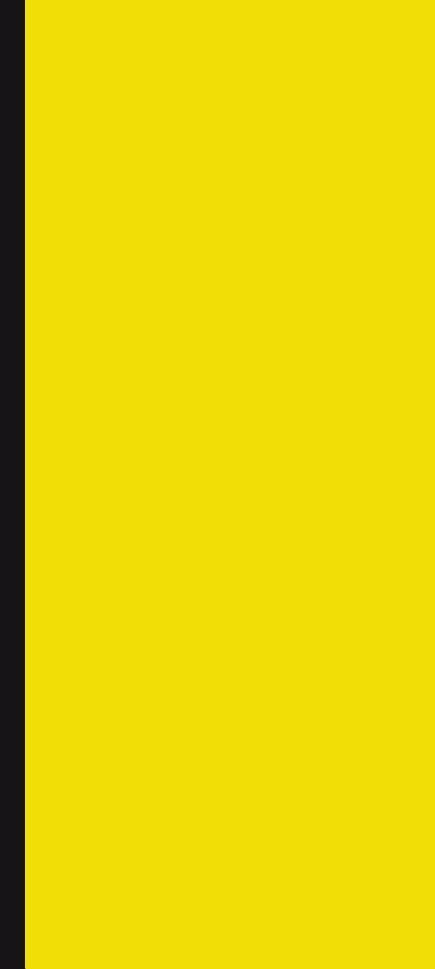




Data Science & Machine Learning



CRISP-DM Deployment

Consultor: Murilo Mendonça



Murilo Mendonça

ML Engineer
@Ambev Tech



Eng. Mecânico

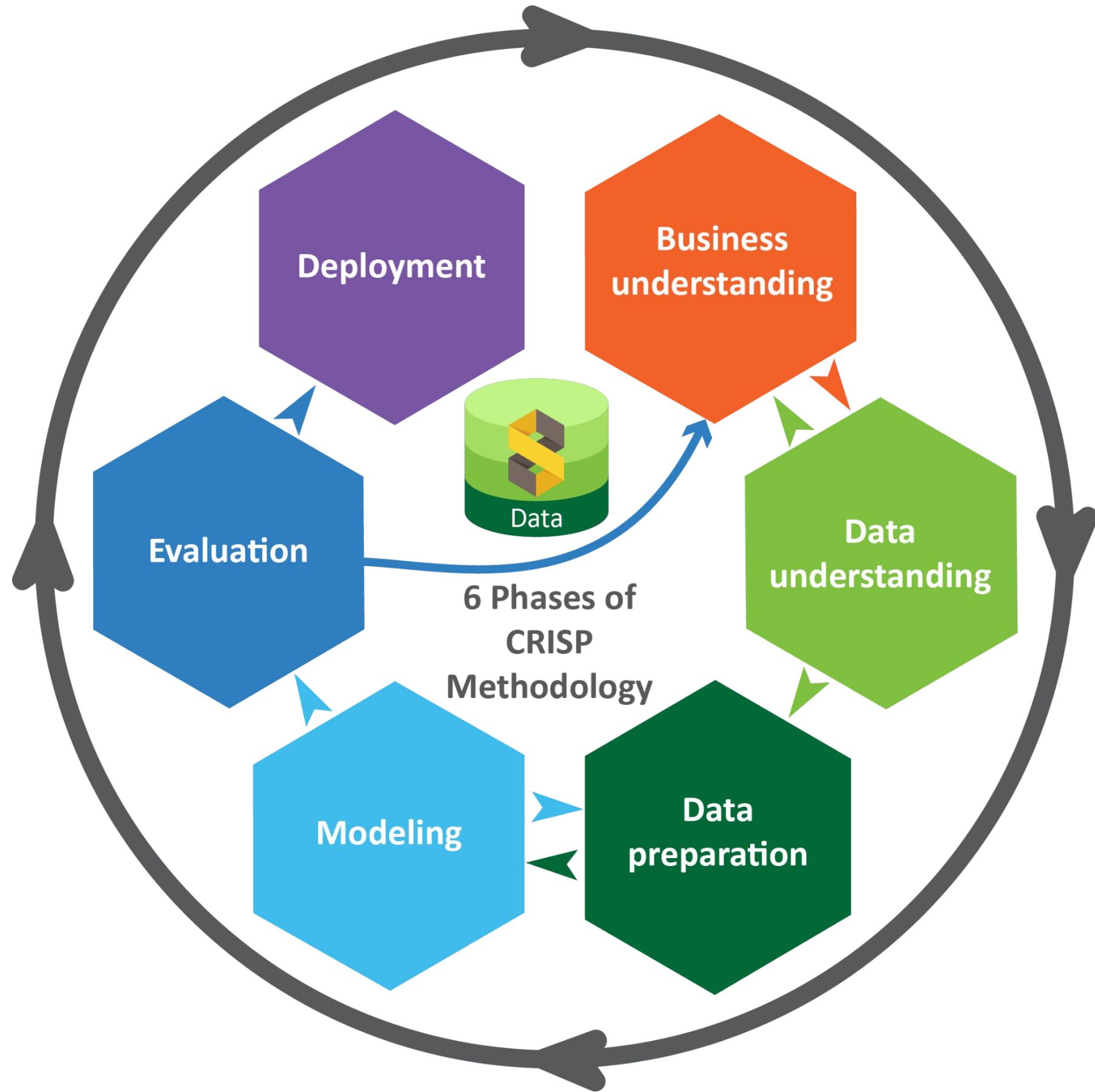
Fraunhofer Institute

Lactec

Ambev

Ambev Tech



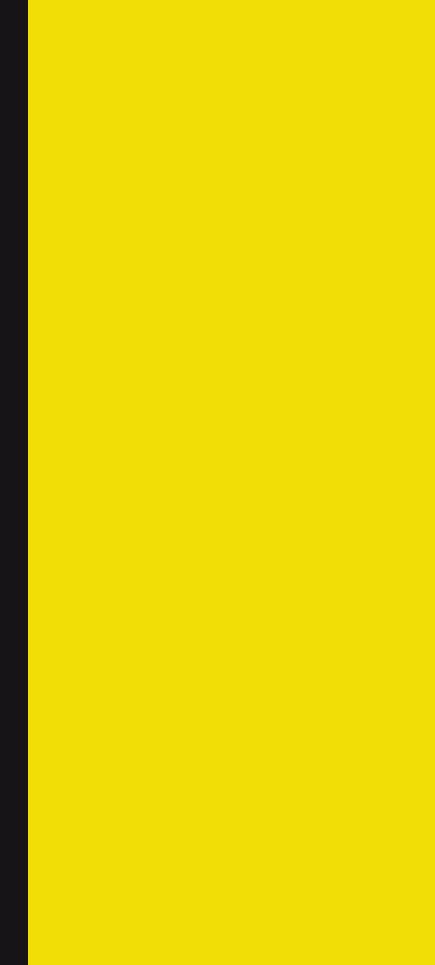


Tópicos do Módulo:





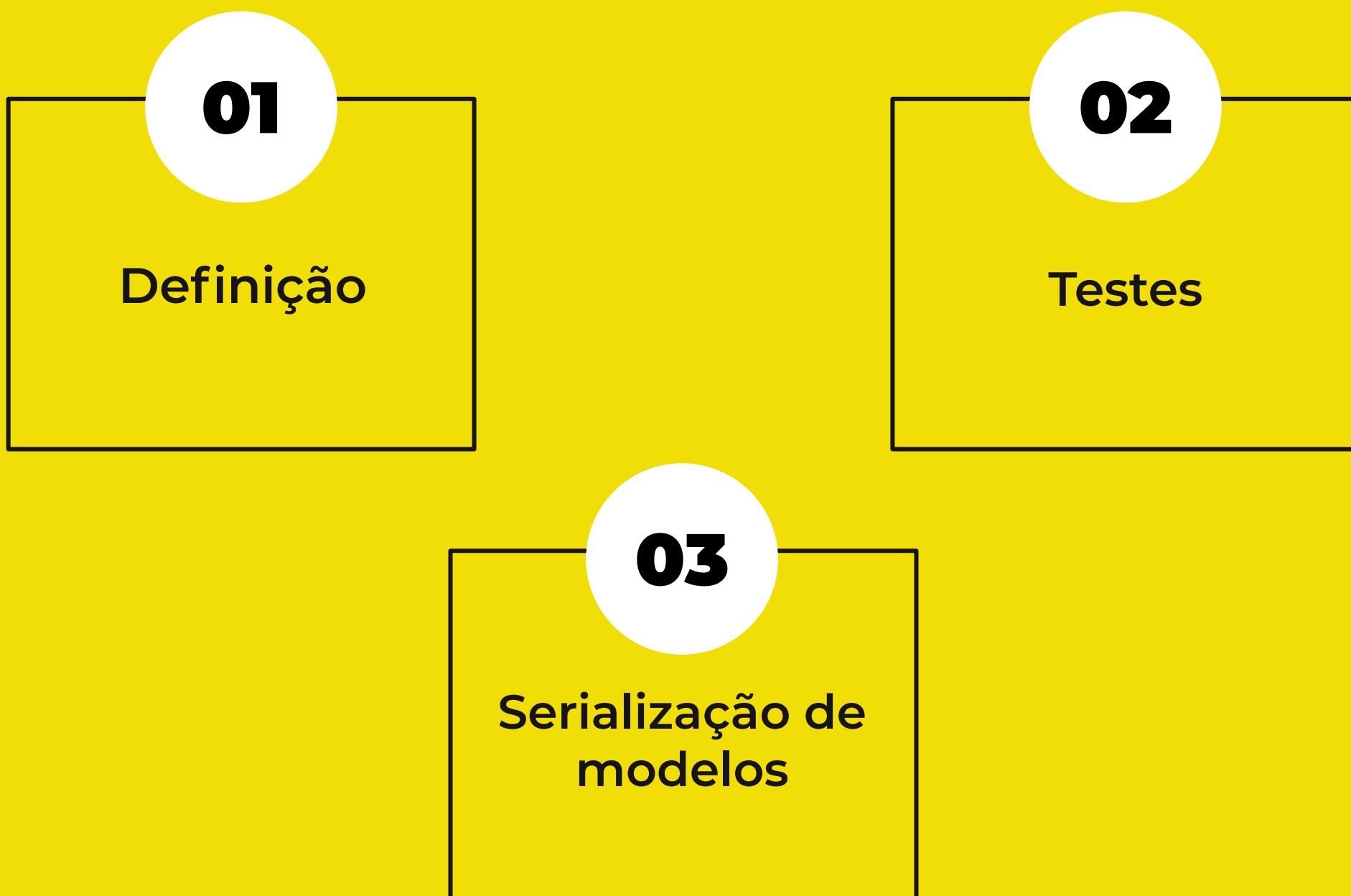
Data Science & Machine Learning



02. Ambiente de Produção

Consultor: Murilo Mendonça

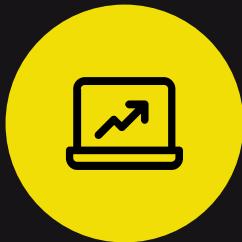
O que Veremos Nesta Aula:



O que é?

Usuário:
é o que ele enxerga e interage

Desenvolvedor:
é uma réplica do seu ambiente de trabalho

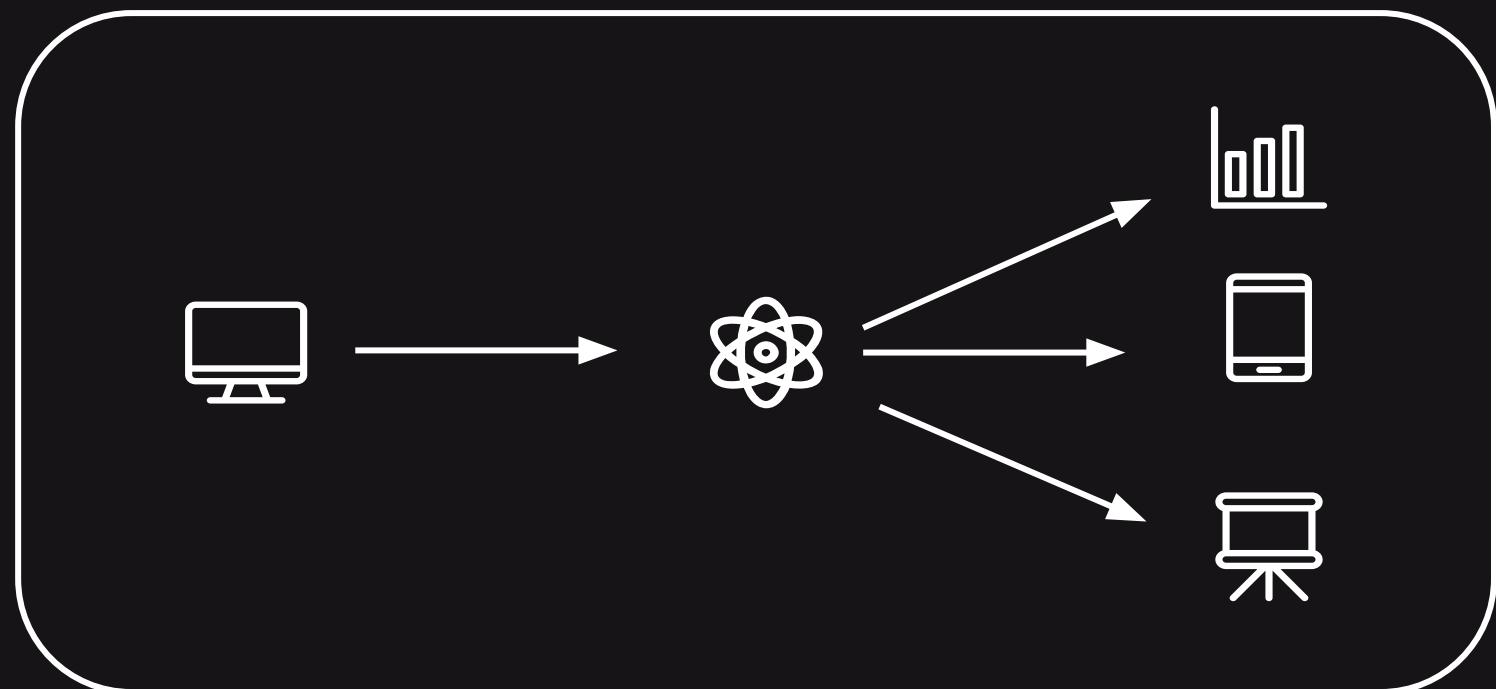
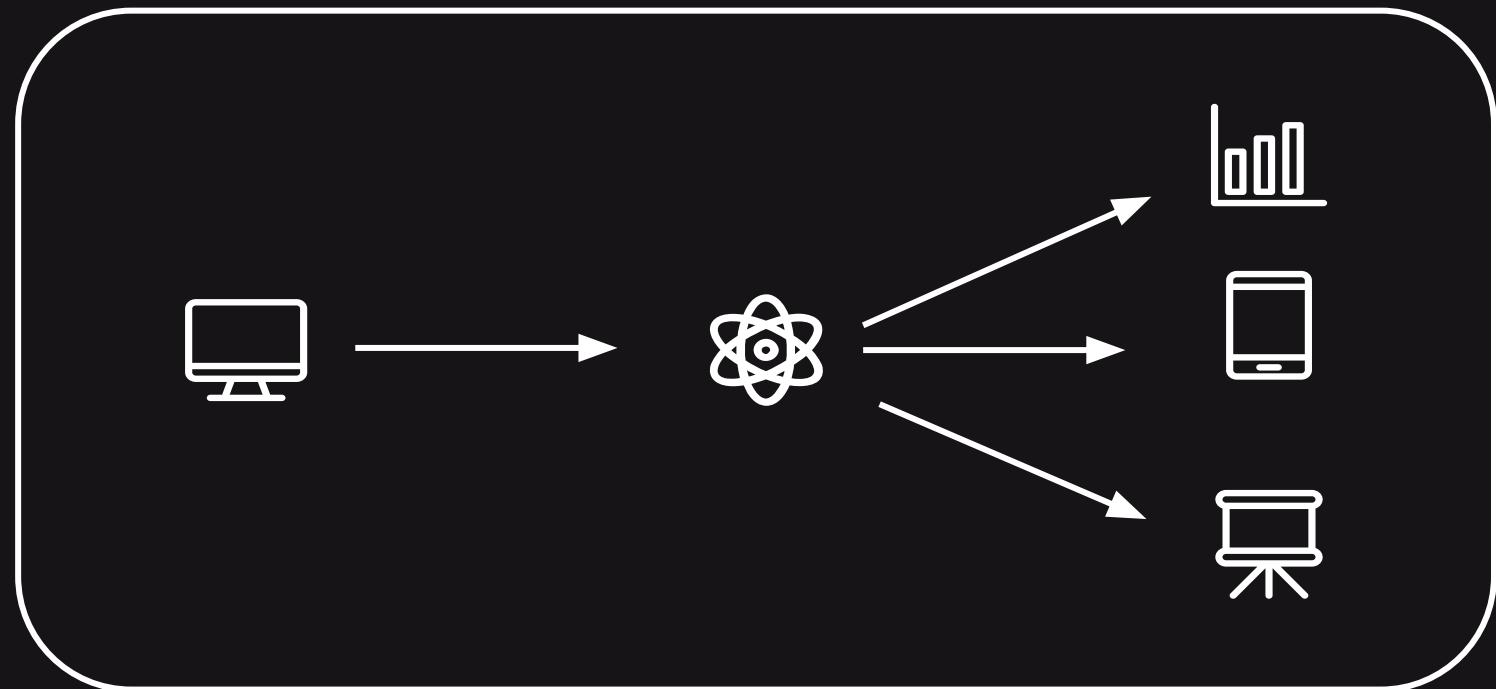


Como?

NONPROD



PROD



Por que?



Controle de Versão



Governança



Entregas Faseadas

Testes

Unitários

Integrados

Teste Unitário

```
def custom_sum(number):
    return number + 10

def test_custom_sum():
    number = 10

    actual_result = custom_sum(number)
    expected_result = 20

    assert isinstance(actual_result, int)
    assert actual_result == expected_result
```

Teste Unitário

```
===== test session starts =====
platform darwin -- Python 3.7.12, pytest-6.2.5, py-1.10.0, pluggy-1.0.0
rootdir: /Users/murilomendonca/Documents/repos/dnc-classes
collected 1 item

test_custom_sum.py .

===== 1 passed in 0.01s =====
```

Teste Integrado

```
✓ def test_load_model(connection_string):  
  
    loaded_model = MyClass().load_model(connection_string)  
  
    assert isinstance(loaded_model, mlflow.pyfunc.PythonModel)  
  
  
def test_make_predictions(loaded_model, data):  
  
    predictions = MyClass().make_predictions(loaded_model, data)  
  
    assert isinstance(predictions, list)  
    assert predictions[0] == 1
```

Serialização

"Ok, treinei meu modelo, as métricas estão ok. "

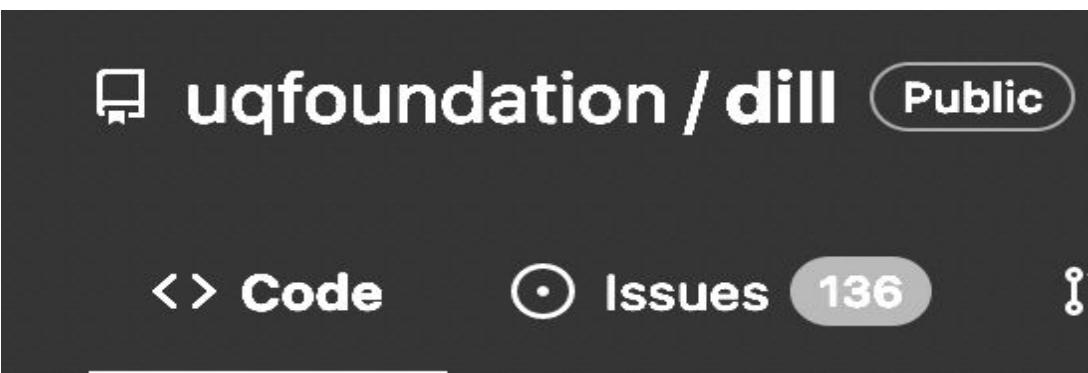
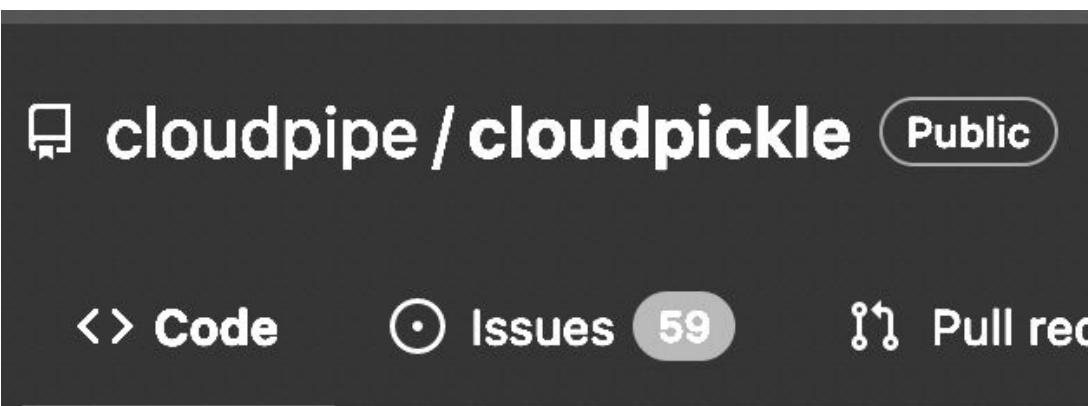
E agora?



Serialização

[pickle — Python object serialization](#)

Source code: [Lib/pickle.py](#)



Serialização

```
# Train model
model = DecisionTreeClassifier(max_depth=max_depth)
model.fit(X, y)

# Save model to disk
with open("tree_classifier.pkl", "wb") as file:
    pickle.dumps(model, file)
```



≡ tree_classifier.pkl

Serialização

```
sklearn.tree._classes??DecisionTreeClassifier??})??}??criterion??gini??splitter??best??max_depth??min_samples_split??K??min_samples_leaf??K??min_weight_fraction_leaf??G??max_features??N??max_leaf_nodes??N??random_state??N??min_impurity_decrease??G??min_impurity_split??N??class_weight??N??ccp_alpha??G??n_features_in_??K??n_features_??K??n_outputs_??K??classes_??numpy.core.multiarray??_reconstruct??numpy??ndarray??K??Cb??R??(KK??h??dtype??i8??R??(K??<?NNNJ????J????Kt??b??C??t??b??n_classes_??h??scalar??h?'C??R??max_features_??K??tree_??sklearn.tree._tree??Tree??KhhK??h??R??(KK??h?'C??t??bK??R??}??(h??K??node_count??K??nodes??hhK??h??R??(KK??h$??V56????R??(K??|??N??left_child??right_child??feature??threshold??impurity??n_node_samples??weighted_n_node_samples??t??}??(hJh$??i8????R??(Kh(NNNJ????J????Kt??bK??hKhUK??hLhUK??hMh$??f8????R??(Kh(NNNJ????J????Kt??bK??hNh\K ??h0hUK(??hPh\K0??uK8KKt??b??B??@VUUUUU@??b@?????????????????????????????2I@??dY@????????????????????h@WH??6K@?????????????????????????????????<b??p?.Ga@t@b@values??hhK??h??R??(KKKK??h\?CxI@I@I@I@I@H@?F@?t@bub?_sklearn_version??0.24.2@ub.
```

Serialização

≡ tree_classifier.pkl



```
def load_model():
    with open("tree_classifier.pkl", "rb") as file:
        model = pickle.load(file)
    return model
```

Onde Armazenar?

File System

Cloud Storage

Docker image

Banco de dados



Resumo

Definição

Testes

Serialização





Data Science & Machine Learning

03. Planejando o Deployment (BATCH)

Consultor: Murilo Mendonça

O que Veremos nesta Aula:

01

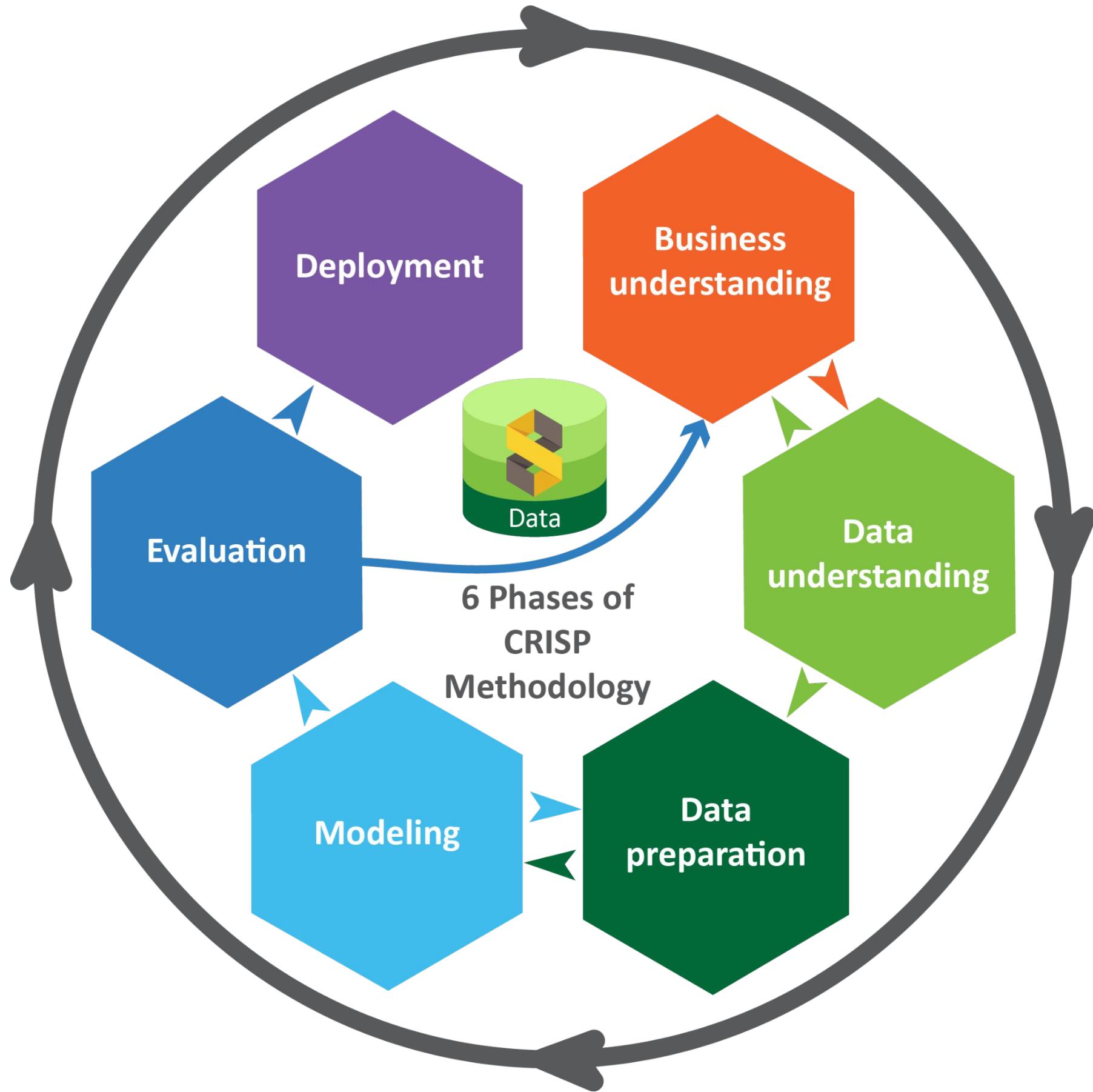
Pipelines

02

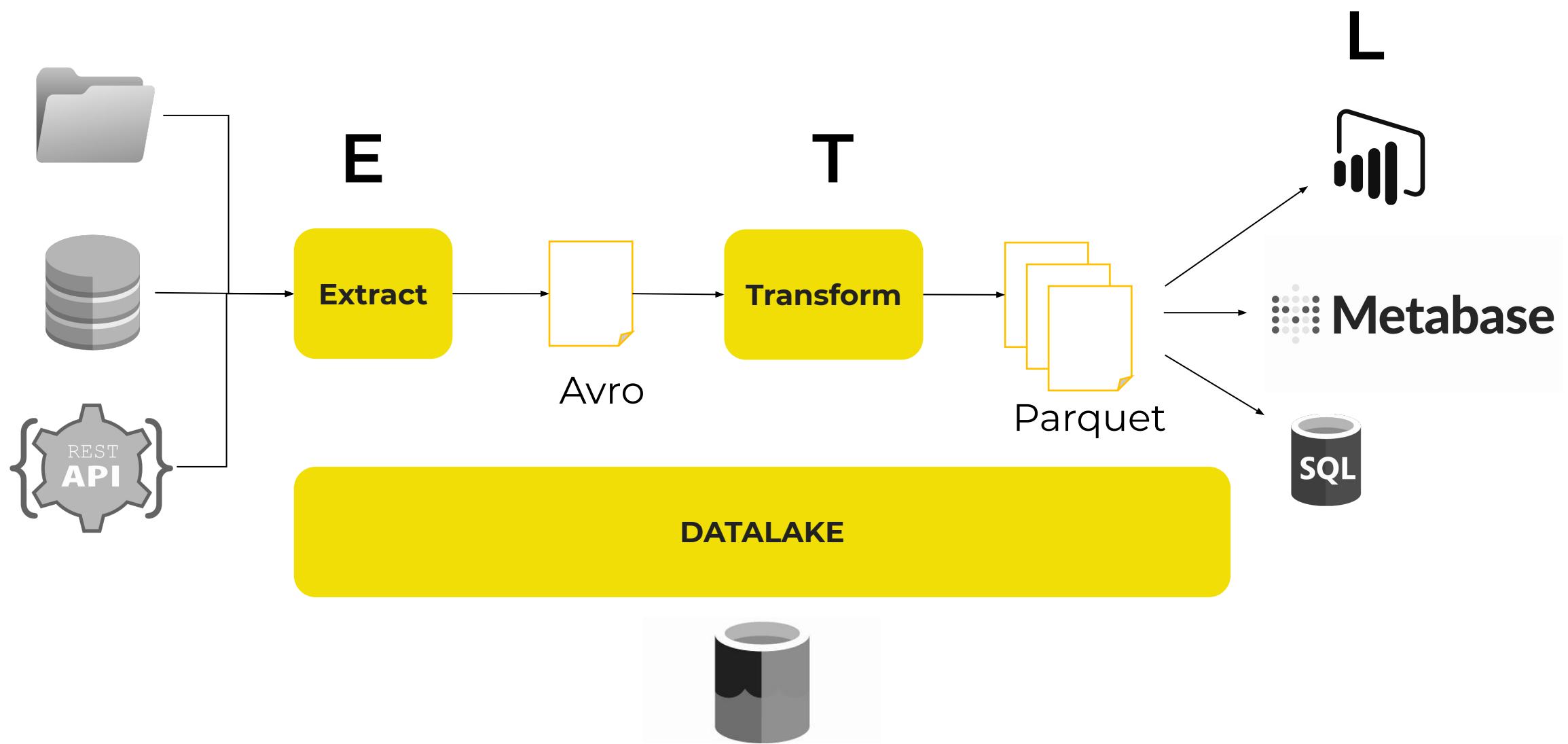
Spark + Mlflow

03

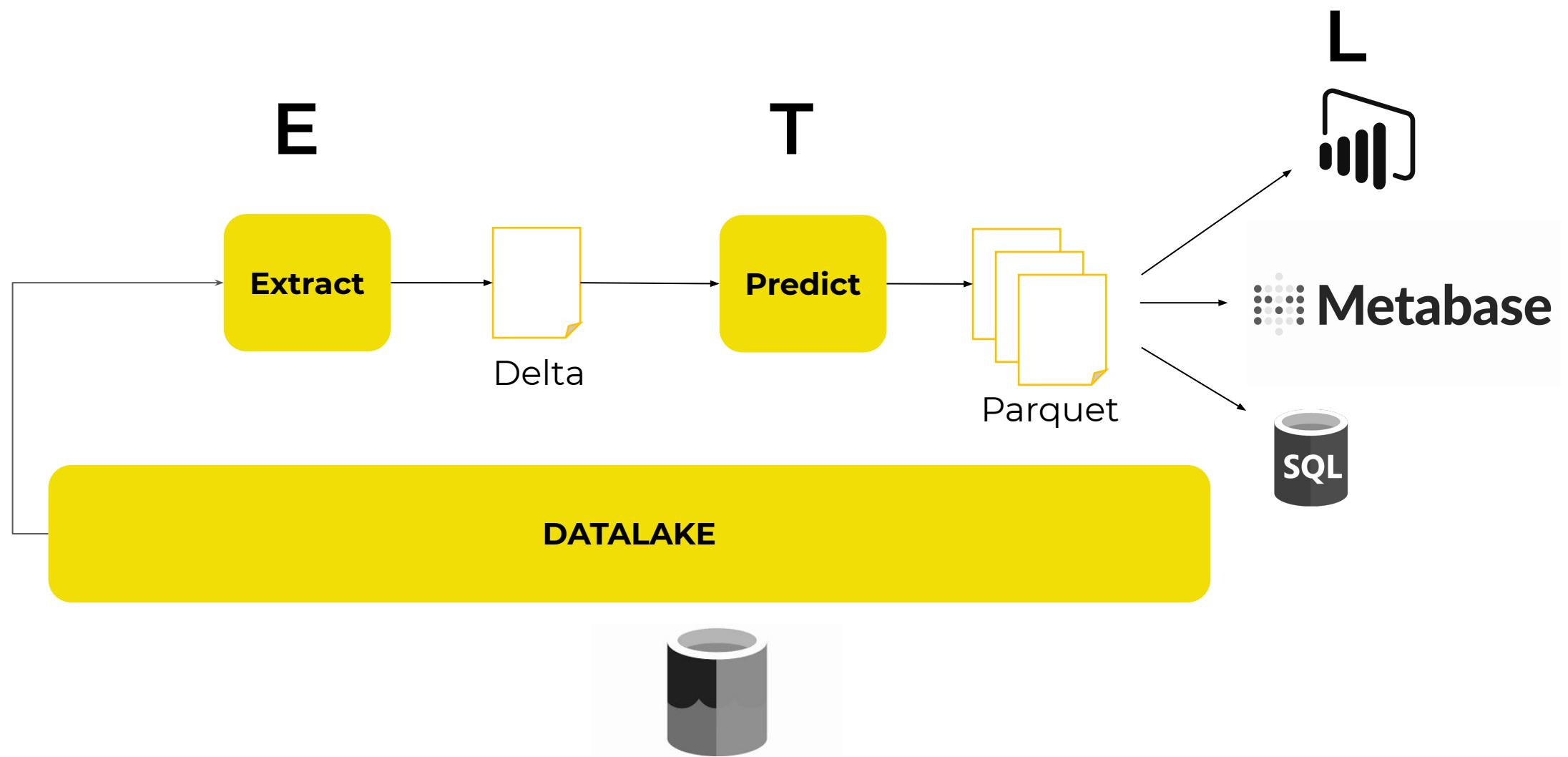
Periodicidade e
Orquestração



Pipeline de Dados



Pipeline ML



MLflow + Spark



Ferramenta de
análise,
treinamento e
processamento de
dados em escala



Plataforma de
gerenciamento do
ciclo de vida de
modelos

MLflow + Spark

```
logged_model = 'runs:/<run_id>/model'

# Load model as a Spark UDF.
loaded_model = mlflow.pyfunc.spark_udf(spark, model_uri=logged_model)

# Predict on a Spark DataFrame.
columns = list(df.columns)
df.withColumn('predictions', loaded_model(*columns)).collect()
```

Periodicidade

De quanto em quanto tempo?

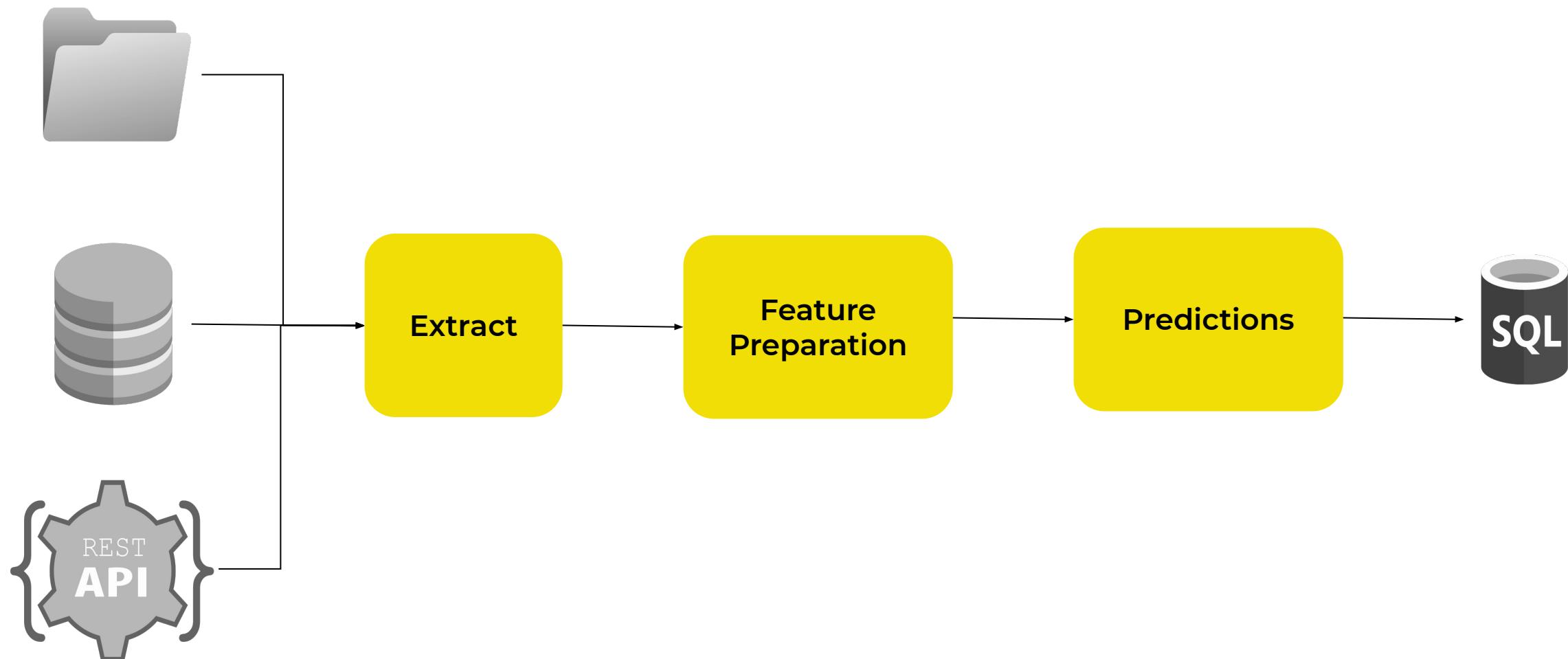
Onde?

Quais são as etapas do seu pipeline?



DAG

Directed Acyclic Graph



Orquestração

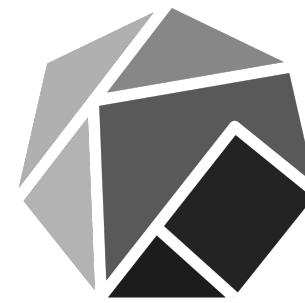
Código

Crontab

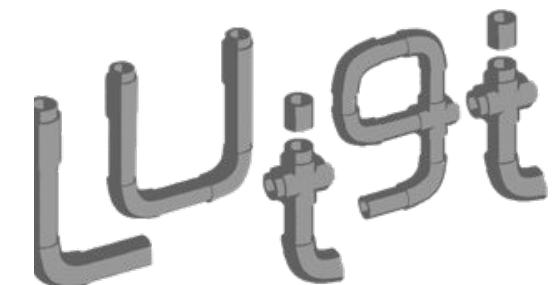
Ferramentas de mercado



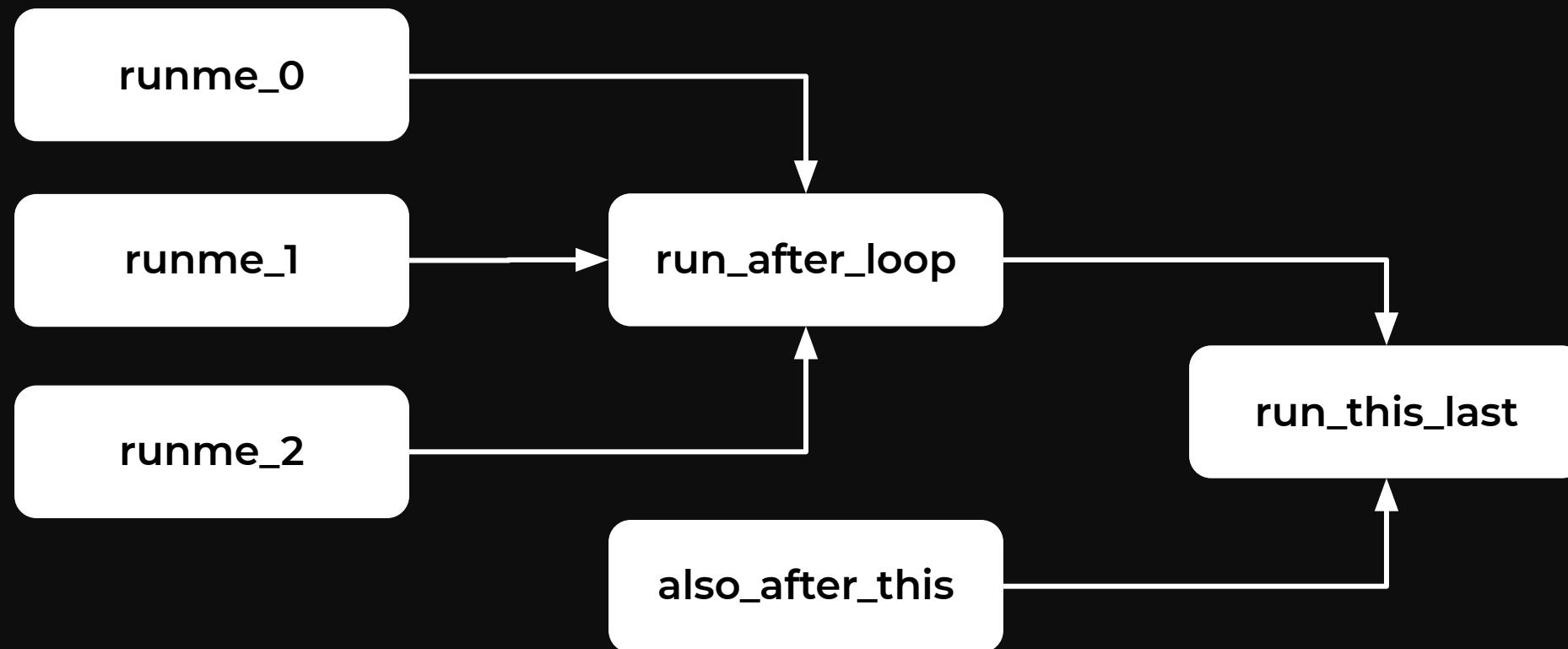
Orquestração



Kubeflow



Orquestração



Resumo





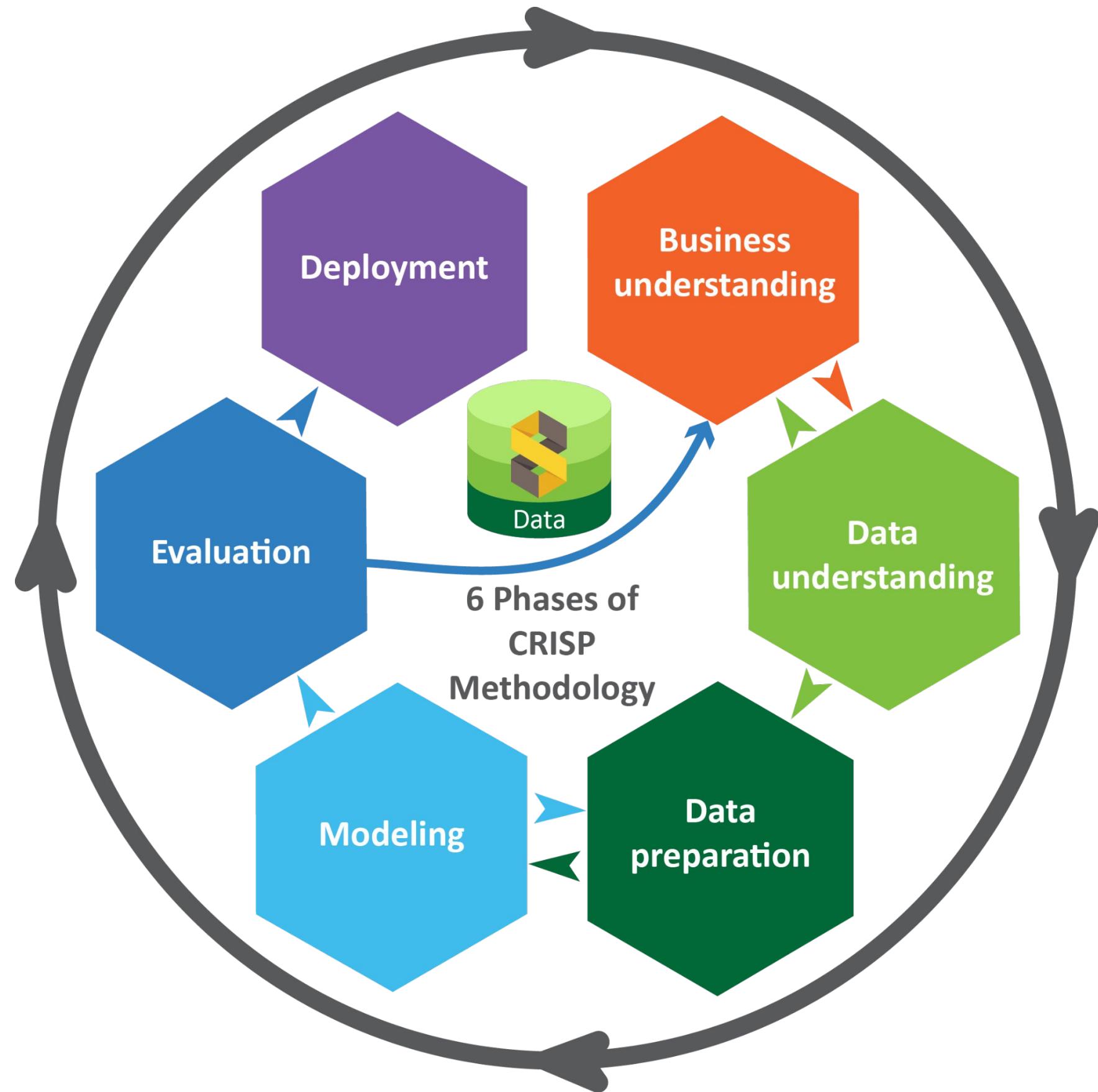
Data Science & Machine Learning

04. Planejando o Deployment (API)

Consultor: Murilo Mendonça

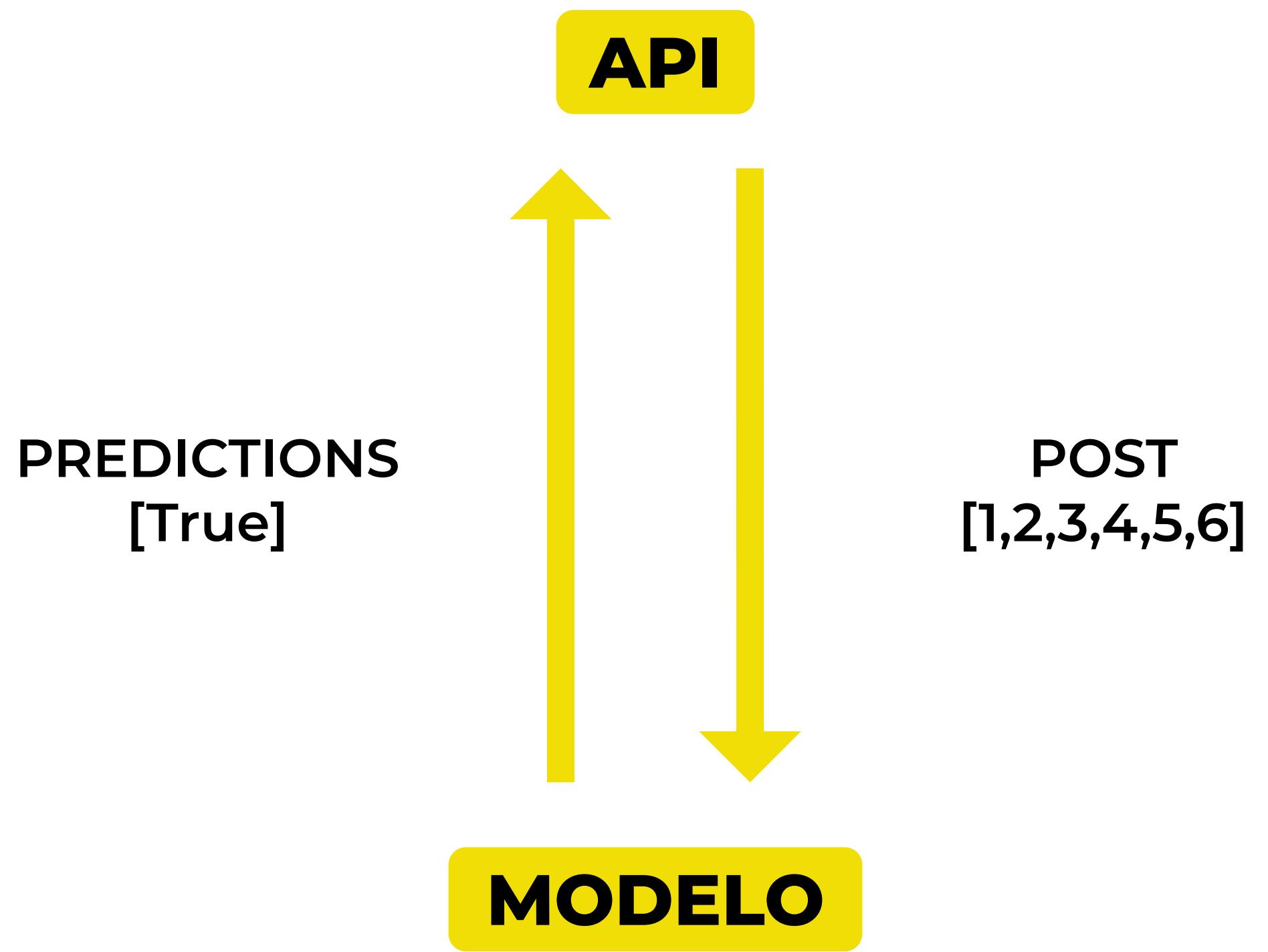
O que Veremos nesta Aula:





API





```
1  {
2      "model": "iris_flower_prediction",
3      "data": {
4          "SEPAL_LENGTH": [5.3],
5          "SEPATL_WIDTH": [6.7],
6          "PETAL_LENGTH": [2.1],
7          "PETAL_WIDTH": [2.3]
8      }
9  }
```

Body ▾



200 OK

3.72 s

256 B

Pretty

Raw

Preview

Visualize

JSON ▾



```
1  {
2      "predictions": [
3          2
4      ]
5  }
```

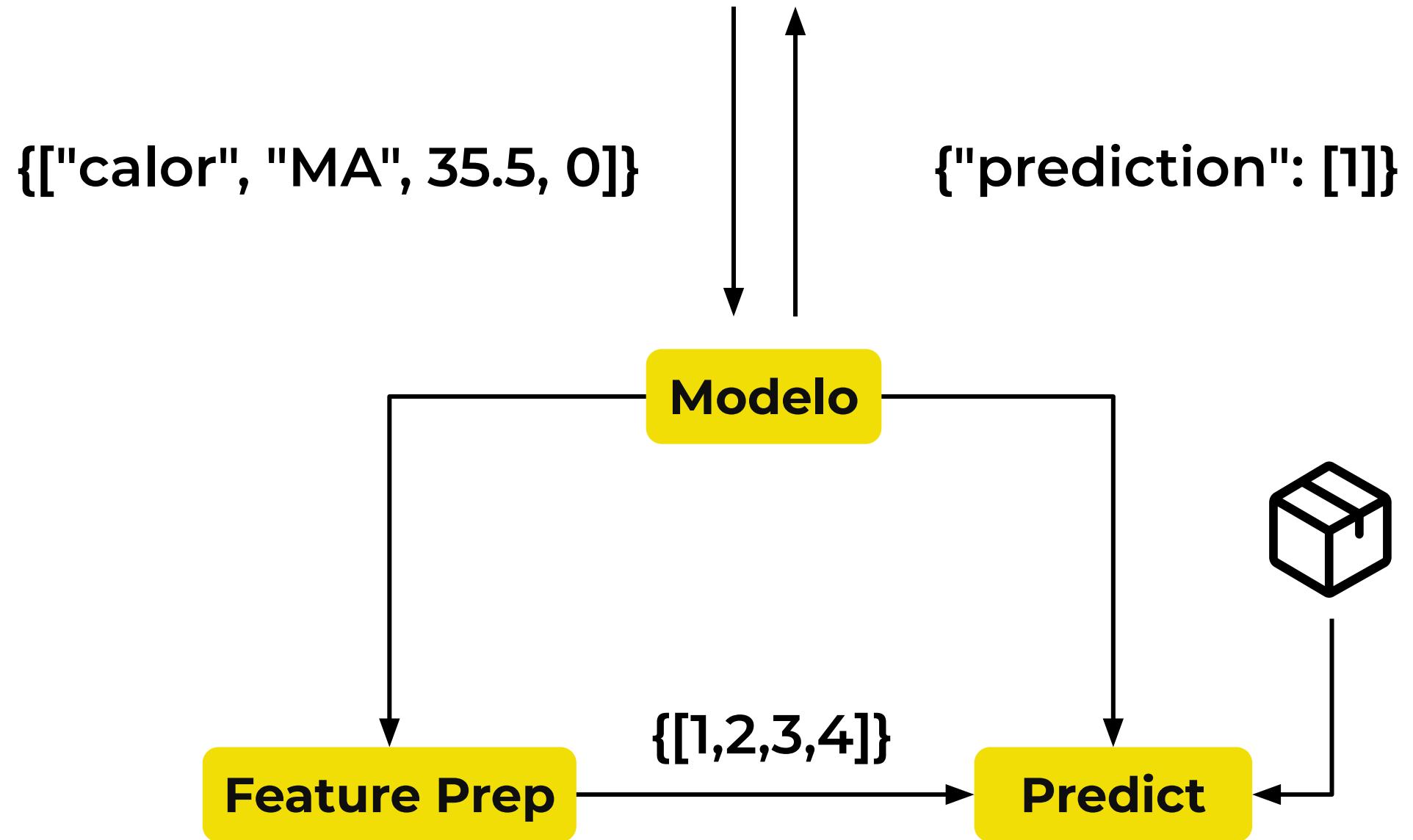
Complexidade

- Intermittência da aplicação
- Contrato
- Documentação
- Autenticação
- Infraestrutura
- Performance



É de uma API
que você
precisa?

Mini-batch



Ambiente

Orçamento

Trade-off

Requerimento de projeto



Ambiente de Execução

- Máquina virtual
- Serviços gerenciados
 - Databricks
 - Sagemaker, Azure ML, etc.
 - Algorithmia, Seldon, etc.
- Serverless
 - Functions
 - Container Instances
- Swarm ou Kubernetes

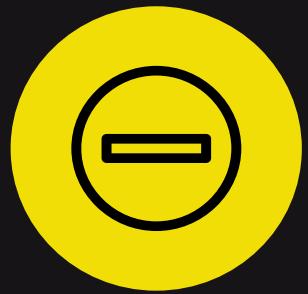
Métricas de App



Latência



Tráfego



Erros



Saturação

Métricas de ML

Definir as métricas de sucesso (F1, RMSE, etc.)

Como medir?

Métricas de negócio (Rotação, tempo médio, etc.)

Qualidade de dados



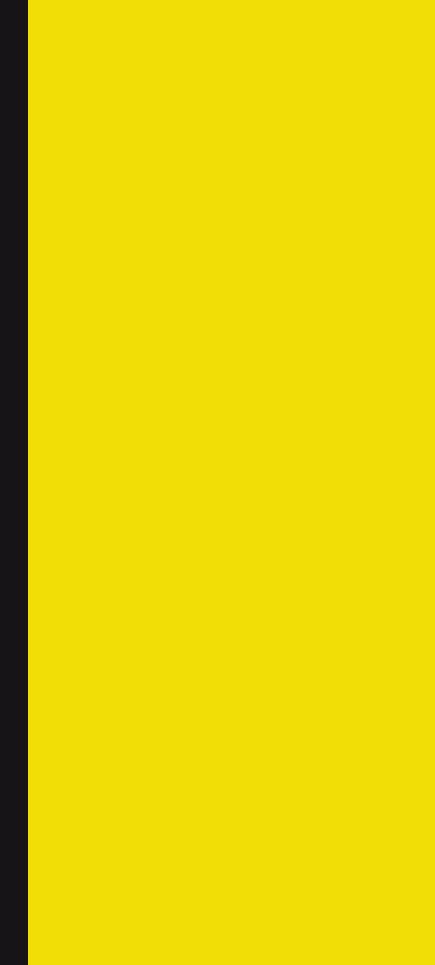
Resumo

- API
- Complexidade
- Mini-Batch
- Ambiente de execução
- Métricas





Data Science & Machine Learning



05. Planejando o Monitoramento

Consultor: Murilo Mendonça

O que Veremos nesta Aula:

01

Métricas de ML

02

Métricas de
aplicação

03

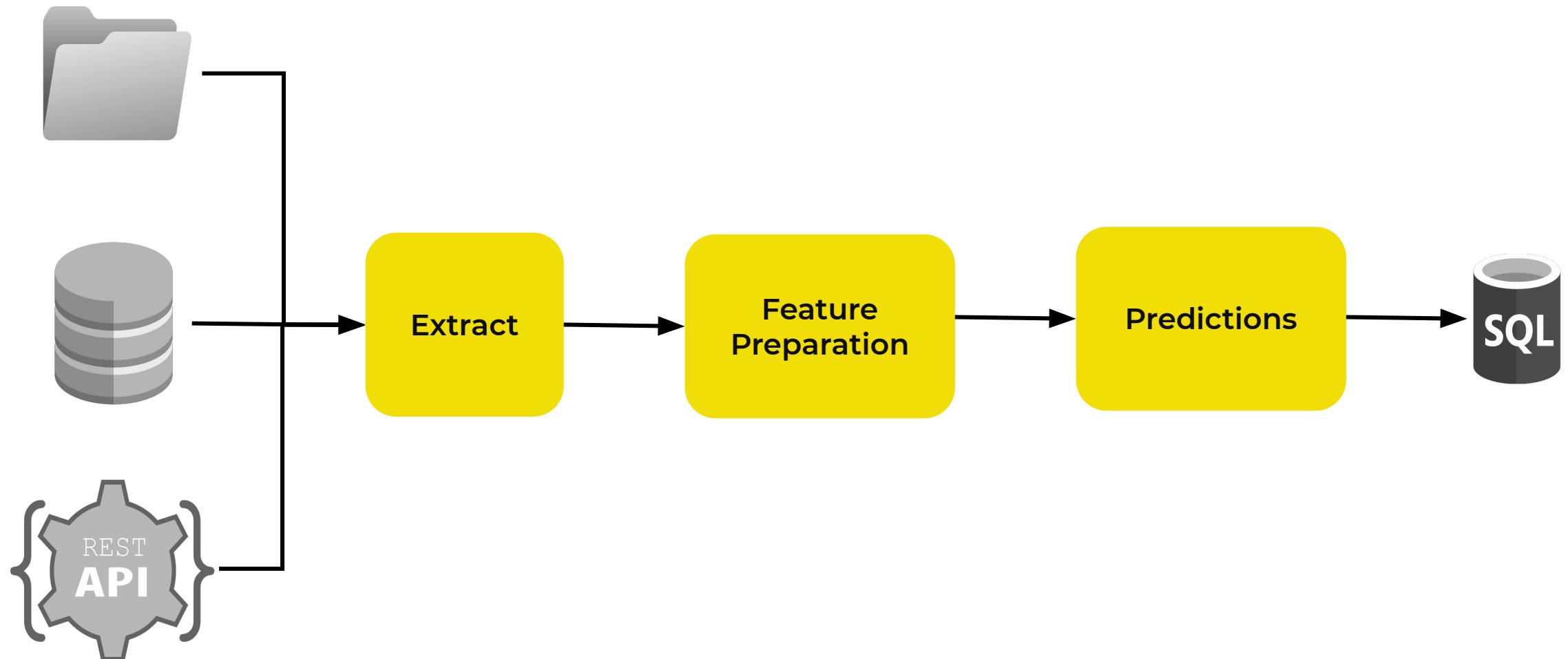
Ferramentas de
mercado

04

Caso real

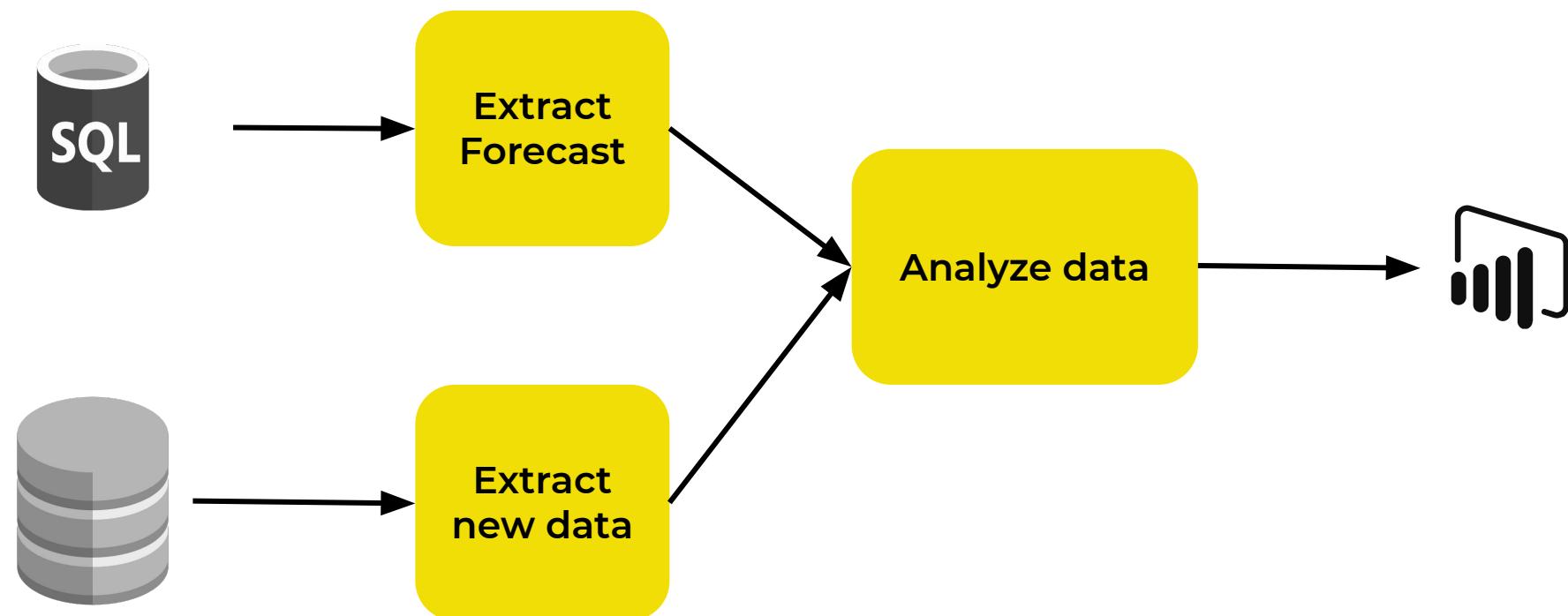
Métricas de ML

Prediction DAG

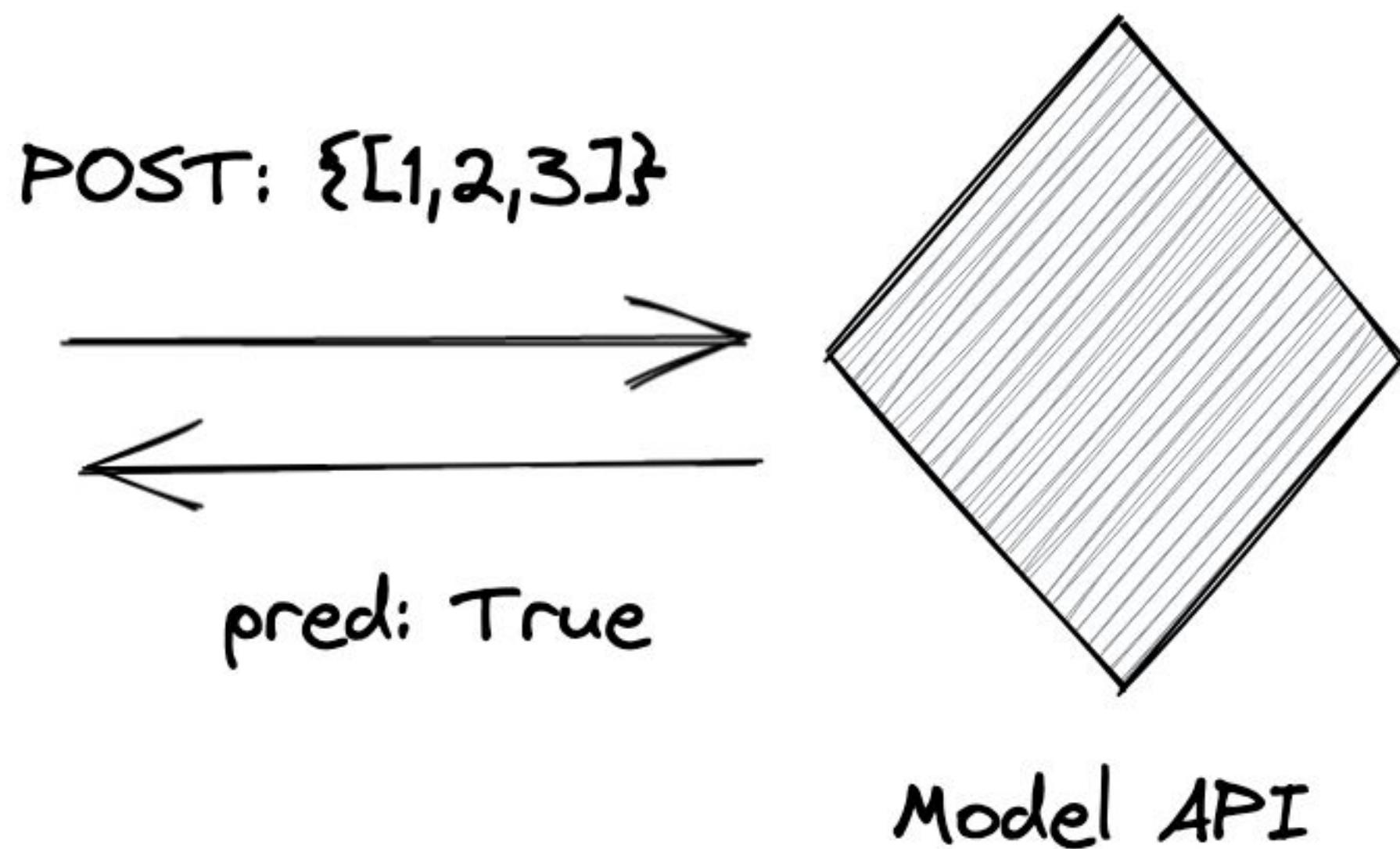


Métricas de ML

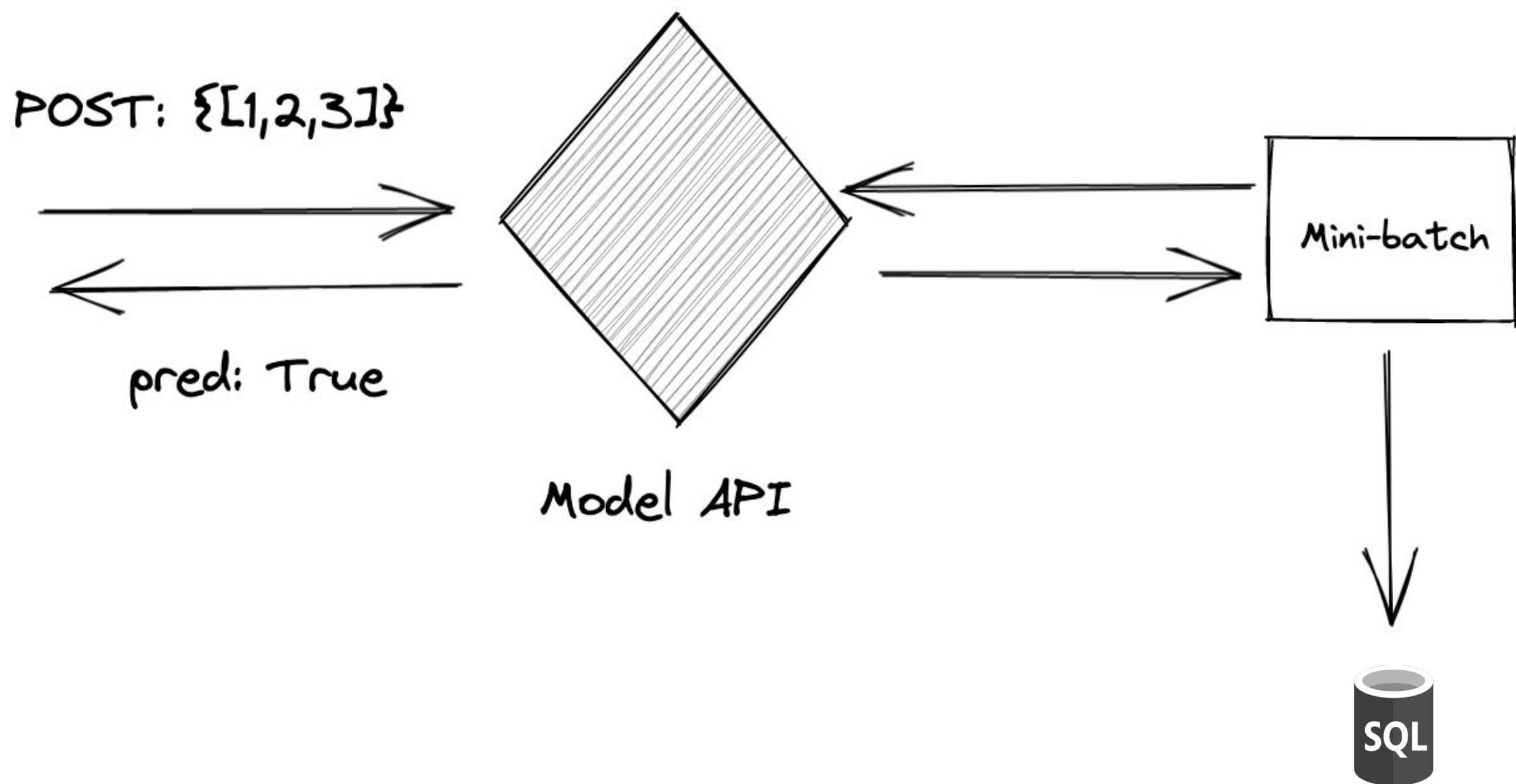
Metrics DAG



Métricas de ML



Métricas de ML



Métricas de App

Tivemos um pico de latência às 4:00am.

O que mais aconteceu?



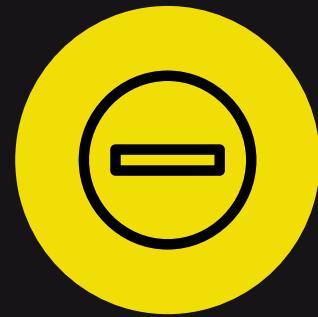
Métricas de App



Latência



Tráfego



Erros



Saturação

Métricas de App



Latência

- 500ms
- Timeout
- Percentil

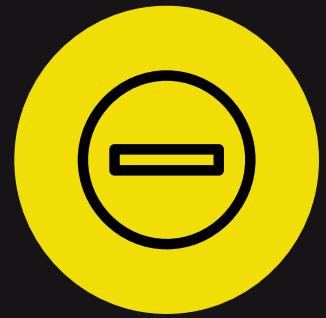
Métricas de App



Tráfego

- Request
- Request/Second
- Request/Response

Métricas de App



Erros

- 3xx - Redirect
- 4xx - Bad Request
- 5xx - Internal Error

Métricas de App



Saturação

- CPU Utilization
- Memory Usage
- Available pods, nodes etc.

Métricas de App

Google Site Reliability Engineering

Home

Spotlight

Resources 

Books

Mobaa

Classroom

Careers

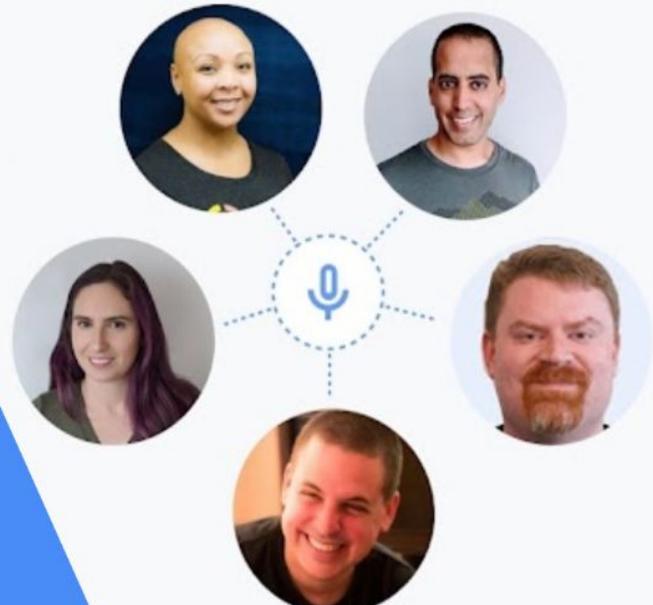
What is Site Reliability Engineering (SRE)?

SRE is what you get when you treat operations as if it's a software problem. Our mission is to protect, provide for, and progress the software and systems behind all of Google's public services — Google Search, Ads, Gmail, Android, YouTube, and App Engine, to name just a few — with an ever-watchful eye on their availability, latency, performance, and capacity.

SRE Experts

Hear veteran Googlers describe their experiences as SREs: how their backgrounds led them to their current roles, and what their day-to-day work looks like.

[Read more →](#)



< >

Métricas de App



Service Level Objetive



Service Level Agreement



Service Level Indicator

Métricas de App

Availability

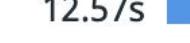
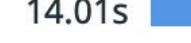
Past 7 Days

100.00%

100% (31.5k reqs) budget

96% target

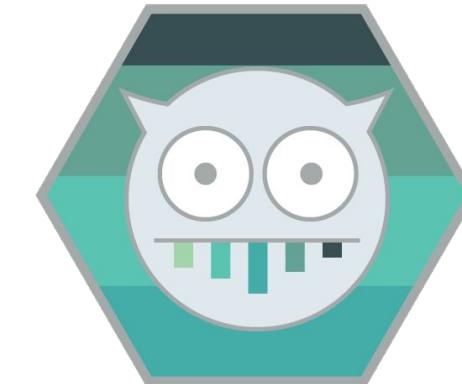
Métricas de App

RESPONSE CODE	↓ COUNT	AVG:DURATION	PC95:DURATION	PC99:DURATION
200	995 	4.32s 	7.31s 	12.57s 
200	91 	8.86s 	14.01s 	22.43s 
200	36 	10.47s 	18.06s 	26.69s 
400	6 	521ms 	821ms 	821ms 

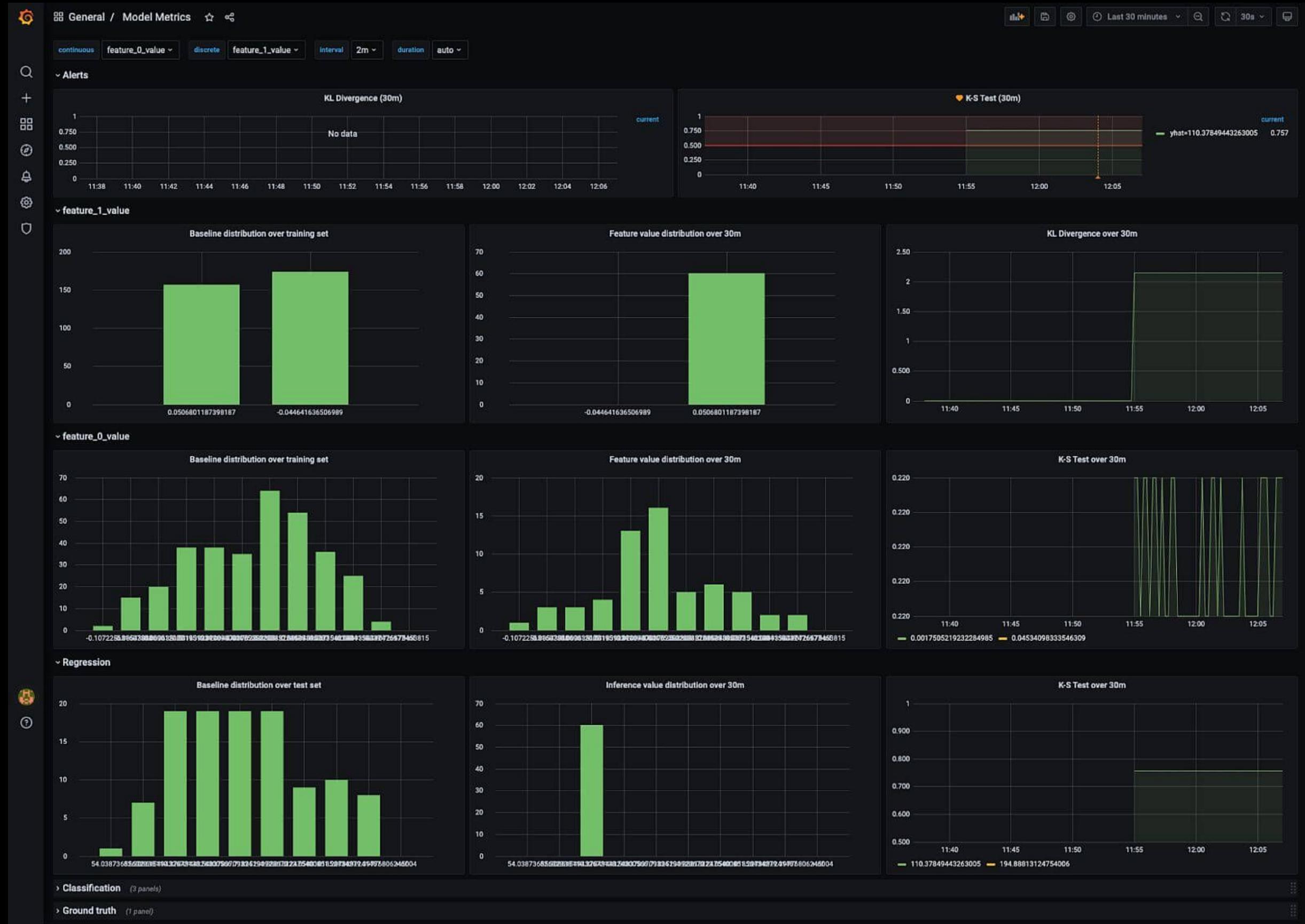
Ferramentas

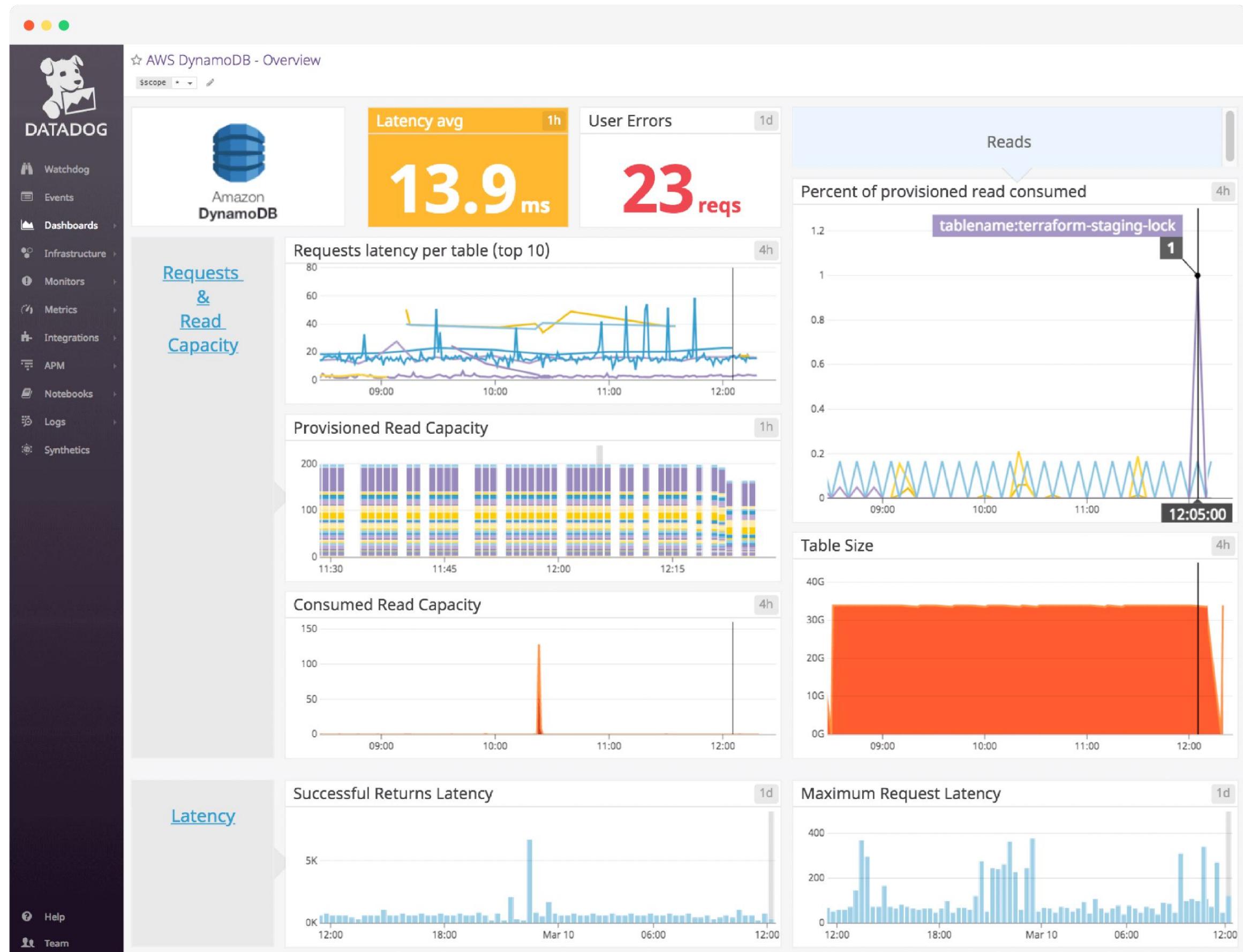


DATADOG



Application Insights

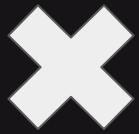




Ferramentas

TRADE-OFF

Orçamento



Requerimento
de projeto



Antes



Depois



Resumo

- Métricas de ML
- Métricas de Aplicação
- Ferramentas de Mercado
- Caso Real





Data Science & Machine Learning

06. Planejando a Manutenção

Consultor: Murilo Mendonça

O que veremos nesta aula:



Manutenção

**Programada: Preciso retreinar o modelo de
série temporal toda semana**

Manutenção

**Não-programada: Preciso retreinar meu
modelo porque o negócio percebeu**

Manutenção

**Não-programada: Preciso retreinar meu
modelo porque recebi um alerta**

Manutenção

Programada: Vou subir uma nova versão da API com novas features

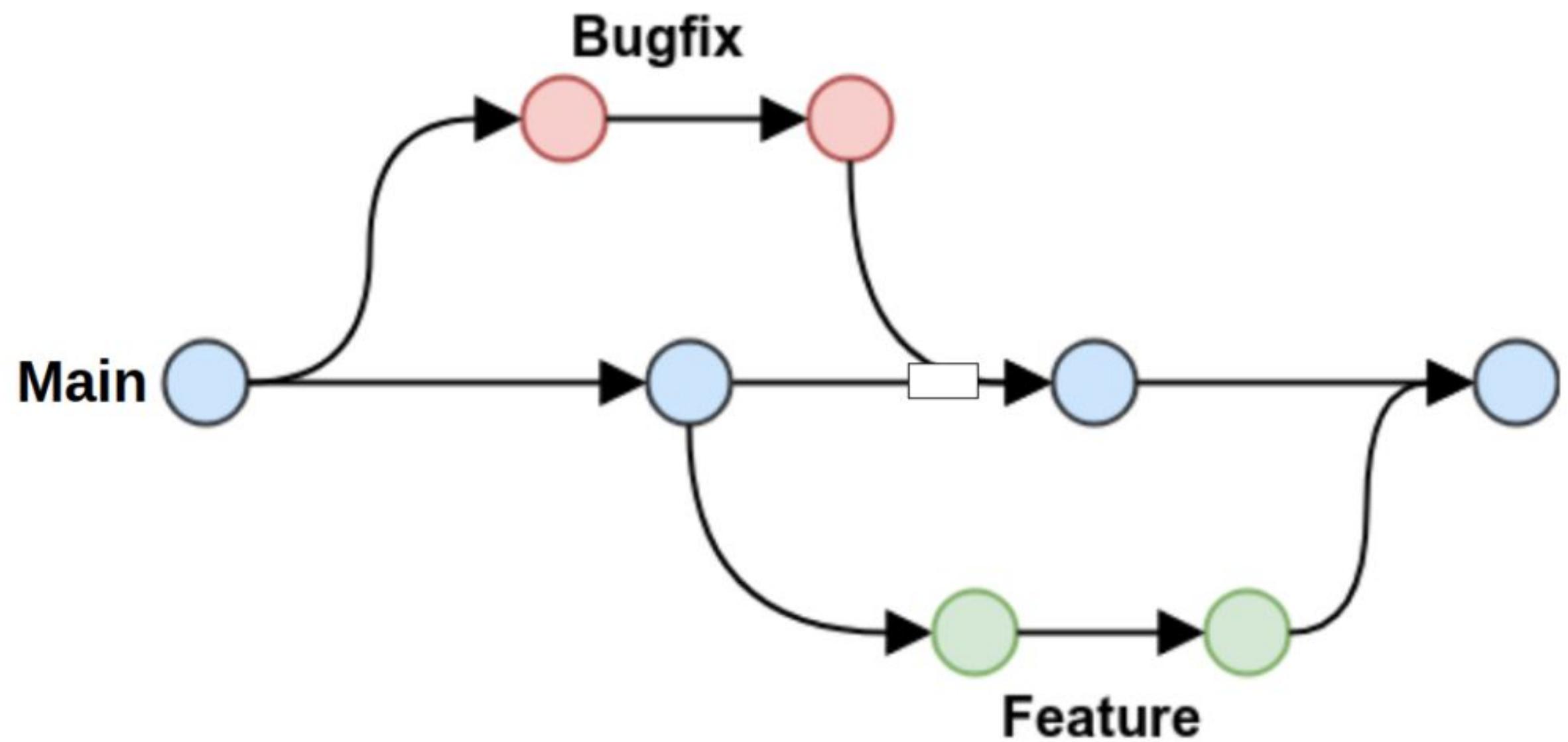
Manutenção

**Não-Programada: Minha API está retornando
erro 500 em 10% das vezes**

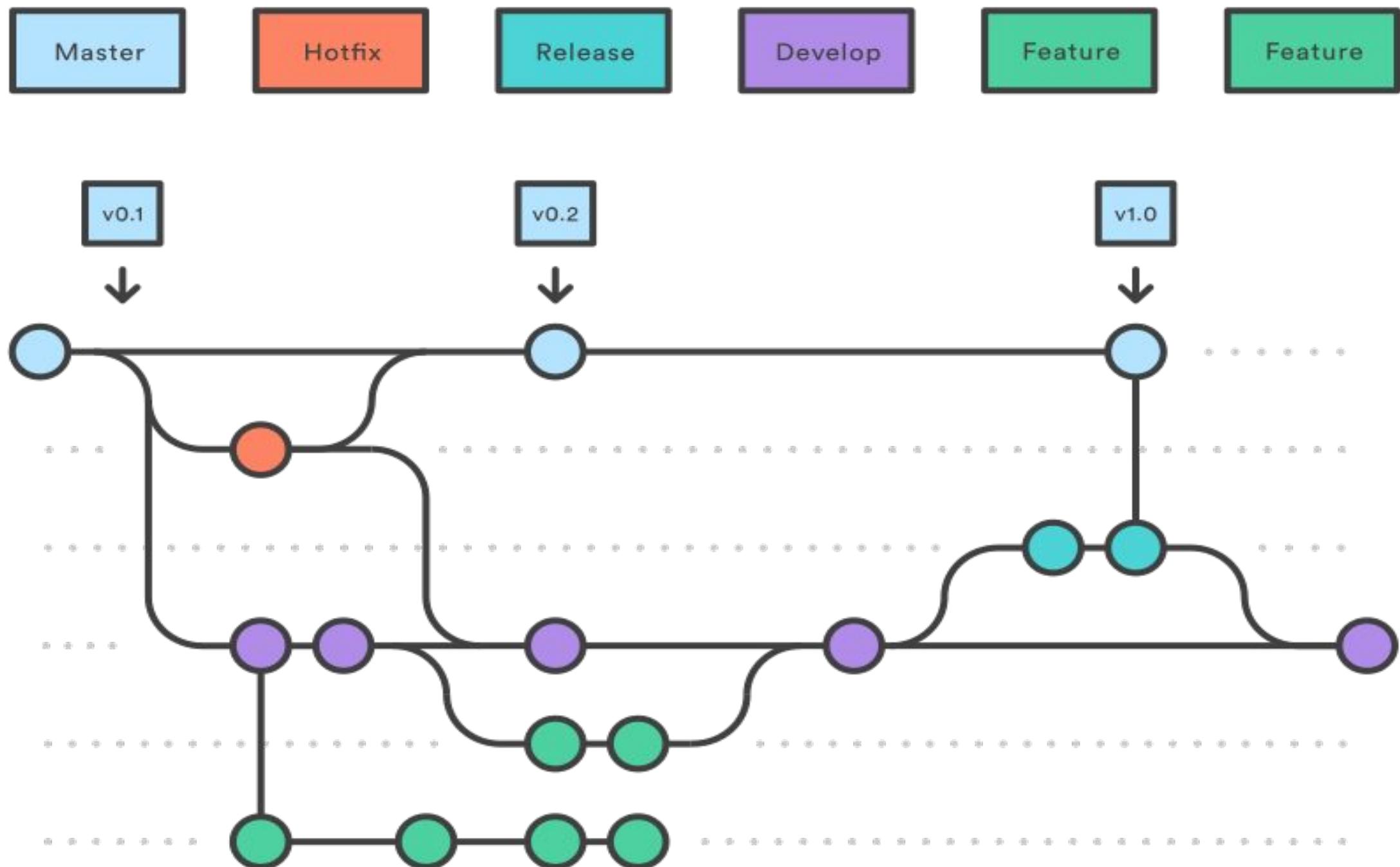
Repositório

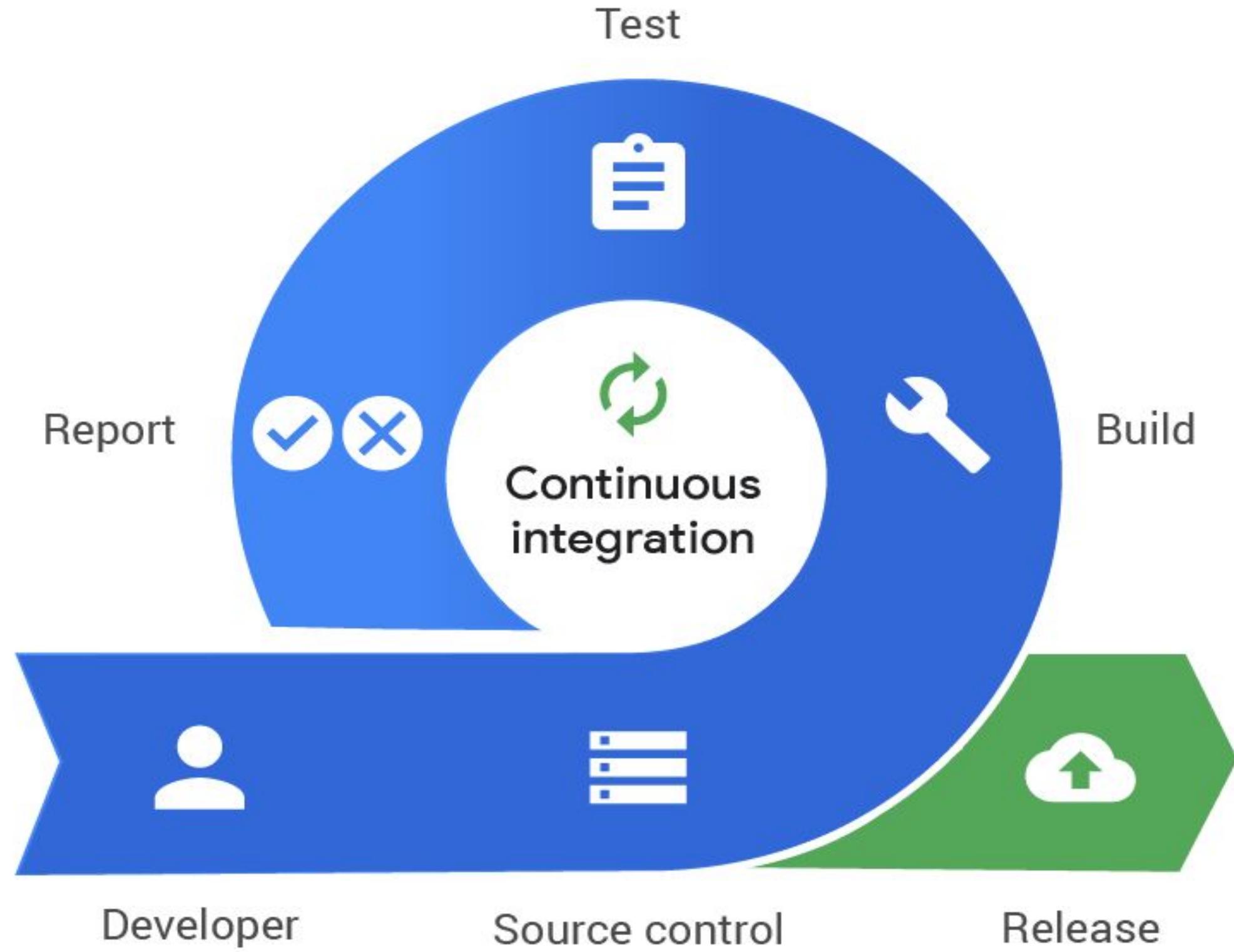
	analysis
	data
	docs
	log
	src
	tests
	.gitignore
	Dockerfile
	README.md
	requirements.txt

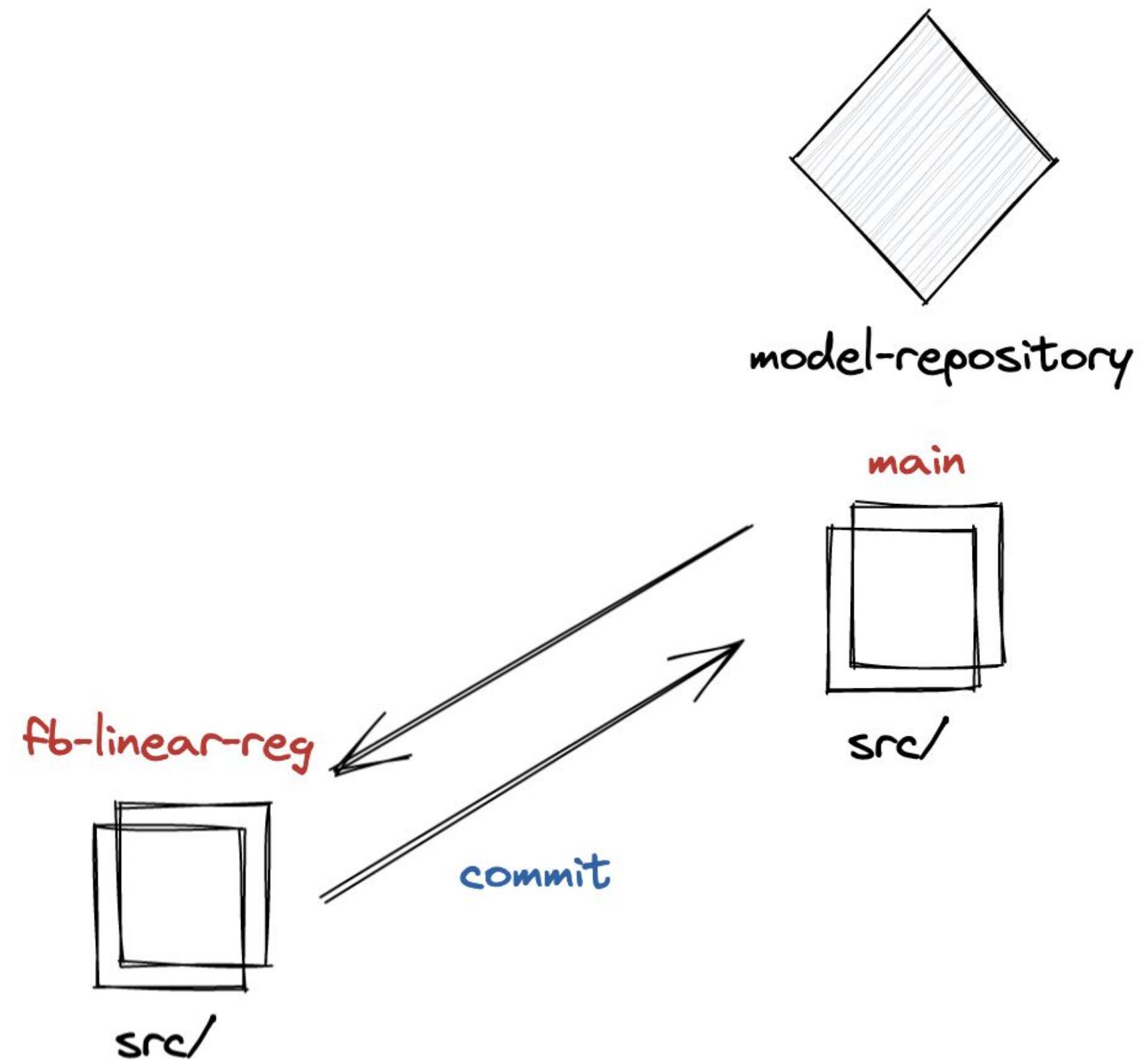
Repositório

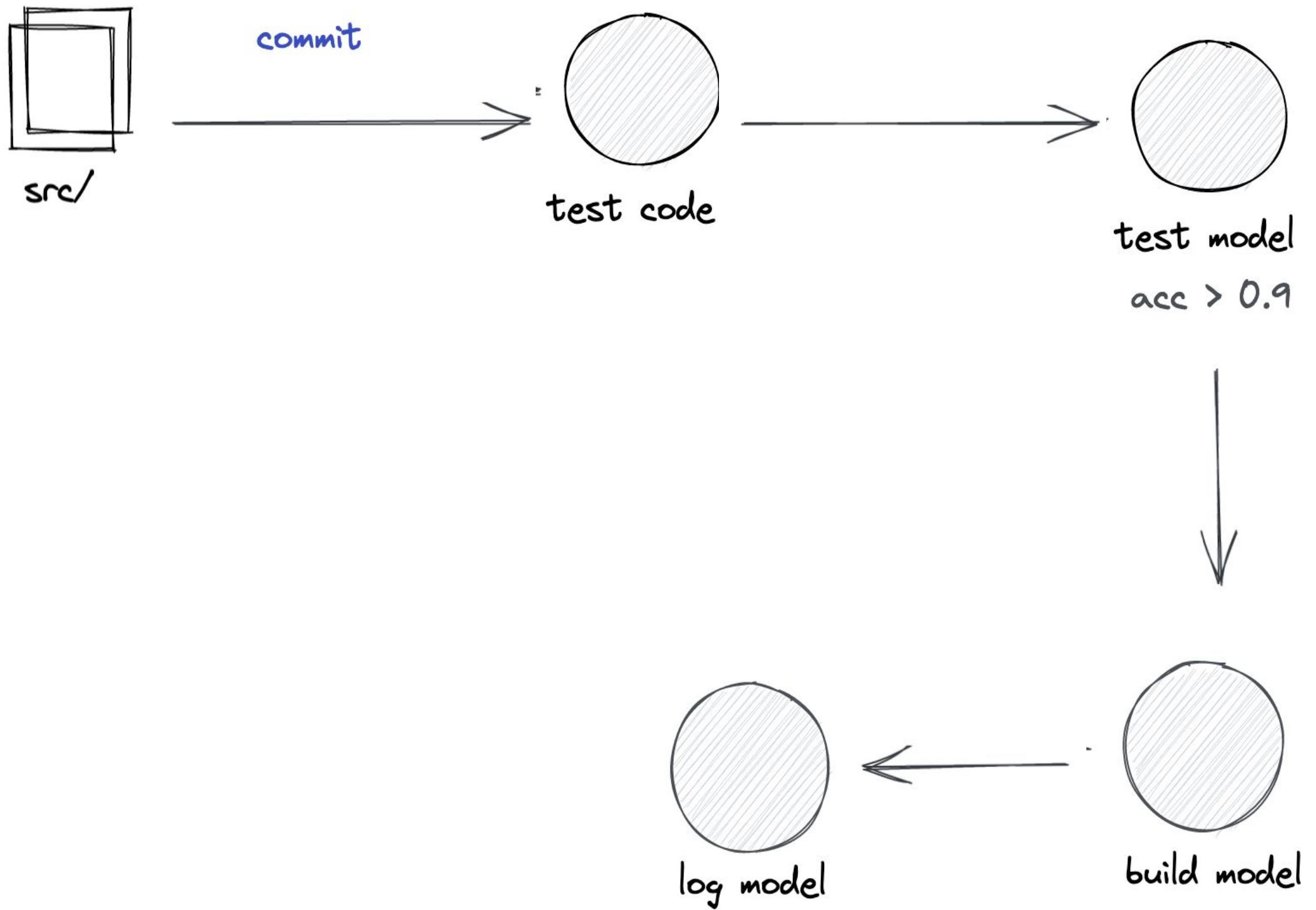


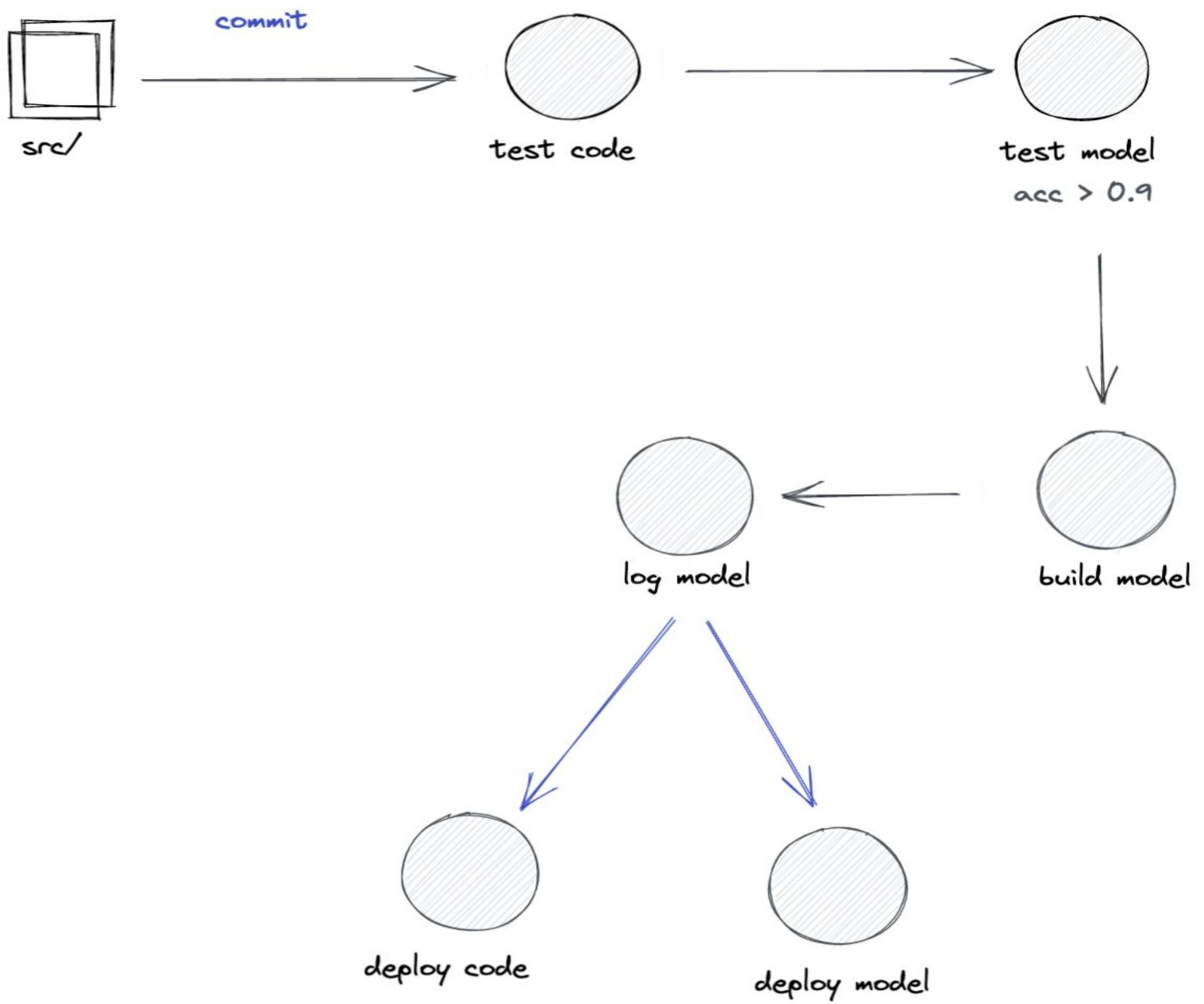
Repositório

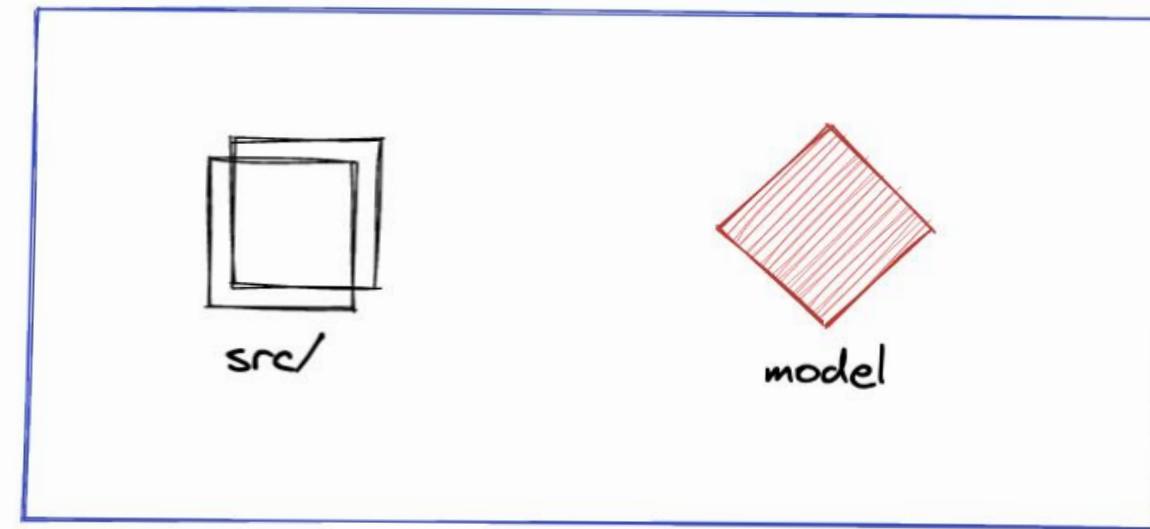




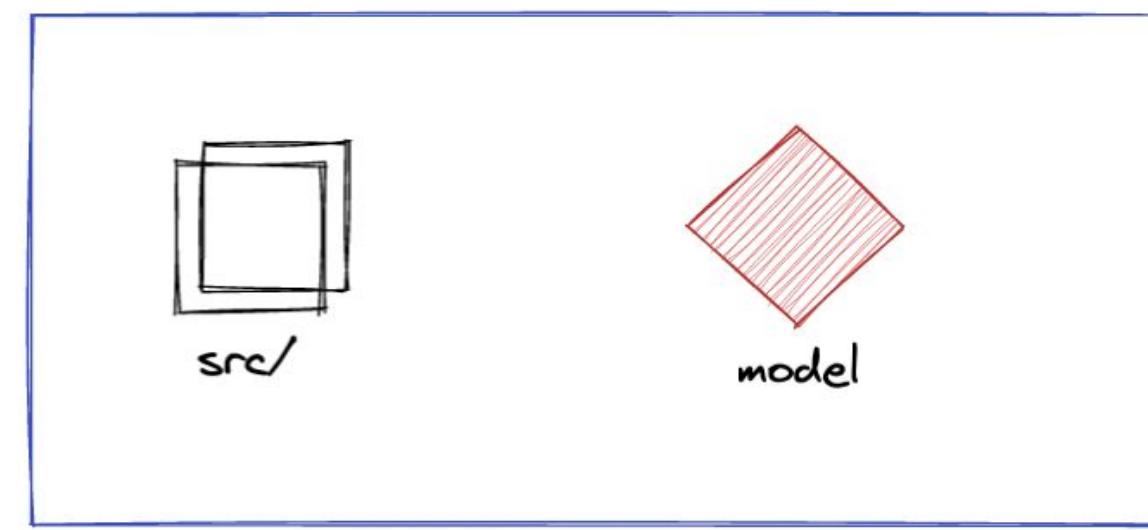








NONPROD



PROD

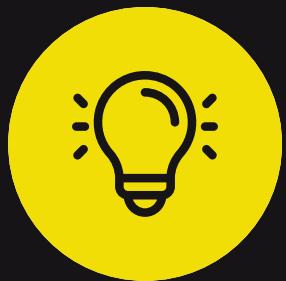
O que vimos nesta aula:



Resumo



Manutenção



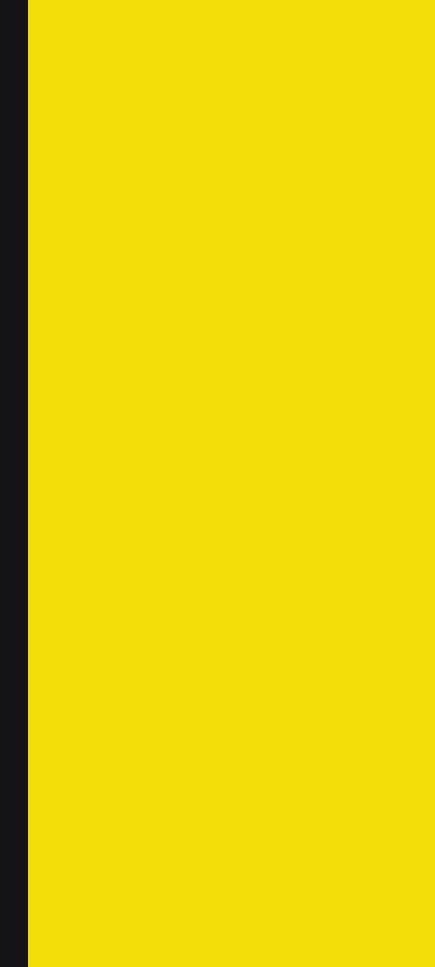
Repositório



ci/CD



Data Science & Machine Learning



07. Report Final

Consultor: Murilo Mendonça

O que veremos nesta aula:



Report final



**Documentar tudo.
Business understanding ao Deploy.**



**Quem contatar para dúvidas,
melhorias ou problemas?**



**Tornar explícito como e qual o
problema o modelo resolve.**

Checklist

- Pre-processing
- Feature selection
- Contracts
- Model selection
- Model validation
- Model optimization
- Model stage
- Drifting and retraining
- Deployment flavor
- Foreseen improvs.



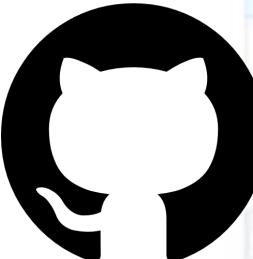
Formatos e Ferramentas

Checklist

Mark which tasks have been performed

- **Summary:** you have included a description, usage, output, accuracy and metadata of your model.
- **Pre-processing:** you have applied pre-processing to your data and this function is reproducible to new datasets.
- **Feature selection:** you have performed feature selection while modeling.
- **Modeling dataset creation:** you have well-defined and reproducible code to generate a modeling dataset that reproduces the behavior of the target dataset. This pipeline is also applicable to generate the deploy dataset.
- **Model selection:** you have chosen a suitable model according to the project specification.
- **Model validation:** you have validated your model according to the project specification.
- **Model optimization:** you have defined functions to optimize hyper-parameters and they are reproducible.
- **Peer-review:** your code and results have been verified by your colleagues and pre-approved by them.
- **Acceptance:** this model report has been accepted by the IRIS Platform team.

Formatos e Ferramentas



A screenshot of a GitHub repository page for 'basecamp / bc3-api'. The page shows basic repository statistics: 213 commits, 4 branches, 0 releases, and 17 contributors. It also displays a list of recent commits, including one from 'georgeclaghorn' and others from 'sections' and 'CONDUCT.md'. Below the stats, there's a section titled 'The Basecamp 3 API' with a welcome message and information about compatibility.

basecamp / bc3-api

Code Issues 27 Pull requests 1 Projects 0 Insights

Watch 50 Star 207 Fork 46

API documentation for Basecamp 3 <https://basecamp.com>

213 commits 4 branches 0 releases 17 contributors

Branch: master New pull request Create new file Upload files Find file Clone or download

georgeclaghorn Chatbot lines may contain rich text Latest commit 0d875bf on Jun 17

sections Chatbot lines may contain rich text 2 months ago

CONDUCT.md Update to current project maintainer 2 years ago

README.md Document /subscription.json endpoint 5 months ago

README.md

The Basecamp 3 API

Welcome to the Basecamp 3 API! If you're looking to integrate your application with Basecamp 3 or create your own application in concert with data inside of Basecamp 3, you're in the right place. We're happy to have you!

Compatibility with previous Basecamp APIs

The Basecamp 3 API is not compatible with the [Basecamp Classic API](#) or the [Basecamp 2 API](#). All integrations will start fresh with the new API. The core ingredients are the same, though: Basecamp 3 is a REST-style API that uses JSON for

Formatos e Ferramentas

bitmovin / [github_wiki_index](#)

Code Issues 0 Pull requests 0 Wiki Pulse Graphs Settings

[Edit](#) [New Page](#)

Christopher Mueller edited this page 2 minutes ago · 1 revision

Welcome to the github_wiki_index wiki!

▶ Pages 4

- folder1
 - page1
- folder2
 - page2
 - page3
- folder3

Clone this wiki locally

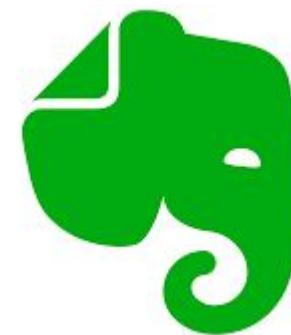
[https://github.com/bitmovin/](https://github.com/bitmovin/github_wiki_index)

[Clone in Desktop](#)

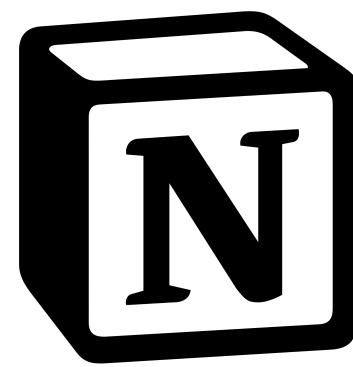
+ Add a custom footer



Formatos e Ferramentas



Wiki.js



Resumo

- Definição
- Checklist
- Formatos e Ferramentas



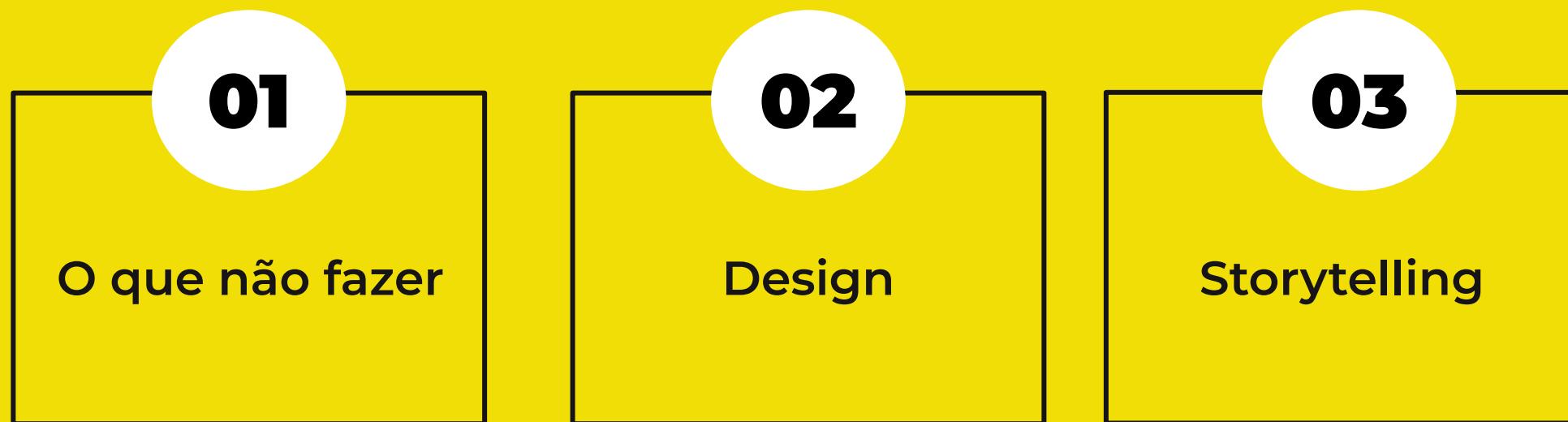


Data Science & Machine Learning

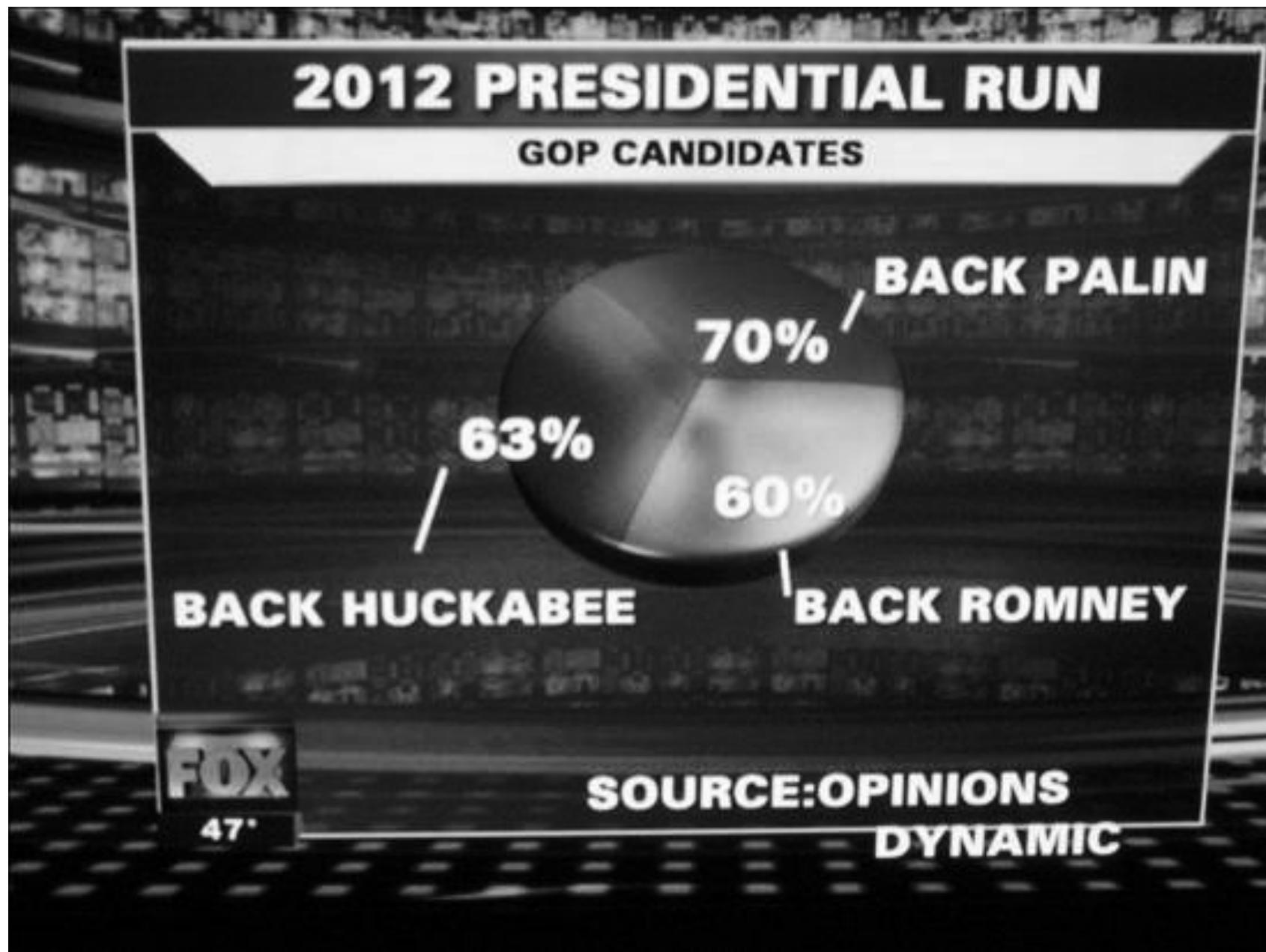
08. Apresentação Final

Consultor: Murilo Mendonça

O que veremos nesta aula:



O que não fazer



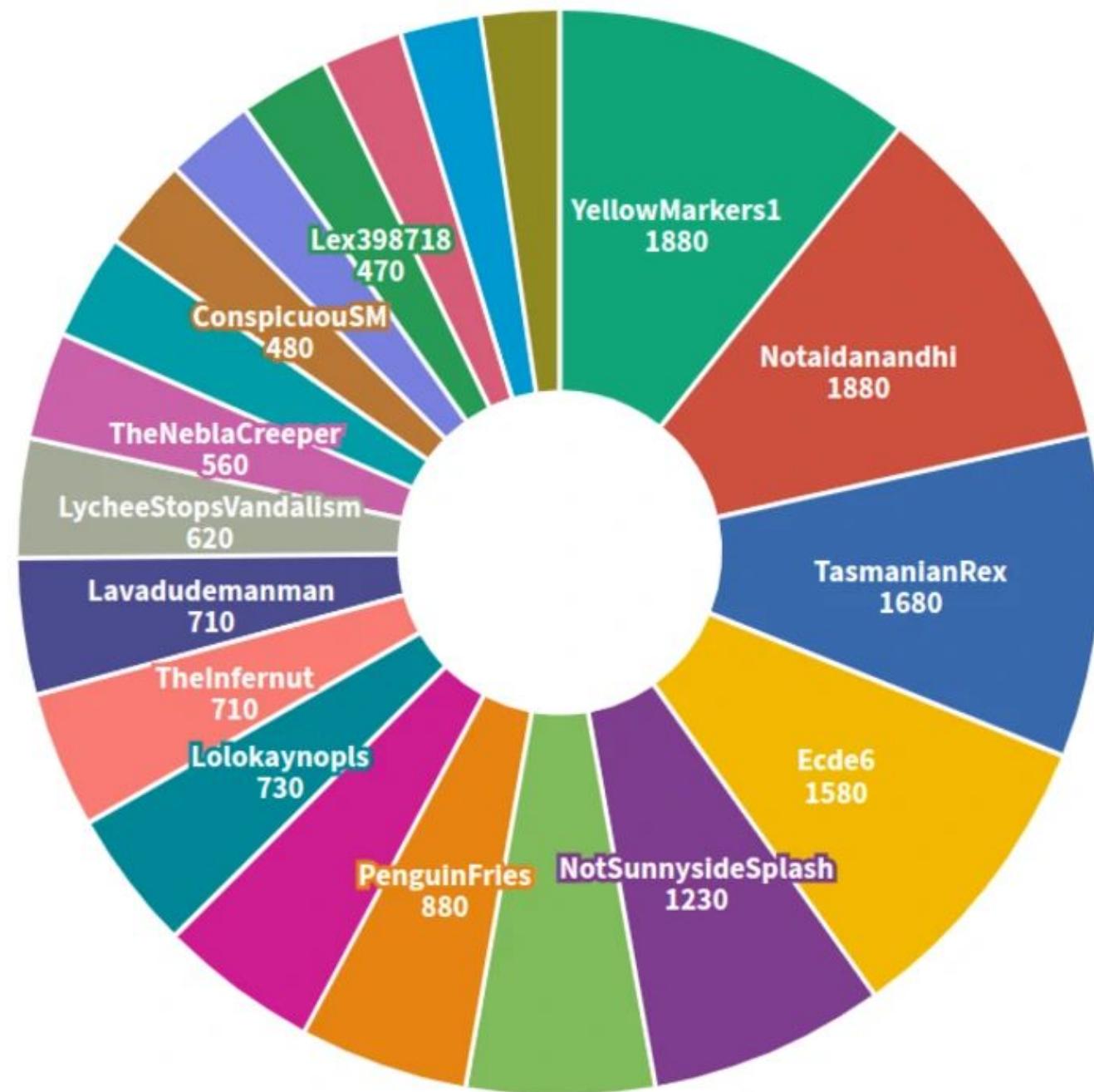
Google: "fox news data visualization"

O que não fazer



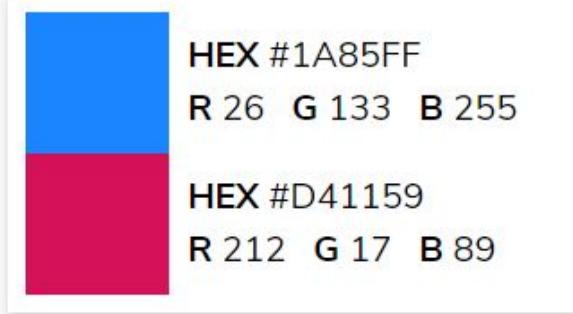
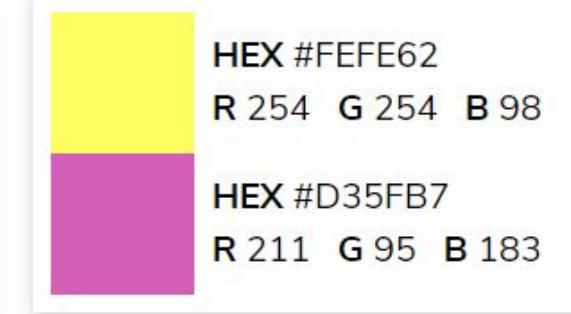
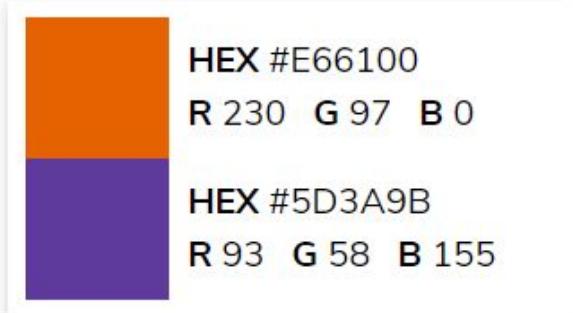
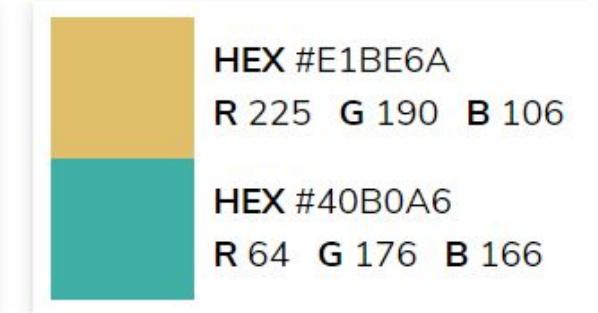
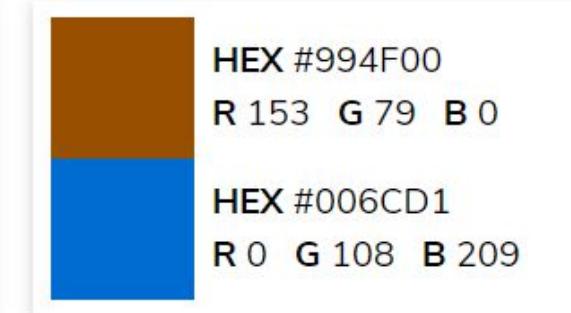
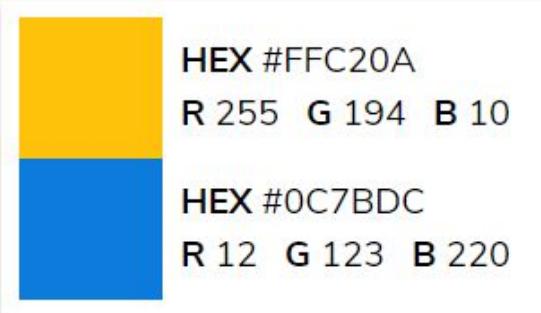
Google: "João Dória números"

O que não fazer

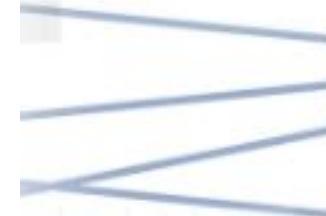
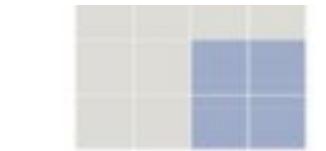


Google: "Pie Chart Hell"

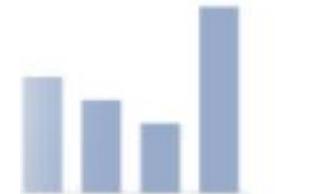
Design



Design



	A	B	C
1	10%	20%	40%
2	40%	20%	20%
3	20%	10%	30%
4	30%	20%	20%
5	50%	30%	60%
6	10%	20%	30%



	A	B	C
Category 1	10%	20%	30%
Category 2	40%	20%	20%
Category 3	20%	10%	30%
Category 4	30%	20%	20%
Category 5	50%	30%	40%
Category 6	10%	20%	30%

240/

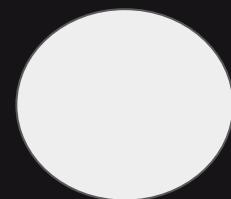
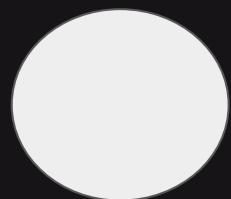
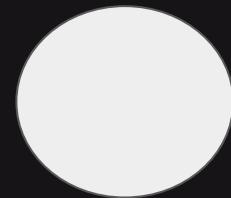
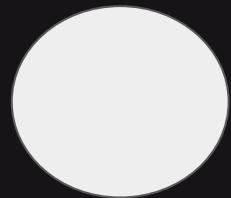
cole nussbaumer knaflic

storytelling with data

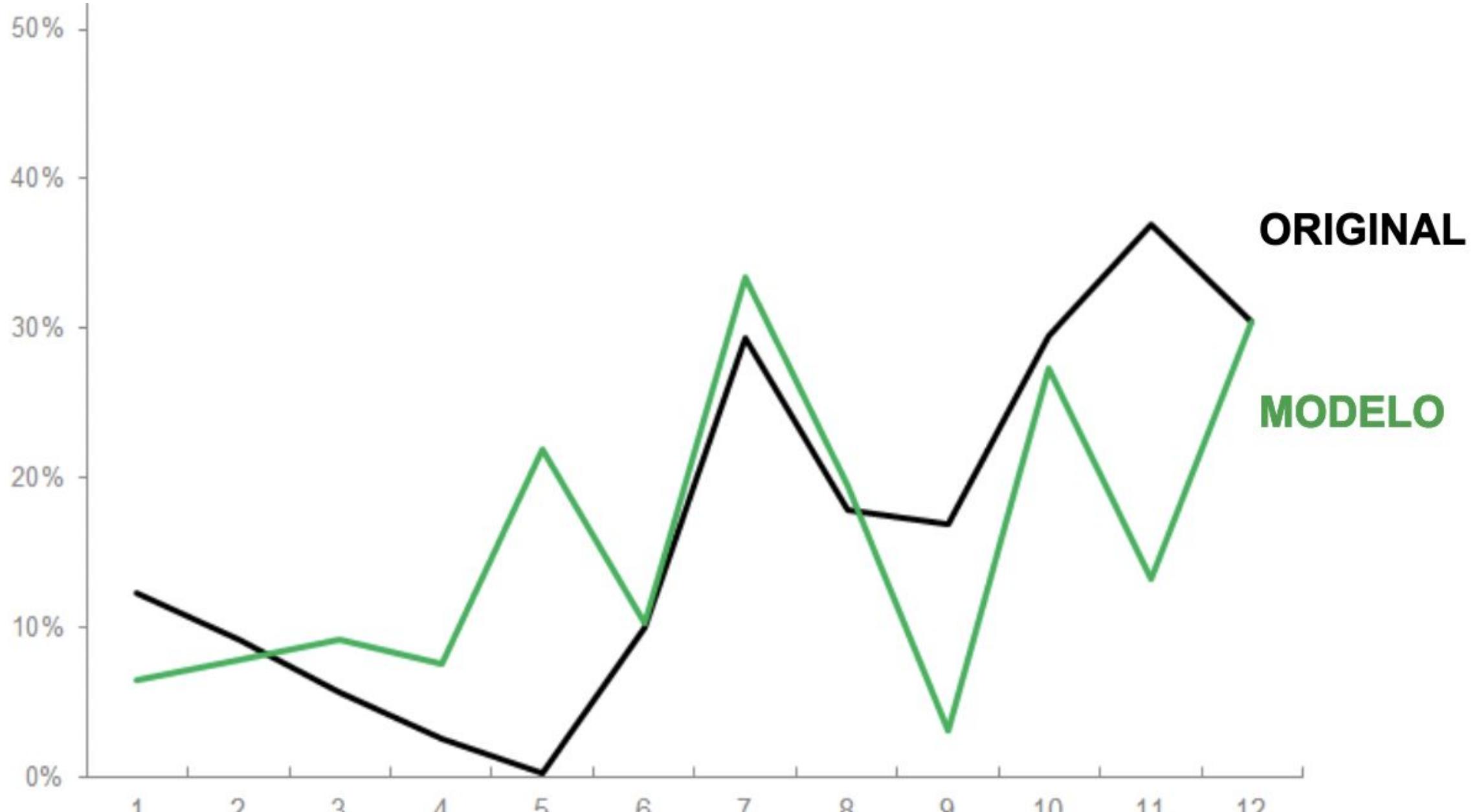
a data
visualization
guide for
business
professionals

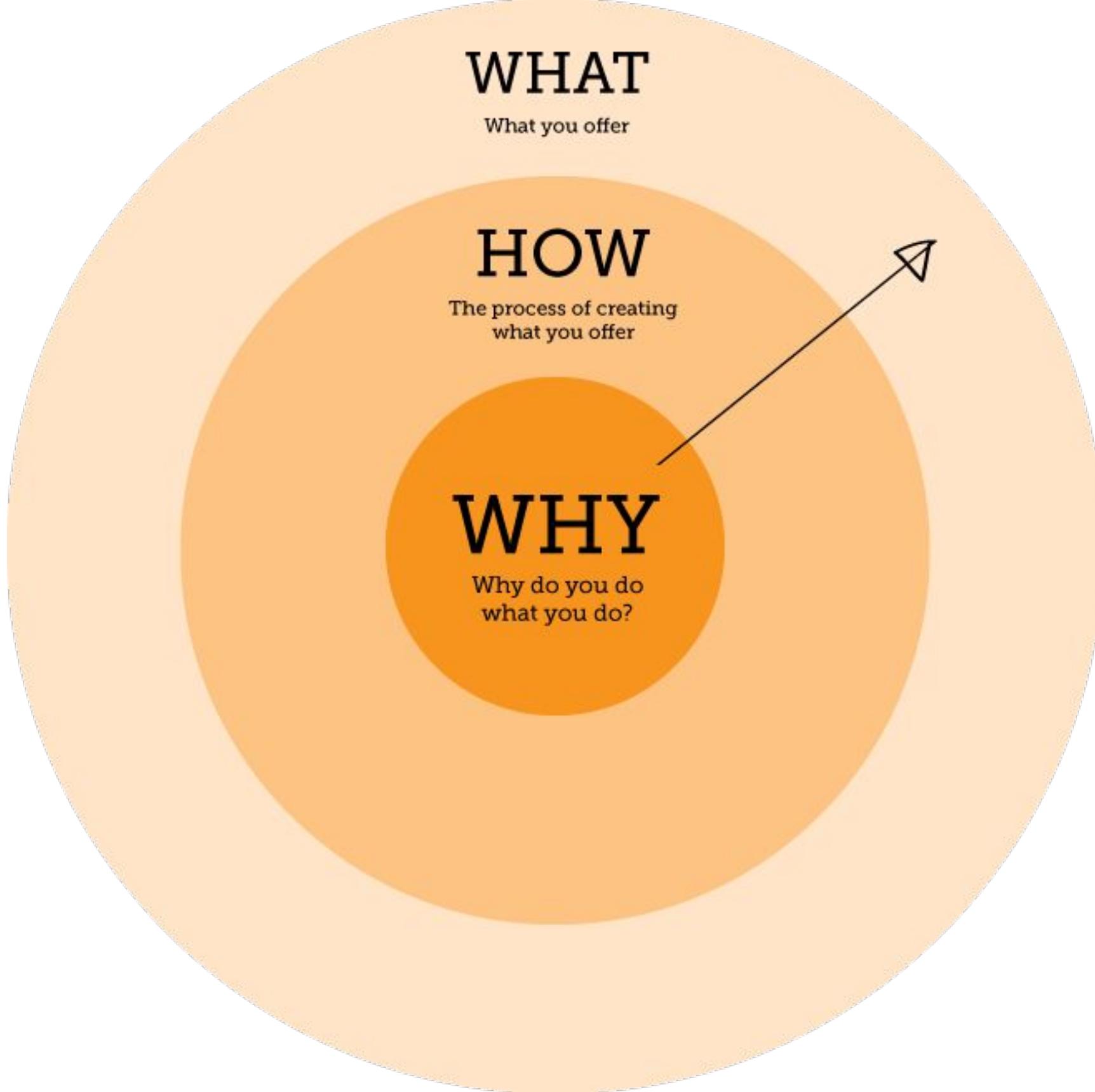
WILEY

Design



Design





WHAT

What you offer

HOW

The process of creating
what you offer

WHY

Why do you do
what you do?

Sem storytelling

Vamos criar um modelo de Regressão Logística com otimização de hiperparâmetros, usando a infraestrutura do time de plataforma para entregar uma previsão de fraude para o time de crédito.

Com storytelling e *design*

O time de crédito reportou um aumento de 50% de fraude no ano passado com a pandemia.

Por isso, vamos criar um modelo com a infraestrutura de plataforma para rapidamente começar um MVP com uma regressão logística.

Com Storytelling e *Design*

Aumento de
50%
de **Fraude**

Plataforma



Resumo

- O que Não Fazer
- Design
- Storytelling



Resumo



Design



Storytelling



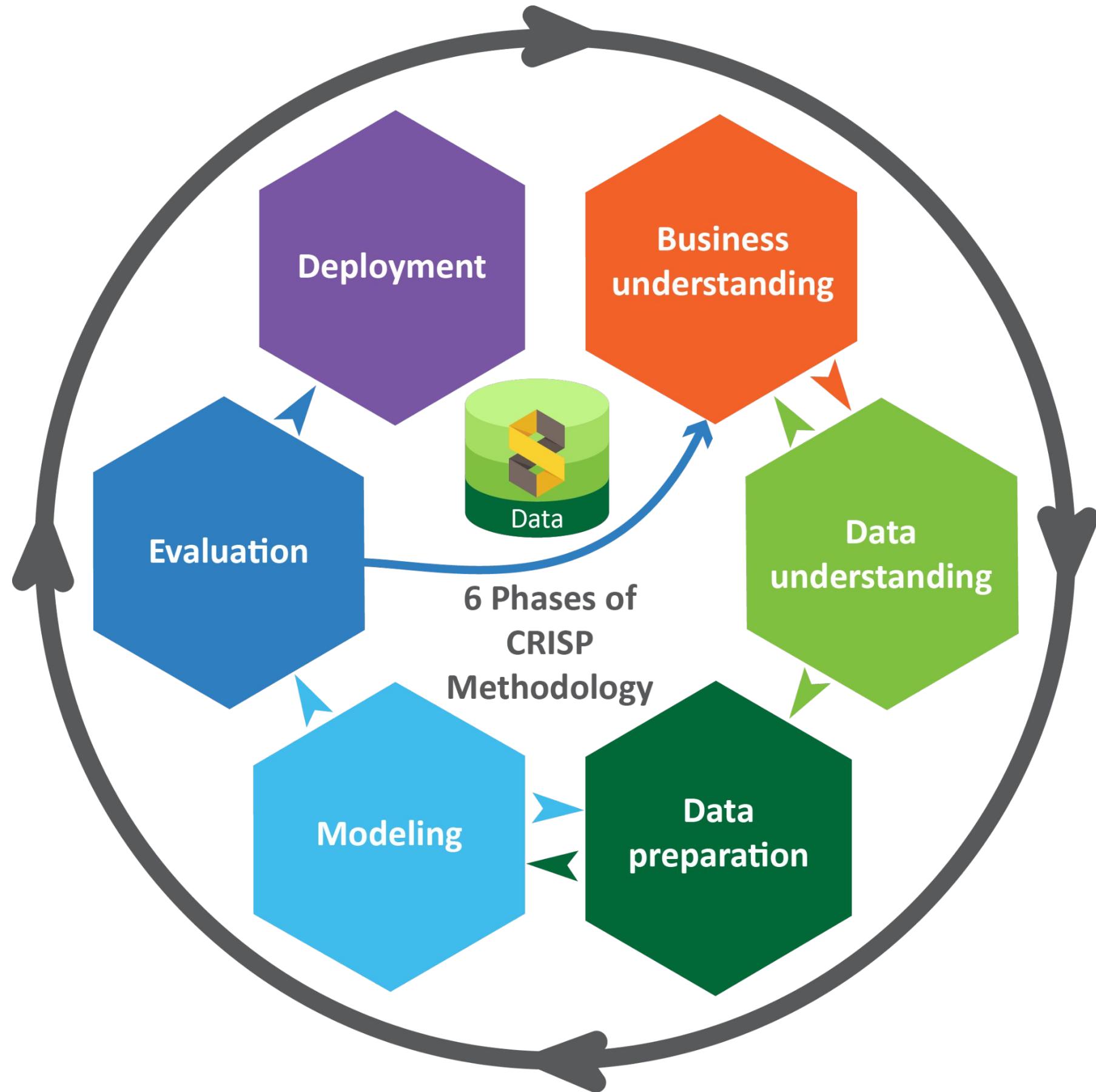
O que Não Fazer



Data Science & Machine Learning

11. Resumo do Módulo

Consultor: Murilo Mendonça



Resumo deste módulo:

