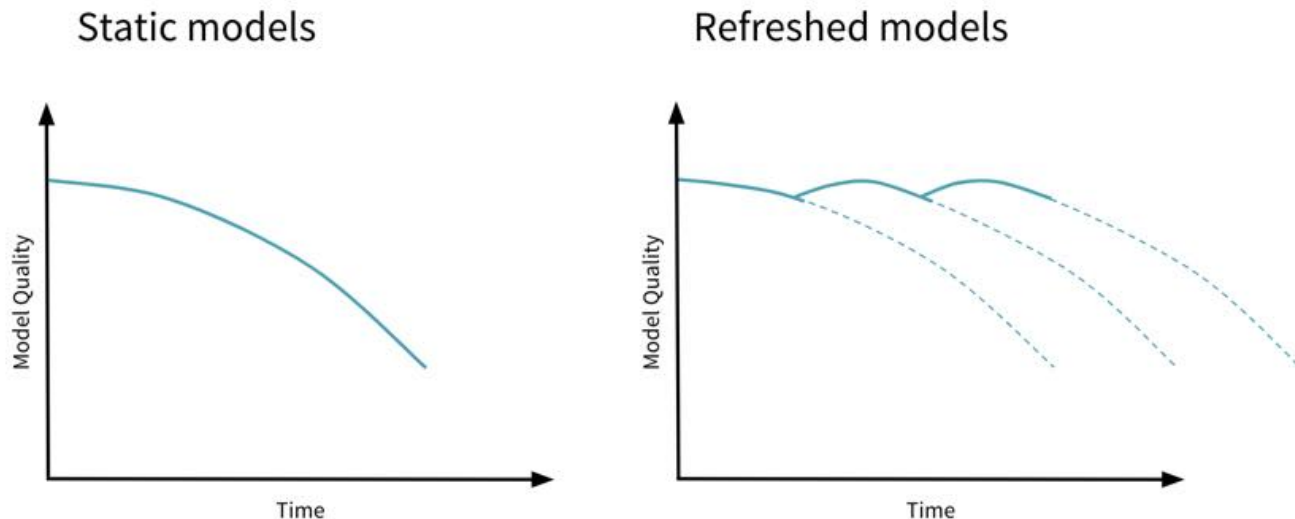
# Outline

- Motivation
- Machine Learning System Life Cycle
- Why Monitor?
  - Types of drift
- What to Monitor?
- How to Monitor?
- Demo

# Why do 96% of ML projects fail in production?

*Neglect maintenance: Lack of re-training and testing*



Static models / Refreshed models — Model Quality vs Time

Sources:

https://databricks.com/blog/2019/09/18/productionizing-machine-learning-from-deployment-to-drift-detection.html

https://www.datanami.com/2020/10/01/most-data-science-projects-fail-but-yours-doesnt-have-to/

# This talk focuses on two questions:

# This talk focuses on two questions:



What are the statistical tests to use when monitoring models in production?

# This talk focuses on two questions:

What are the statistical tests to use when monitoring models in production?
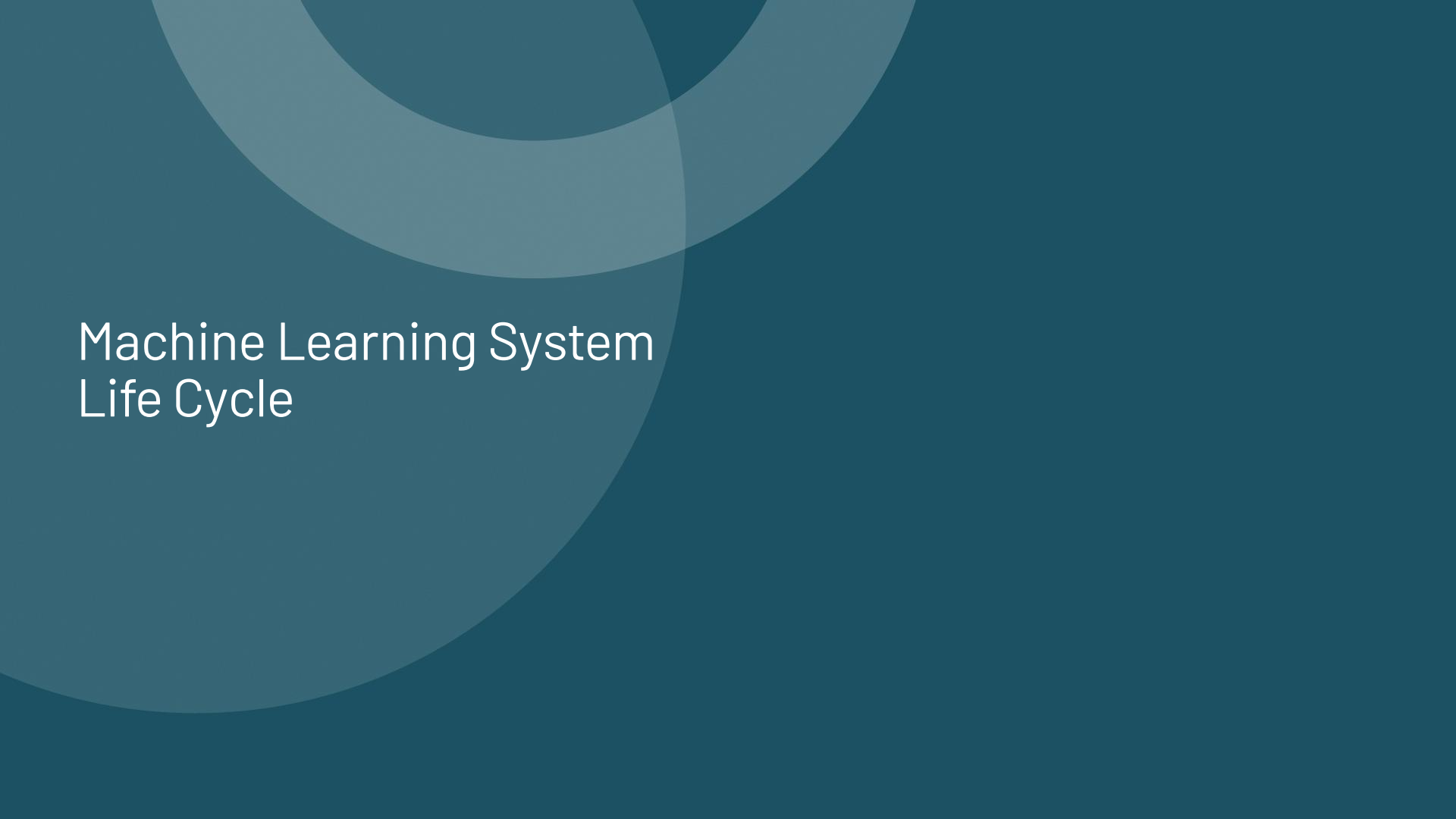
What tools can I use to coordinate the monitoring of data and models?
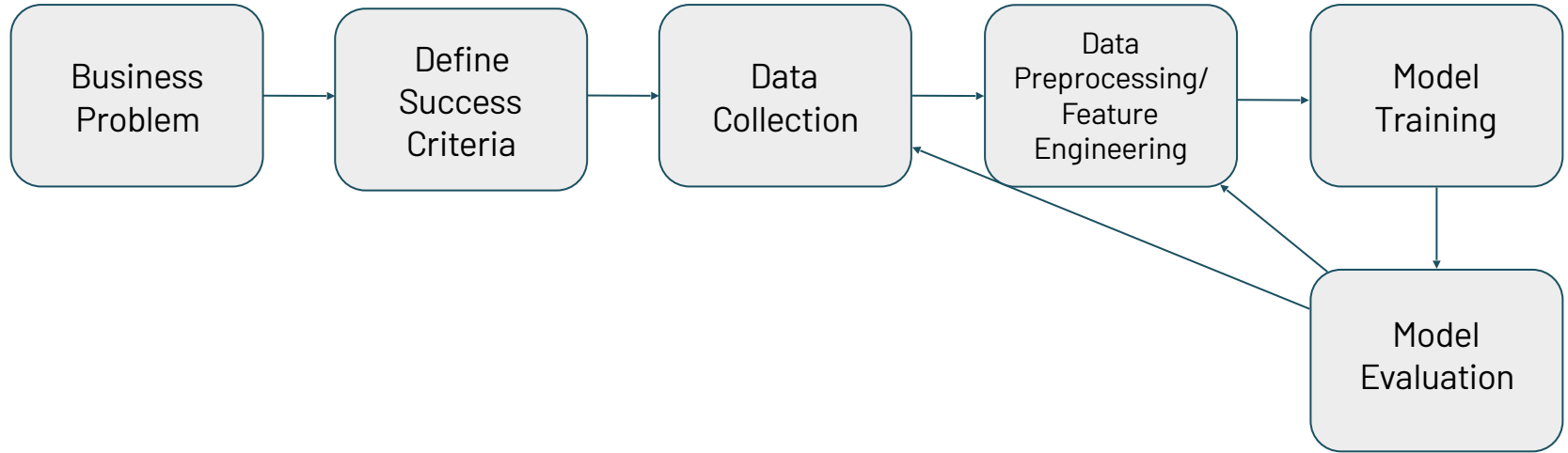
# What this talk is *not*

- A tutorial on model deployment strategies

- An exhaustive walk through of how to robustly test your production ML code

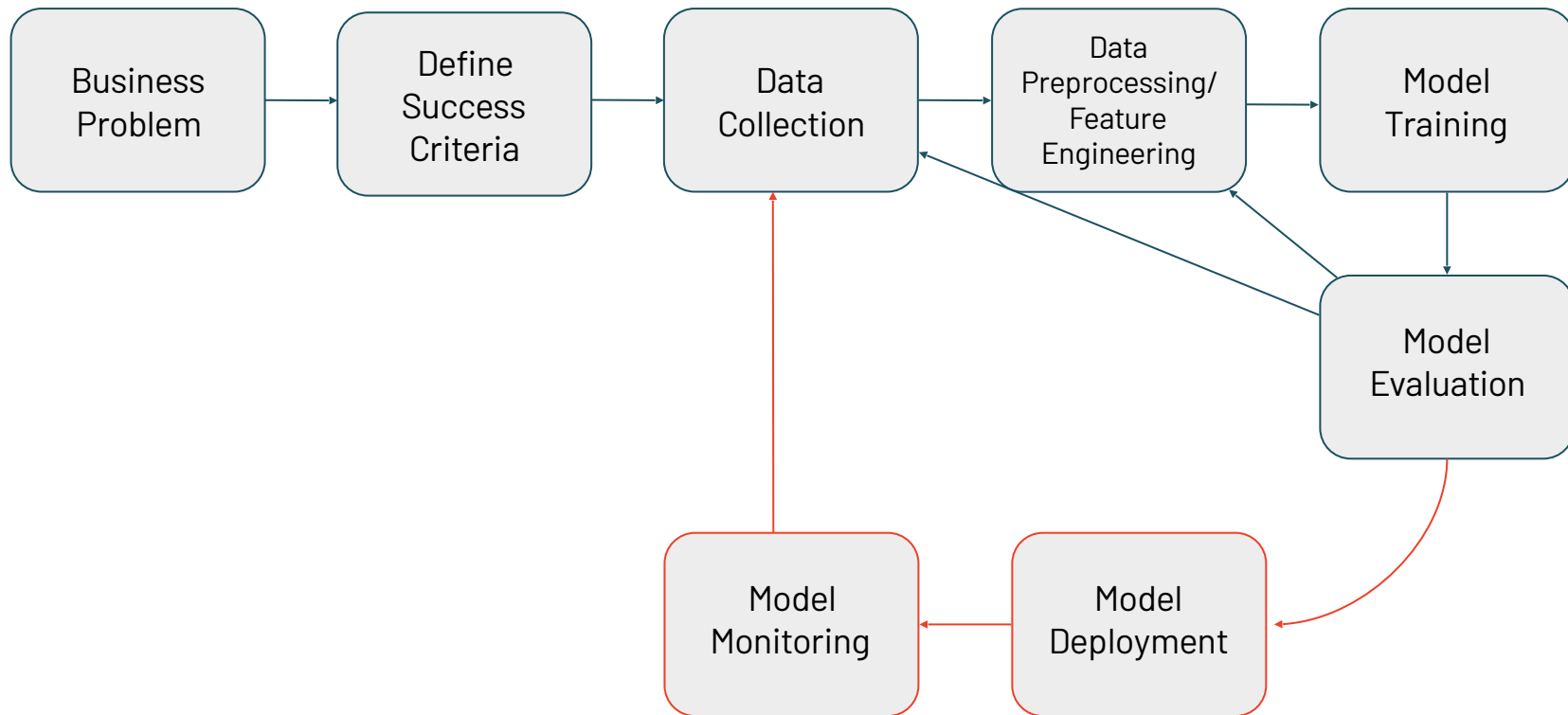- A prescriptive list of *when* to update a model in production

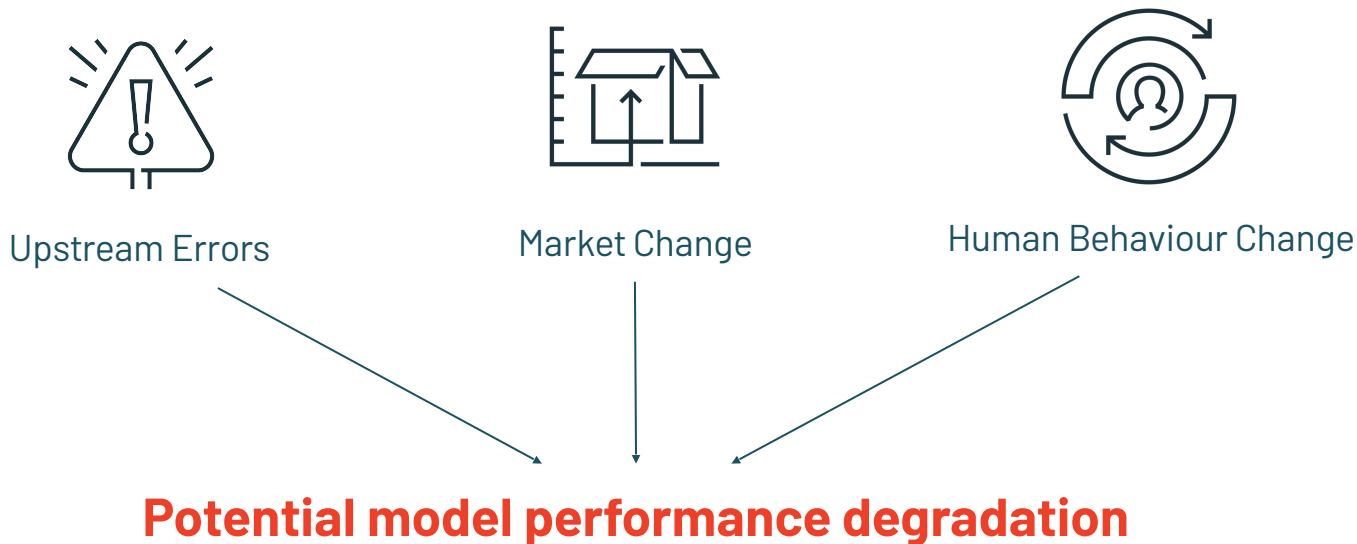# Machine Learning System Life Cycle

# ML system life cycle

```
┌──────────┐     ┌──────────┐     ┌──────────┐     ┌──────────────┐     ┌──────────┐
│ Business │ ──> │  Define  │ ──> │   Data   │ ──> │     Data     │ ──> │  Model   │
│ Problem  │     │ Success  │     │Collection│     │Preprocessing/│     │ Training │
│          │     │ Criteria │     │          │     │   Feature    │     │          │
│          │     │          │     │          │     │ Engineering  │     │          │
└──────────┘     └──────────┘     └──────────┘     └──────────────┘     └──────────┘
                                                                             │
                                                         ┌──────────┐        │
                                                         │  Model   │ <──────┘
                                                         │Evaluation│
                                                         └──────────┘
```

Business Problem → Define Success Criteria → Data Collection → Data Preprocessing/Feature Engineering → Model Training → Model Evaluation

Model Evaluation → Data Collection
Model Evaluation → Data Preprocessing/Feature Engineering

# ML system life cycle

# Why Monitor?

# Model deployment is not the end

*It is the beginning of model measurement and monitoring*

- Data distributions and feature types can change over time due to:



Upstream Errors

Market Change

Human Behaviour Change

**Potential model performance degradation**

Models *will* degrade over time

**Challenge:** catching this when it happens

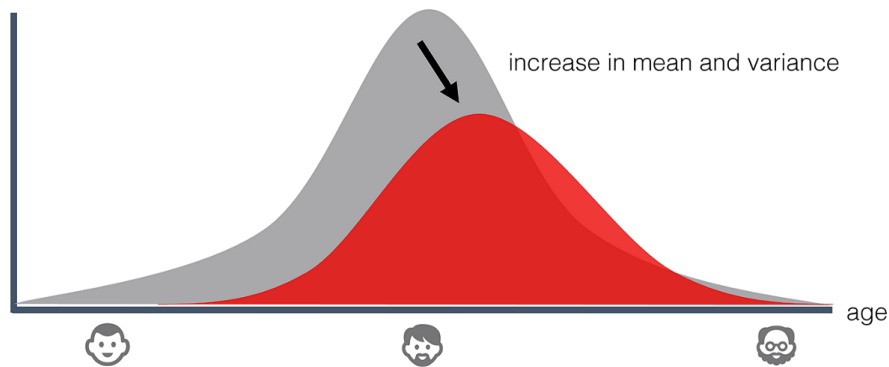# Types of drift

## Data Drift

One of more distributions deviate:

- Input features
- Label
- Model prediction

## Concept Drift

External factors cause the label to evolve

# Data Drift

| Categories | Expected | Observed | Total |
|------------|---------:|---------:|------:|
| A          | 25       | 35       | 60    |
| B          | 25       | 20       | 56    |
| C          | 25       | 25       | 50    |
| D          | 25       | 20       | 45    |
| Total      | 100      | 100      | 100   |

increase in mean and variance

age

Sources:

https://dataz4s.com/statistics/chi-square-test/

https://towardsdatascience.com/machine-learning-in-production-why-you-should-care-about-data-and-concept-drift-d96d0bc907fb

# Concept drift

# Drift types and actions to take

| Drift Type | Retrain using new data | Investigate process | Assess business impact | Consider alternative solutions |
|---|---|---|---|---|
| **Feature Drift** | Y | Y | | |
| **Label Drift** | Y | Y | | |
| **Prediction Drift** | | Y (Model training) | Y | |
| **Concept Drift** | Y (Or tune) | | | Y (Additional feature engineering) |

# What to Monitor?

# What should I monitor?

- Basic summary statistics of features and target

- Distributions of features and target

- Model performance metrics

- Business metrics

# Monitoring tests on data

Numeric Features

- Summary statistics:
  - Median / mean
  - Minimum
  - Maximum
  - Percentage of missing values

- Statistical tests:
  - Mean:
    - Two-sample Kolmogorov-Smirnov (KS) test with Bonferroni correction
    - Mann-Whitney (MW) test
  - Variance:
    - Levene test

# Kolmogorov-Smirnov (KS) test with Bonferroni correction

*Comparison of two continuous distributions*

- Null hypothesis ($H_0$):

  *Distributions x and y come from the same population*

- If the KS statistic has a *p*-value lower than α, reject $H_0$

- Bonferroni correction:
  - Adjusts the α level to reduce false positives
  - $\alpha_{new} = \alpha_{original} / n$, where n = total number of feature comparisons

# Levene test

*Comparison of variances between two continuous distributions*

- Null hypothesis ($H_0$):

$$\sigma^2_1 = \sigma^2_2 = \ldots = \sigma^2_n$$

- If the Levene statistic has a *p*-value lower than $\alpha$, reject $H_0$

# Monitoring tests on data

## Numeric Features

- Summary statistics:
    - Median / mean
    - Minimum
    - Maximum
    - Percentage of missing values

- Statistical tests:
    - Mean:
        - Two-sample Kolmogorov-Smirnov (KS) test with Bonferroni correction
        - Mann-Whitney (MW) test
    - Variance:
        - Levene test

## Categorical Features

- Summary statistics:
    - Mode
    - Number of unique levels
    - Percentage of missing values

- Statistical test:
    - One-way chi-squared test

# One-way chi-squared test

*Comparison of two categorical distributions*

- Null hypothesis ($H_0$):

  Expected distribution = observed distribution

- If the Chi-squared statistic has a *p*-value lower than α, reject $H_0$

# Monitoring tests on models

- Relationship between target and features
  - Numeric Target: Pearson Coefficient
  - Categorical Target: Contingency tables


- Model Performance
  - Regression models: MSE, error distribution plots etc
  - Classification models: ROC, confusion matrix, F1-score etc
  - Performance on data slices


- Time taken to train

# How to Monitor?

# Workflow

# Workflow

# mlflow™

An open-source platform for ML lifecycle that helps with operationalizing ML

## mlflow™ Tracking

Record and query experiments: code, metrics, parameters, artifacts, models

## mlflow™ Projects

Packaging format for reproducible runs on any compute platform

## mlflow™ Models

General model format that standardizes deployment options

## mlflow™ Model Registry

Centralized and collaborative model lifecycle management

MLflow documentation linked here: https://www.mlflow.org/docs/latest/index.html

# MLflow Experiment UI

ⓘ Track machine learning training runs in an experiment. Learn more                                                    ✕

Experiment ID : 2808846569707146                    Artifact Location : dbfs:/databricks/mlflow-tracking/2808846569707146

▸ Notes ✎

Showing 3 matching runs

| ⟳ Refresh | Compare | Delete | Download CSV⬇ | ☰ ▦ | ⚙ Columns | ❓ | 🔍 metrics.rmse < 1 and params.model = "tree" | Search | ☰ Filter | Clear |

| | | | | | Parameters ⟩ | | | Metrics ⟩ | | | | Tags | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Start Time | Run Name | User | Source | Version | Models | delta_path | delta_versic | memory | test_mse | test_r2 | training_ma | estimator_c | estimator_n | sparkDatas |
| ☐ | ⊘ 2021-06-12 11:46:5 | month_2 | chengyi... | 📄 mlops2( | - | 🐦 airbnb_haw.../3 | /Users/... | 2 | None | 48573... | 0.043 | 394.7 | sklearn.... | Pipeline | path=d... |
| ☐ | ⊘ 2021-06-12 11:46:3 | month_1 | chengyi... | 📄 mlops2( | - | 🐦 airbnb_haw.../2 | /Users/... | 1 | None | 16962... | 0.113 | 276 | sklearn.... | Pipeline | path=d... |
| ☐ | ⊘ 2021-06-12 11:46:0 | month_0 | chengyi... | 📄 mlops2( | - | 🐦 airbnb_haw.../1 | /Users/... | 0 | None | 24371... | 0.08 | 230.8 | sklearn.... | Pipeline | path=d... |

# MLflow Model Registry

airbnb_hawaii ▾

Notify me about ⓘ    All new activity ⌄

Details    Serving

Created Time : 2021-06-12 11:46:15      Last Modified : 2021-06-12 11:47:11      Creator : chengyin.eng@databricks.com

▸ Description   Edit

▸ Tags

▾ Versions   | All | Active (2) |   Compare

| | | | Version | Registered at ▾ | Created by | Stage | Pending Requests | Description |
|---|---|---|---|---|---|---|---|---|
| ☐ | ✓ | 🔔 | Version 3 | 2021-06-12 11:47:05 | chengyin.eng@databricks.com | Staging | — | |
| ☐ | ✓ | 🔔 | Version 2 | 2021-06-12 11:46:41 | chengyin.eng@databricks.com | Production | — | |
| ☐ | ✓ | 🔔 | Version 1 | 2021-06-12 11:46:15 | chengyin.eng@databricks.com | Archived | — | |

# DELTA LAKE

## An open-source data storage format that allows ACID transaction and metadata handling

Parquet files combined with transaction logs

```
/mytable/_delta_log/00000000000000000000.json
/mytable/_delta_log/00000000000000000001.json
/mytable/_delta_log/00000000000000000003.json
/mytable/_delta_log/00000000000000000003.checkpoint.parquet
/mytable/_delta_log/_last_checkpoint
/mytable/part-00000-3935a07c-416b-4344-ad97-2a38342ee2fc.c000.snappy.parquet
```

Read older versions of data using time travel

```Python
df1 = spark.read.format("delta").option("timestampAsOf", timestamp_string).load("/delta/events")
df2 = spark.read.format("delta").option("versionAsOf", version).load("/delta/events")
```

Delta documentation linked here: https://docs.delta.io/latest/index.html

# Delta Table History

```
1  gold_delta_path = "/Users/chengyin.eng@databricks.com/mlops2021/data/airbnb_hawaii_delta"
2  display(DeltaTable.forPath(spark, gold_delta_path).history())
```

▸ (1) Spark Jobs

| | version | timestamp | userId | userName | operation | operationParameters |
|---|---|---|---|---|---|---|
| 1 | 2 | 2021-06-12T16:46:54.000+0000 | 4470711271069202 | chengyin.eng@databricks.com | WRITE | ▸ {"mode": "Append", |
| 2 | 1 | 2021-06-12T16:46:30.000+0000 | 4470711271069202 | chengyin.eng@databricks.com | WRITE | ▸ {"mode": "Append", |
| 3 | 0 | 2021-06-12T16:46:03.000+0000 | 4470711271069202 | chengyin.eng@databricks.com | WRITE | ▸ {"mode": "ErrorIfExis |

Showing all 3 rows

# Demo Notebook

http://bit.ly/mlops2021-drifting-away

# Conclusion

- Model measurement and monitoring are crucial when operationalizing ML models
- No one-size fits all
  - Domain & problem specific considerations
- Reproducibility
  - Enable rollbacks and maintain record of historic performance

# Literature resources

- Paleyes et al 2021. Challenges in Deploying ML

- Klaise et al. 2020 Monitoring and explainability of models in production

- Rabanser et al 2019 Failing Loudly: An Empirical Study of Methods for Detecting Dataset Shift

- Martin Fowler: Continuous Delivery for Machine Learning

# Emerging open-source monitoring packages

- EvidentlyAI

- Data Drift Detector

- Alibi Detect