

Projeto da Disciplina: Engenharia de Machine Learning

- A solução criada nesse projeto deve ser disponibilizada em repositório git e disponibilizada em servidor de repositórios (Github (recomendado), Bitbucket ou Gitlab). O projeto deve obedecer o Framework TDSP da Microsoft (estrutura de arquivos, arquivo requirements.txt e arquivo README - com as respostas pedidas nesse projeto, além de outras informações pertinentes). Todos os artefatos produzidos deverão conter informações referentes a esse projeto (não serão aceitos documentos vazios ou fora de contexto). Escreva o link para seu repositório.

- Repos:

<https://github.com/pmneto/engenhariaml-docker.git>

https://github.com/pmneto/pd_engenharia_ml.git

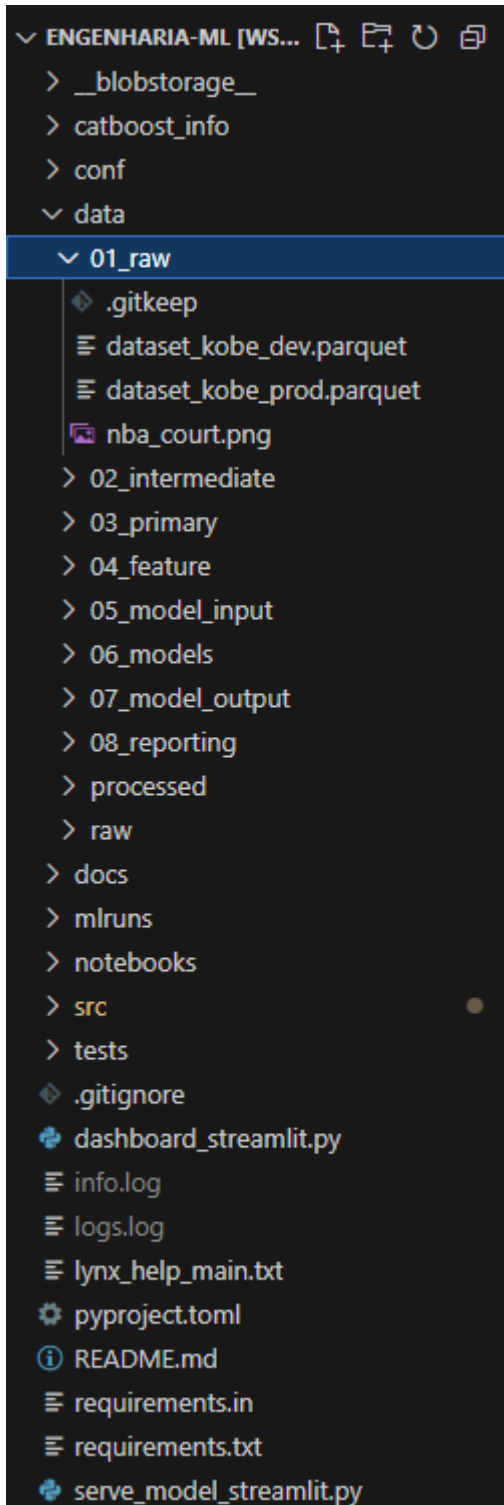
- Iremos desenvolver um preditor de arremessos usando duas abordagens (regressão e classificação) para prever se o "Black Mamba" (apelido de Kobe) acertou ou errou a cesta.

Baixe os dados de desenvolvimento e produção [aqui](#) (datasets: dataset_kobe_dev.parquet e dataset_kobe_prod.parquet). Salve-os numa pasta /data/raw na raiz do seu repositório.

Para começar o desenvolvimento, desenhe um diagrama que

demonstra todas as etapas necessárias para esse projeto, desde a aquisição de dados, passando pela criação dos modelos, indo até a operação do modelo.

[illegible]



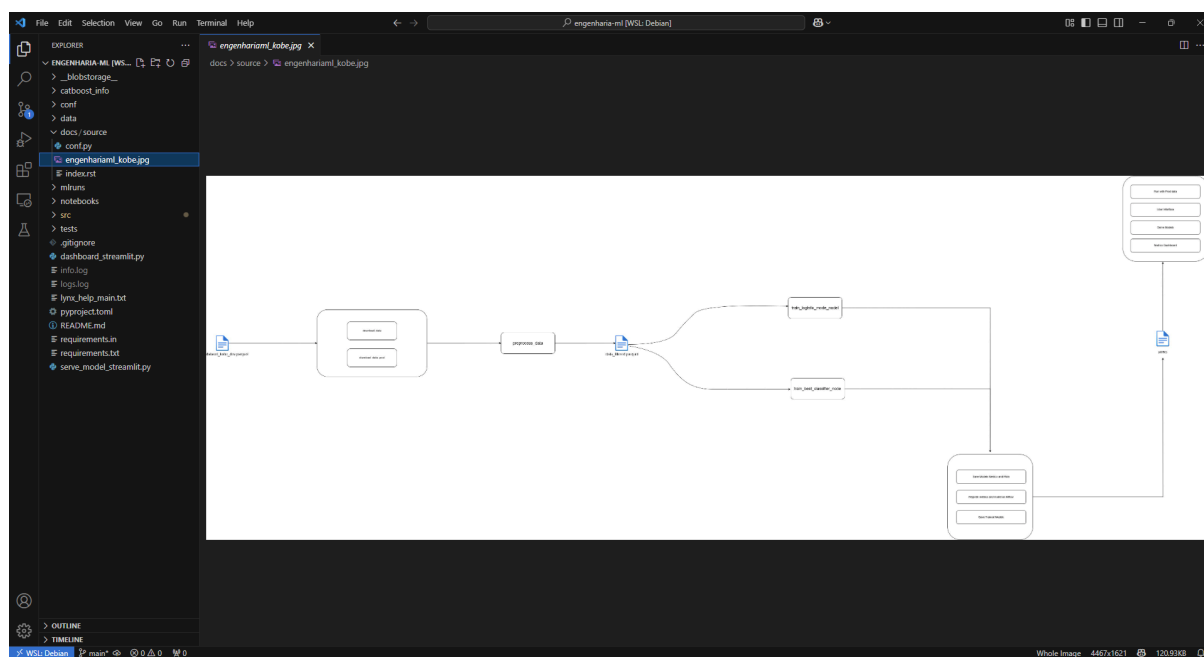
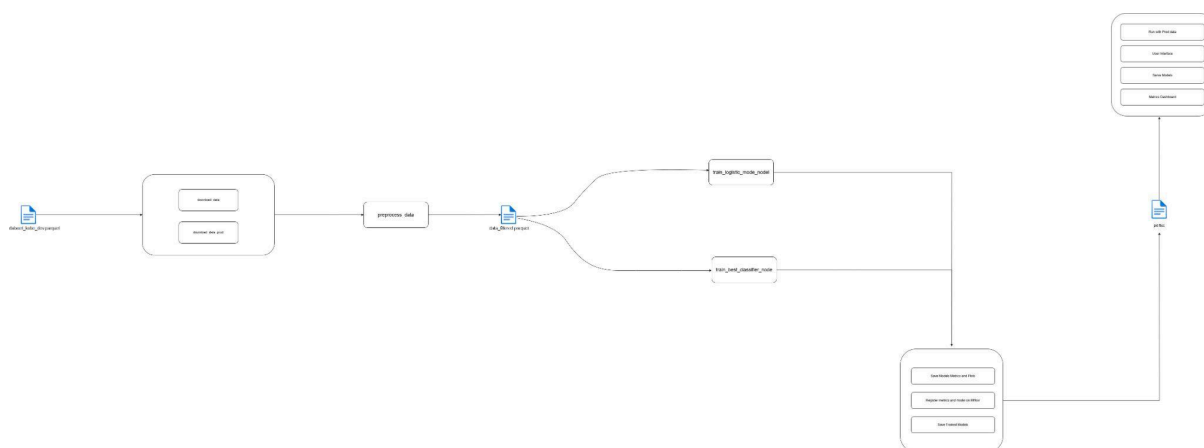


Diagrama com fluxo das implementações:



- Como as ferramentas Streamlit, MLFlow, PyCaret e Scikit-Learn auxiliam na construção dos pipelines descritos anteriormente? A resposta deve abranger os seguintes aspectos:

a. Rastreamento de experimentos;

O MLFlow tem a capacidade de fazer todo o rastreamento do ciclo de vida da aplicação, no contexto de modelagem de inteligência artificial. As ferramentas PyCaret, Scikit-Learn e Streamlit facilitam este processo também por possuírem integração com o MLFlow, permitindo inclusão de objetos e metadados e muitas outras facilidades na correta implementação de processos de engenharia de software dedicados a machine learning.

b. Funções de treinamento;

O PyCaret ajuda no contexto de simplificação dos processos para emprego de modelagem de machine learning. Usando a biblioteca, é possível selecionar de forma bastante otimizada modelos e parametros mais adequados para seu caso de uso, permitindo uma diminuição de esforço para testagem de suposições e implementações complexas. Neste sentido, o MLFlow colabora por conseguir manter históricos e facilidades, além de apis dedicadas para finalidades como servir modelos.

c. Monitoramento da saúde do modelo;

O MLflow permite o registro de métricas como *log loss*, *F1-score* e outras que refletem o desempenho dos modelos ao longo do tempo. Isso facilita a identificação de possíveis quedas de performance (model drift).

Além disso, os dados de produção podem ser comparados com os dados de treino por meio de análises em dashboards construídos no Streamlit, que oferece uma maneira rápida e interativa de visualizar a "saúde" do modelo em operação. Assim, é possível criar alertas e mecanismos para reavaliação do modelo com base em métricas e distribuição de dados.

d. Atualização de modelo;

Com os experimentos e métricas versionadas no MLflow, é possível revalidar ou retreinar modelos com novas amostras de dados com facilidade. O pipeline de engenharia permite reexecutar etapas específicas, como pré-processamento, seleção de features ou tuning de hiperparâmetros, minimizando esforço e maximizando rastreabilidade. O PyCaret, com suas funções de reavaliação automática e tuning, facilita a atualização sem necessidade de reescrever código. Esse ciclo de retreinamento pode ser tanto reativo (após perda de performance) quanto preditivo (antecipando possíveis desvios detectados no monitoramento).

e. Provisionamento (Deployment).

O MLflow oferece uma interface para servir modelos como APIs REST de forma simples e rápida com o comando `mlflow models serve`. Essa funcionalidade permite integrar o modelo em qualquer aplicação que consiga fazer requisições HTTP. Já o Streamlit pode ser usado para criar interfaces interativas para usuários finais acessarem os resultados dos modelos ou simular previsões em tempo real. Por fim, com o uso de Docker, toda a infraestrutura do projeto (incluindo Kedro, MLflow, e Streamlit) pode ser provisionada e mantida com consistência em qualquer ambiente, seja local ou em nuvem.

- Com base no diagrama realizado na questão 2, aponte os artefatos que serão criados ao longo de um projeto. Para cada artefato, a descrição detalhada de sua composição.

dataset_kobe_dev.parquet - Dataset de desenvolvimento baixado do link indicado do Github no descritivo do projeto.

dataset_kobe_prod.parquet - Dataset de produção baixado do link indicado do Github no descritivo do projeto.

data_filtered.parquet - filtragem solicitada em 5.b

confusion_matrix_reg_log.png - Matriz de Confusão da Regressão logística.

confusion_matrix.png - Matriz de confusão do modelo de classificação.

prediction_map_prod.png - Mapa de predição com os resultados de predict.

prediction_map_raw_data.png - Mapa de predição usando os dados originais.

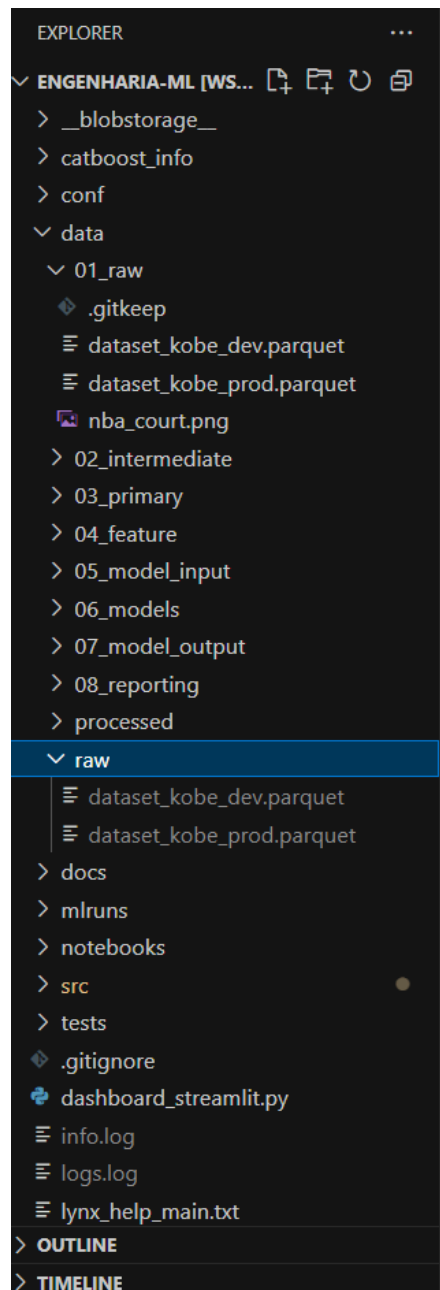
prediction_map_reg_log.png - Mapa de predição levando em consideração a regressão logística.

roc_curve_reg_log.png - Curva ROC da Regressão Logística

roc_curve.png - Curva ROC da Classificação

predictions_prod.parquet - Resultado da inferência em produção

- Implemente o pipeline de processamento de dados com o mlflow, rodada (run) com o nome "PreparacaoDados":
 - a. Os dados devem estar localizados em
"/data/raw/dataset_kobe_dev.parquet" e
"/data/raw/dataset_kobe_prod.parquet"



b. Observe que há dados faltantes na base de dados! As linhas que possuem dados faltantes devem ser desconsideradas. Para esse exercício serão apenas consideradas as colunas:

i. lat

ii. lng

- iii. minutes remaining
- iv. period
- v. playoffs
- vi. shot_distance

A variável shot_made_flag será seu alvo, onde 0 indica que Kobe errou e 1 que a cesta foi realizada. O dataset resultante será armazenado na pasta `"/data/processed/data_filtered.parquet"`. Ainda sobre essa seleção, qual a dimensão resultante do dataset?

```
def preprocess_data(df: pd.DataFrame, is_prod=False) -> pd.DataFrame:
    """Pré-processa os dados (dev ou prod) e salva resultado no caminho correto."""
    mlflow.set_experiment("engenharia_ml")
    run_name = "PreparacaoDados"

    with mlflow.start_run(run_name=run_name):
        filename = (
            "data/processed/data_filtered_prod.parquet"
            if is_prod
            else "data/processed/data_filtered.parquet"
        )

        os.makedirs(os.path.dirname(filename), exist_ok=True)

        features = ['lat', 'lng', 'minutes_remaining', 'period', 'playoffs', 'shot_distance', 'shot_made_flag']

        df.columns = df.columns.str.lower().str.replace(' ', '_').str.replace('(', '').str.replace(')', '')

        if 'lon' in df.columns:
            df = df.rename(columns={'lon': 'lng'})

        df = df.dropna().drop_duplicates()
        df = df[features]

        df.to_parquet(filename)
        mlflow.log_artifact(filename)

        print(f"[{'PROD' if is_prod else 'DEV'}] Dataset salvo com {df.shape[0]} linhas e {df.shape[1]} colunas")
        return df
```

- vii. Separe os dados em treino (80%) e teste (20 %) usando uma escolha aleatória e estratificada.

Armazene os datasets resultantes em
"/Data/processed/base_{train|test}.parquet .
Explique como a escolha de treino e teste
afetam o resultado do modelo final. Quais
estratégias ajudam a minimizar os efeitos de
viés de dados.

```
def setup_experiment(data, session_id: int, cv_folds: int) -> ClassificationExperiment:
    exp = ClassificationExperiment()
    exp.setup(
        data=data,
        target="shot_made_flag",
        session_id=session_id,
        train_size=0.8,
        fold=cv_folds,
        fold_shuffle=True,
        html=False,
        log_experiment=False,
        log_plots=False,
        verbose=False
    )
    return exp
```

> engenharia_ml_kobe > pipelines > data_science > nodes.py > train_decision_tree_model

```
14 def train_logistic_model(data, session_id: int, cv_folds: int):
15
16     mlflow.set_experiment("kobe_classificacao")
17     with mlflow.start_run(run_name="Treinamento"):
18         exp = setup_experiment(data, session_id, cv_folds)
19
20         logistic = exp.create_model("lr", fold=cv_folds)
21         tuned_logistic = exp.tune_model(logistic, n_iter=10, fold=cv_folds, optimize="AUC")
22         pred = exp.predict_model(tuned_logistic)
23         plot_previsoes_modelo('1', pred)
24
25
26
27
28     y_test, y_pred_label, y_pred_proba = extract_metrics(pred)
29     mlflow.log_metric("log_loss", log_loss(y_test, y_pred_proba))
30     mlflow.log_metric("f1_score", f1_score(y_test, y_pred_label))
31
32
33     metricas = {
34         "model_type": type(tuned_logistic).__name__,
35         "log_loss": log_loss(y_test, y_pred_proba),
36         "accuracy": accuracy_score(y_test, y_pred_label),
37         "precision": precision_score(y_test, y_pred_label, zero_division=0),
38         "recall": recall_score(y_test, y_pred_label, zero_division=0),
39         "f1_score": f1_score(y_test, y_pred_label),
40         "cv_folds": cv_folds,
41         "train_rows": len(data),
42         "timestamp": datetime.now().isoformat()
43     }
44
45     salvar_metricas_csv('1', metricas)
46
47     plot_confusion_matrix('1', y_test, y_pred_label)
48     plot_roc_curve('1', y_test, y_pred_proba)
49
50     os.makedirs("data/06_models", exist_ok=True)
51     X_input = pred.drop(columns=["prediction_label", "prediction_score", "Label", "Score", "shot_made_flag"], errors="ignore")
52     exp.save_model(tuned_logistic, "data/06_models/logistic_model")
53
54
55     mlflow.sklearn.log_model(
56         sk_model=tuned_logistic,          # seu modelo treinado
57         artifact_path="model",            # o nome do diretório no MLflow
58         input_example=X_input.iloc[:1],    # opcional mas recomendado
59         signature=mlflow.models.infer_signature(X_input, tuned_logistic.predict(X_input))
60     )
61
62     mlflow.log_param("model_type", "LogisticRegression")
63     mlflow.log_param("cv_folds", cv_folds)
64     mlflow.log_metric("train_rows", len(data))
65
66     return tuned_logistic
67
```

```

def train_best_classifier(data, session_id: int, cv_folds: int):

    mlflow.set_experiment("kobe_classificacao")
    with mlflow.start_run(run_name="Treinamento"):
        exp = setup_experiment(data, session_id, cv_folds)

        best_model = exp.compare_models(exclude=["lr"], fold=cv_folds, n_select=1)
        tuned_best = exp.tune_model(best_model, n_iter=10, fold=cv_folds, optimize="AUC")
        pred = exp.predict_model(tuned_best)
        plot_previsoes_modelo('c', pred)

        y_test, y_pred_label, y_pred_proba = extract_metrics(pred)
        mlflow.log_metric("log_loss", log_loss(y_test, y_pred_proba))
        mlflow.log_metric("f1_score", f1_score(y_test, y_pred_label))

        metricas = {
            "model_type": type(tuned_best).__name__,
            "log_loss": log_loss(y_test, y_pred_proba),
            "accuracy": accuracy_score(y_test, y_pred_label),
            "precision": precision_score(y_test, y_pred_label, zero_division=0),
            "recall": recall_score(y_test, y_pred_label, zero_division=0),
            "f1_score": f1_score(y_test, y_pred_label),
            "cv_folds": cv_folds,
            "train_rows": len(data),
            "timestamp": datetime.now().isoformat()
        }

        salvar_metricas_csv('c', metricas)

        plot_confusion_matrix('c', y_test, y_pred_label)
        plot_roc_curve('c', y_test, y_pred_proba)

        os.makedirs("data/06_models", exist_ok=True)
        # Salva o modelo treinado
        X_input = pred.drop(columns=["prediction_label", "prediction_score", "Label", "Score", "shot_made_flag"], errors="ignore")
        _, model_path = exp.save_model(tuned_best, "data/06_models/best_classifier")

        # Faz o log do artefato corretamente
        mlflow.sklearn.log_model(
            sk_model=best_model,
            artifact_path="model",
            input_example=X_input.iloc[:1],
            signature=infer_signature(X_input, best_model.predict(X_input))
        )

        mlflow.log_param("model_type", type(tuned_best).__name__)
        mlflow.log_param("cv_folds", cv_folds)
        mlflow.log_metric("train_rows", len(data))

    return tuned_best

```

viii. Registre os parâmetros (% teste) e métricas (tamanho de cada base) no MIFlow

- Implementar o pipeline de treinamento do modelo com o MLFlow usando o nome "Treinamento"
 - a. Com os dados separados para treinamento, treine um modelo com regressão logística do sklearn usando a biblioteca pyCaret.

```
def train_logistic_model(data, session_id: int, cv_folds: int):

    mlflow.set_experiment("kobe_classificacao")
    with mlflow.start_run(run_name="Treinamento"):
        exp = setup_experiment(data, session_id, cv_folds)

        logistic = exp.create_model("lr", fold=cv_folds)
        tuned_logistic = exp.tune_model(logistic, n_iter=10, fold=cv_folds, optimize="AUC")
        pred = exp.predict_model(tuned_logistic)
        plot_previsoes_modelo('1', pred)

        y_test, y_pred_label, y_pred_proba = extract_metrics(pred)
        mlflow.log_metric("log_loss", log_loss(y_test, y_pred_proba))
        mlflow.log_metric("f1_score", f1_score(y_test, y_pred_label))

        metricas = {
            "model_type": type(tuned_logistic).__name__,
            "log_loss": log_loss(y_test, y_pred_proba),
            "accuracy": accuracy_score(y_test, y_pred_label),
            "precision": precision_score(y_test, y_pred_label, zero_division=0),
            "recall": recall_score(y_test, y_pred_label, zero_division=0),
            "f1_score": f1_score(y_test, y_pred_label),
            "cv_folds": cv_folds,
            "train_rows": len(data),
            "timestamp": datetime.now().isoformat()
        }

        salvar_metricas_csv('1', metricas)

        plot_confusion_matrix('1', y_test, y_pred_label)
        plot_roc_curve('1', y_test, y_pred_proba)

    os.makedirs("data/06_models", exist_ok=True)
    X_input = pred.drop(columns=["prediction_label", "prediction_score", "Label", "Score", "shot_made_flag"], errors="ignore")
    exp.save_model(tuned_logistic, "data/06_models/logistic_model")

    mlflow.sklearn.log_model(
        sk_model=tuned_logistic,          # seu modelo treinado
        artifact_path="model",            # o nome do diretório no MLflow
        input_example=X_input.iloc[:1],   # opcional mas recomendado
        signature=mlflow.models.infer_signature(X_input, tuned_logistic.predict(X_input))
    )

    mlflow.log_param("model_type", "LogisticRegression")
    mlflow.log_param("cv_folds", cv_folds)
    mlflow.log_metric("train_rows", len(data))
    mlflow.log_metric("test_rows", len(y_test))

    return tuned_logistic
```


b. Registre a função custo "log loss" usando a base de teste

```
y_test, y_pred_label, y_pred_proba = extract_metrics(pred)
mlflow.log_metric("log_loss", log_loss(y_test, y_pred_proba))
mlflow.log_metric("f1_score", f1_score(y_test, y_pred_label))
```

- c. Com os dados separados para treinamento, treine um modelo de árvore de decisão do sklearn usando a biblioteca pyCaret.

```
21
22 def train_decision_tree_model(data: pd.DataFrame, session_id: int, cv_folds: int):
23     """Treina e registra um modelo de árvore de decisão com PyCaret e MLflow."""
24
25
26     mlflow.set_experiment("kobe_classificacao")
27
28     with mlflow.start_run(run_name="Treinamento_ArvoreDecisao"):
29         # Setup do experimento
30         exp = ClassificationExperiment()
31         exp.setup(
32             data=data,
33             target="shot_made_flag",
34             session_id=session_id,
35             train_size=0.8,
36             fold=cv_folds,
37             fold_shuffle=True,
38             html=False,
39             log_experiment=False,
40             verbose=False
41         )
42
43         # Criação e tuning do modelo de árvore de decisão
44         dt_model = exp.create_model("dt", fold=cv_folds)
45         tuned_dt = exp.tune_model(dt_model, n_iter=10, fold=cv_folds, optimize="F1")
46
47         # Predição
48         pred = exp.predict_model(tuned_dt)
49
50         # Extração de métricas
51         y_test = pred["shot_made_flag"]
52         y_pred_label = pred["prediction_label"]
53         y_pred_proba = pred["prediction_score"]
54
55         logloss = log_loss(y_test, y_pred_proba)
56         f1 = f1_score(y_test, y_pred_label)
57
58         # Log de métricas
59         mlflow.log_metric("log_loss", logloss)
60         mlflow.log_metric("f1_score", f1)
61         mlflow.log_metric("train_rows", int(len(data) * 0.8))
62         mlflow.log_metric("test_rows", int(len(data) * 0.2))
63
64         # Log de parâmetros
65         mlflow.log_param("model_type", "DecisionTree")
66         mlflow.log_param("cv_folds", cv_folds)
67         mlflow.log_param("test_size", 0.2)
68
69         # Salvar modelo
70         os.makedirs("data/06_models", exist_ok=True)
71         X_input = pred.drop(columns=["prediction_label", "prediction_score", "Label", "Score", "shot_made_flag"], errors="ignore")
72         exp.save_model(tuned_dt, "data/06_models/decisionTree")
73
74         # Log do modelo completo
75         mlflow.sklearn.log_model(
76             sk_model=tuned_dt,
77             artifact_path="model",
78             input_example=X_input.iloc[:1],
79             signature=infer_signature(X_input, tuned_dt.predict(X_input))
80         )
81
82         print("[OK] Modelo de árvore de decisão treinado e registrado no MLflow.")
83         return tuned_dt
84
```

- d. Registre a função custo "log loss" e F1_score para o modelo de árvore.

```

f1 = f1_score(y_test, y_pred_label)

# Log de métricas
mlflow.log_metric("log_loss", logloss)
mlflow.log_metric("f1_score", f1)
mlflow.log_metric("train_rows", int(len(data)*0.8))
mlflow.log_metric("test_rows", int(len(data)*0.2))

# Log de parâmetros
mlflow.log_param("model_type", "DecisionTree")
mlflow.log_param("cv_folds", cv_folds)
mlflow.log_param("test_size", 0.2)

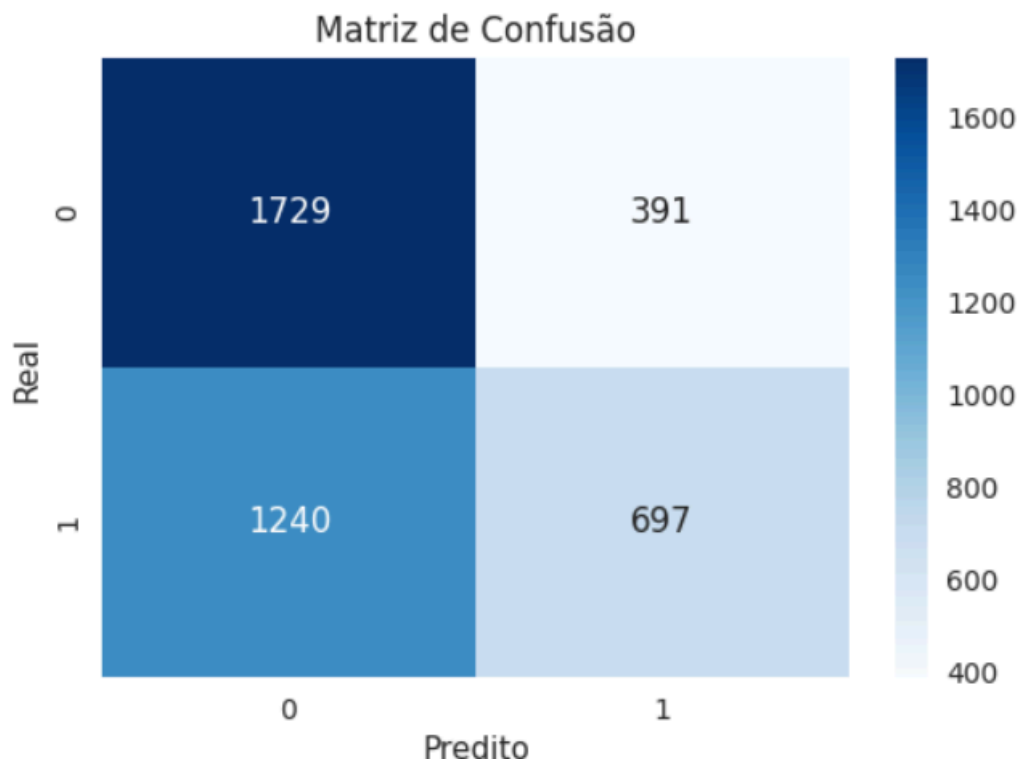
# Salvar modelo
os.makedirs("data/models", exist_ok=True)

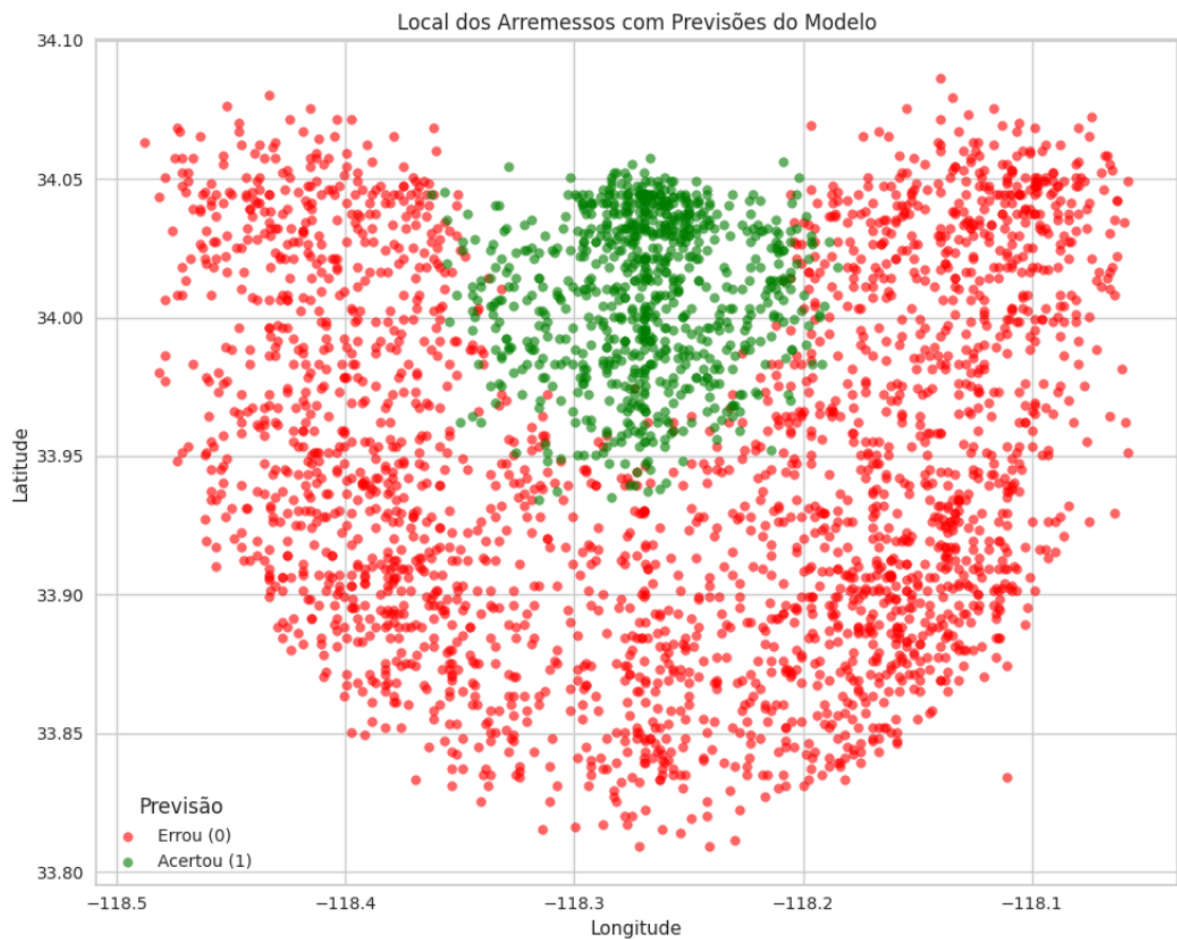
```

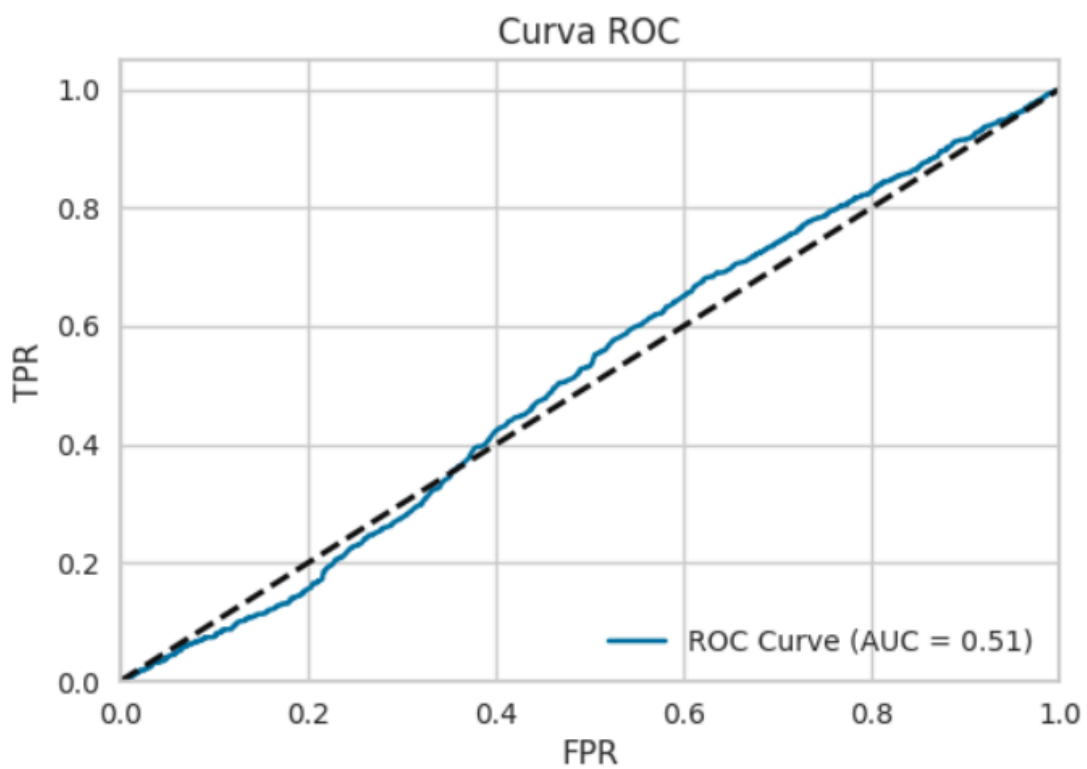
- e. Selecione um dos dois modelos para finalização e justifique sua escolha.

Escolhi a Regressão Logística por que apresentou melhores métricas com o modelo treinado:

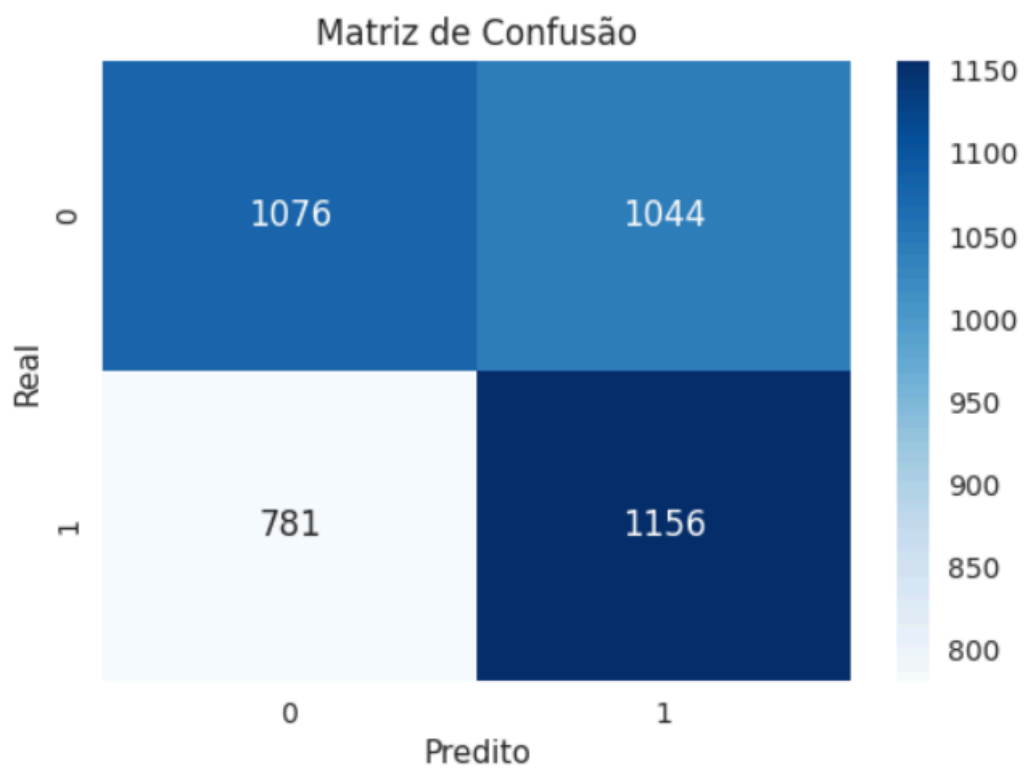
Regressão Logística:

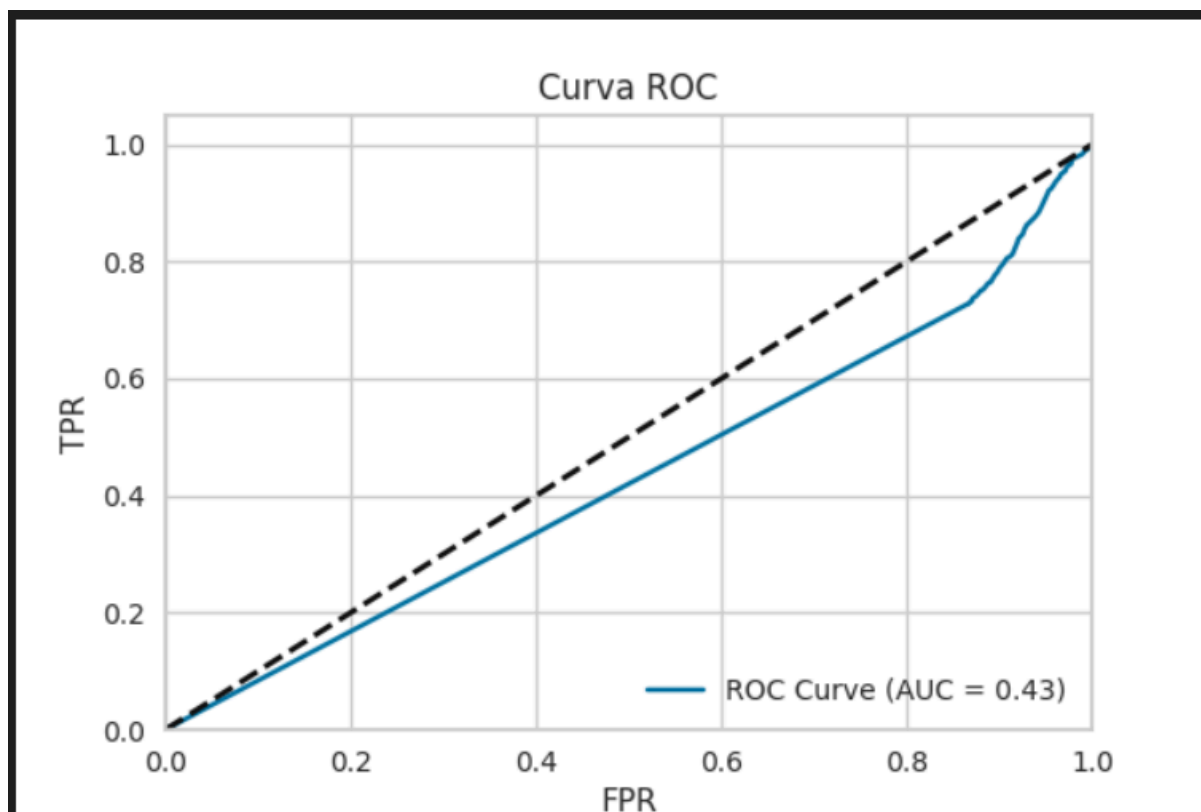
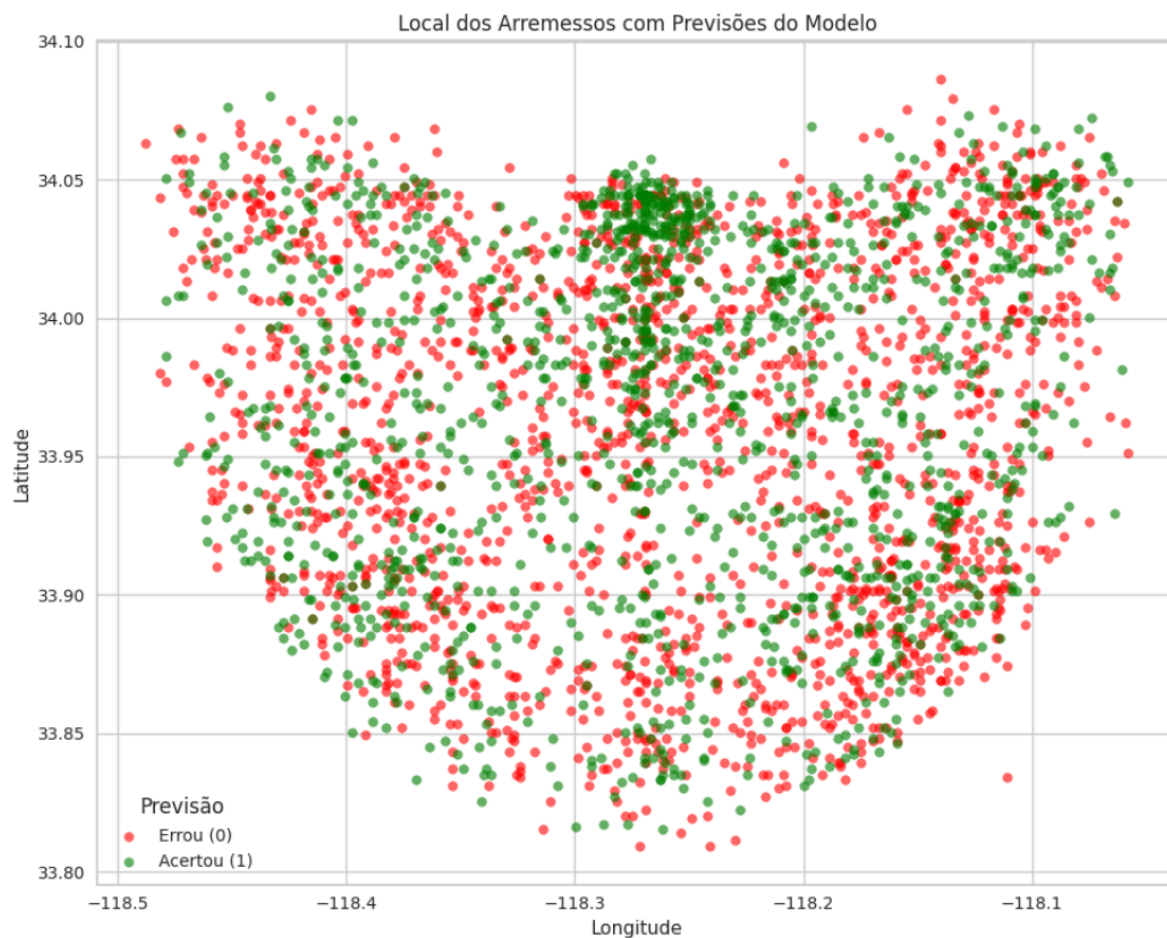






Decision Tree:





Análise:

Desta forma temos que **Árvore de Decisão: AUC = 0.43 e Regressão Logística: AUC = 0.51.**

Claramente a Regressão Logística performa melhor em termos de treinamento do que a Árvore de Decisão, tanto na comparação da curva RoC quanto na Matriz de Confusão, a regressão se mostra superior. Tanto em termos de balanceamento, quanto em menor indicativo de **overfitting**.

Em resumo, a Regressão Logística apresentou:

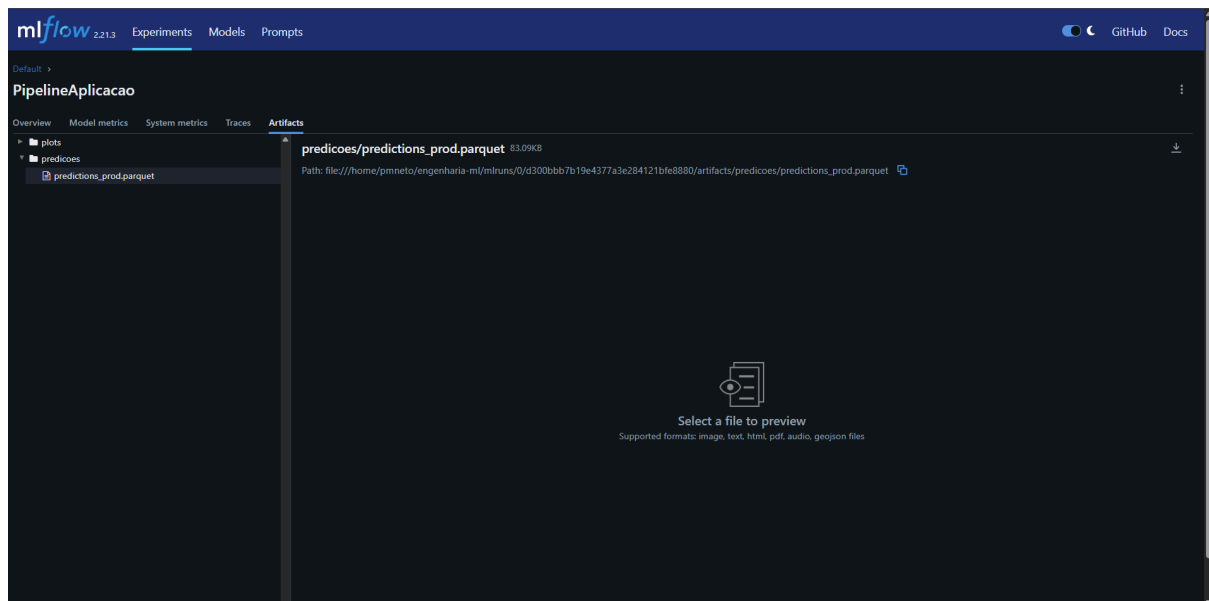
- Melhor **AUC**
- Melhor **distribuição de acertos reais**
- Melhor equilíbrio entre **falsos negativos e positivos**
- **Generalização mais robusta** no mapa de predições

Mesmo que ambos tenham performances medianas, a árvore de decisão claramente sofre de sobreajuste ao treinamento.

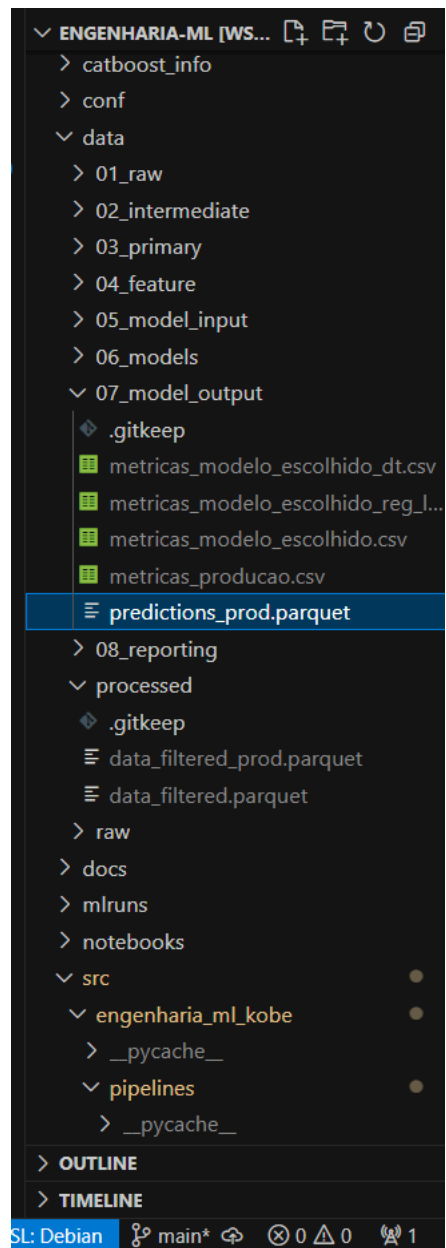
- Registre o modelo de classificação e o sirva através do MLFlow (ou como uma API local, ou embarcando o modelo na aplicação).

Desenvolva um pipeline de aplicação (`aplicacao.py`) para carregar a base de produção (`/data/raw/dataset_kobe_prod.parquet`) e aplicar o modelo. Nomeie a rodada (run) do mlflow como “PipelineAplicacao” e publique, tanto uma tabela com os resultados obtidos (artefato como `.parquet`), quanto log as métricas do novo log loss e `f1_score` do modelo.

Artefato com saída dos dados do modelo:



artefato predictions_prod.parquet




```
1 import requests
2 import json
3 import pandas as pd
4
5 url = "http://127.0.0.1:1234/invocations"
6
7
8 headers = {
9     "Content-Type": "application/json"
10 }
11
12 df_input = pd.DataFrame([
13     {"lat": 34.02,
14      "lng": -118.25,
15      "period": 2,
16      "minutes_remaining": 5,
17      "shot_distance": 15,
18      "playoffs": 1
19 }])
20
21
22 data = {
23     "inputs": df_input.to_dict(orient="records")
24 }
25
26
27 response = requests.post(url, headers=headers, data=json.dumps(data))
28
29
30 print("Status:", response.status_code)
31 print("Resposta:", response.json())
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

```
(engenhariaml-wsl) (engenhariaml-wsl) \[e]0;u@h: \w[a]\[\033[01;22m\]u@h\[\033[00m\]:\[\033[01;34m\]w\[\033[00m\]$ /mnt/c/Users/pmet/Documents/Projetos/pos_graduacao/Engenharia_ML_1_tr1/engenharia-ml/engenhariaml-wsl/bin/python3.11 /home/pmeto/engenharia-ml/tests/test_model_serving_mlflow.py
Status: 200
Resposta: {'predictions': [0.0]}
(engenhariaml-wsl) (engenhariaml-wsl) \[e]0;u@h: \w[a]\[\033[01;32m\]u@h\[\033[00m\]:\[\033[01;34m\]w\[\033[00m\]$
```

Ln 14, Col 10 Spaces: 4 UTF-8 CRLF () Python 3.11.2 (engenhariaml-wsl: venv)

Testando servir como api atraves do mlflow.

- a. O modelo é aderente a essa nova base? O que mudou entre uma base e outra? Justifique.

Não, o perfil dos dados de produção é diferente. São registros com lat e lng diferentes das usadas no treino. No predict de produção, o modelo erra quase tudo.

- b. Descreva como podemos monitorar a saúde do modelo no cenário com e sem a disponibilidade da variável resposta para o modelo em operação.

É possível acompanhar métricas como F1, Log Loss, AUC no tempo e queda de performance diretamente. Tanto pelo MIFlow quanto por um dashboard interativo modelado em streamlit.

- c. Descreva as estratégias reativa e preditiva de retreinamento para o modelo em operação.

Reativa:

- Reavaliar periodicamente o modelo.
- Se F1 ou log loss degradarem acima de um limiar (ex.: 15%), disparar o retreinamento.

Preditiva:

- Implementar alertas de mudança de distribuição nas features (ex.: KS-test).
- Antecipar retreinamento antes que o desempenho piore.

- Implemente um dashboard de monitoramento da operação usando Streamlit.

