



Minimum Viable Data Products

From 0 to Data Science Pipeline

Daniel Imberman

Gonzalo Diaz

What is a MVDP?

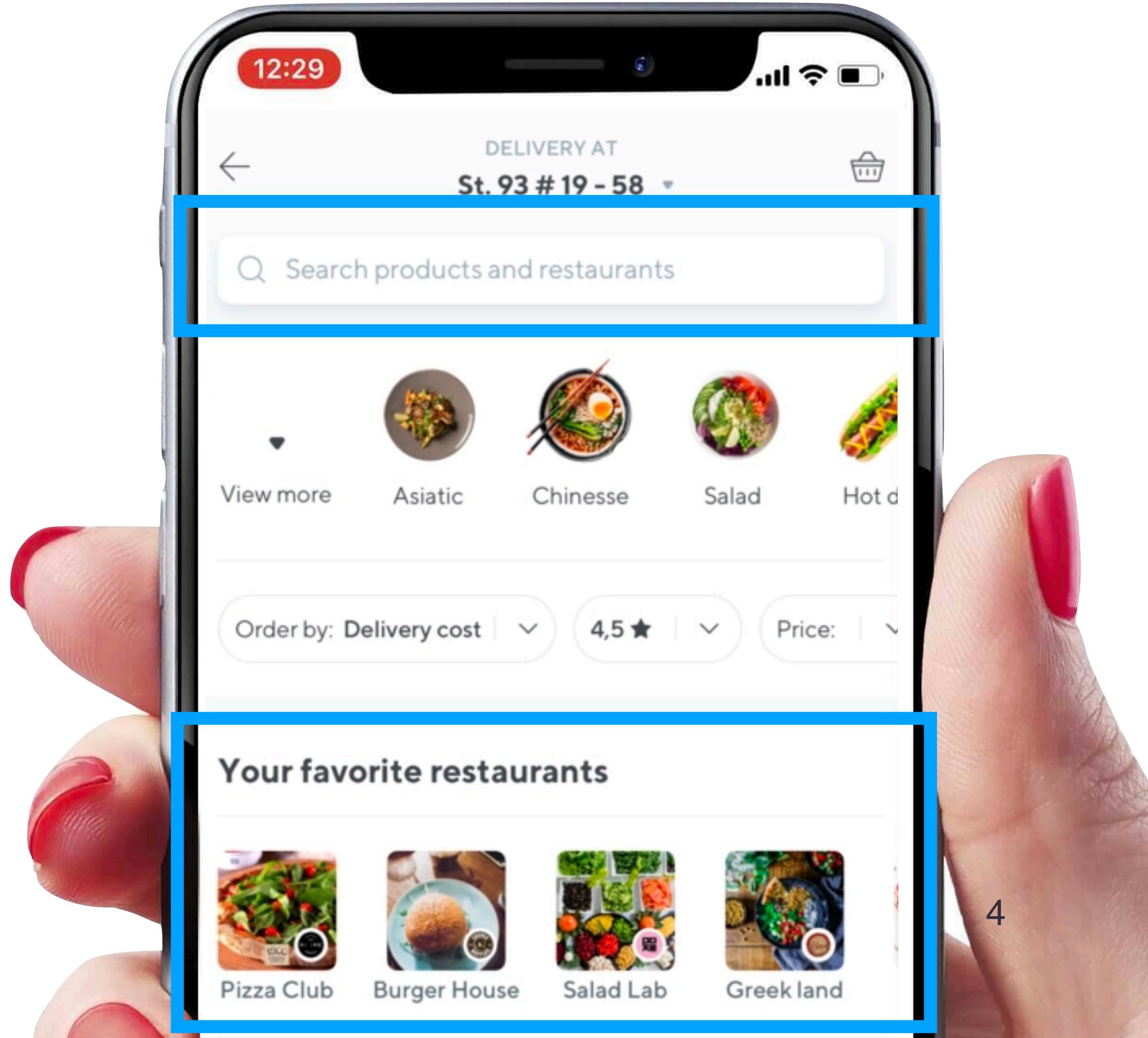
What is a Product?



- Brings value to customers
- Those customers pay companies for those products
- Those companies send engineers to Pydata Cordoba to build better products



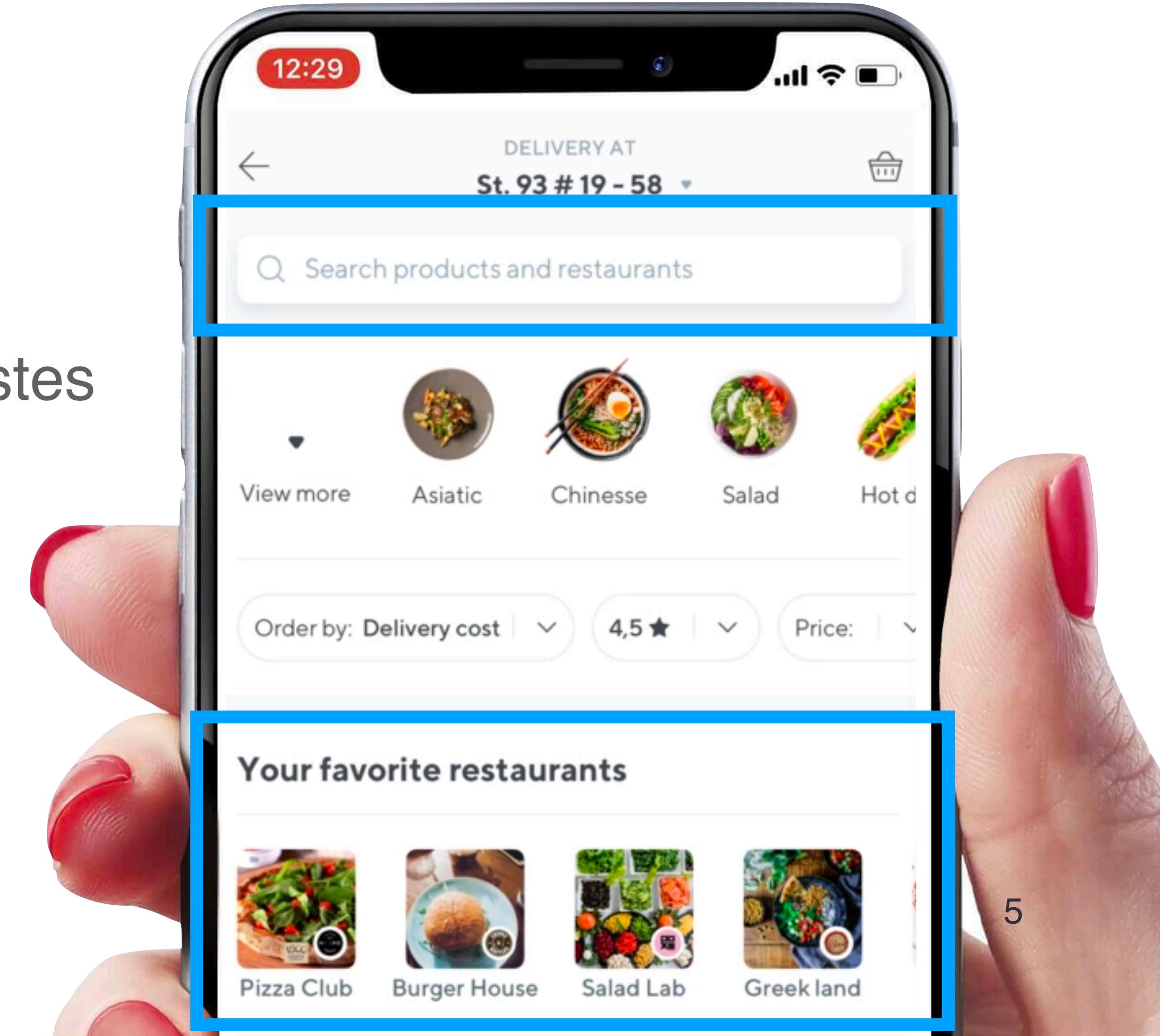
What is a Data Product?



Data Products



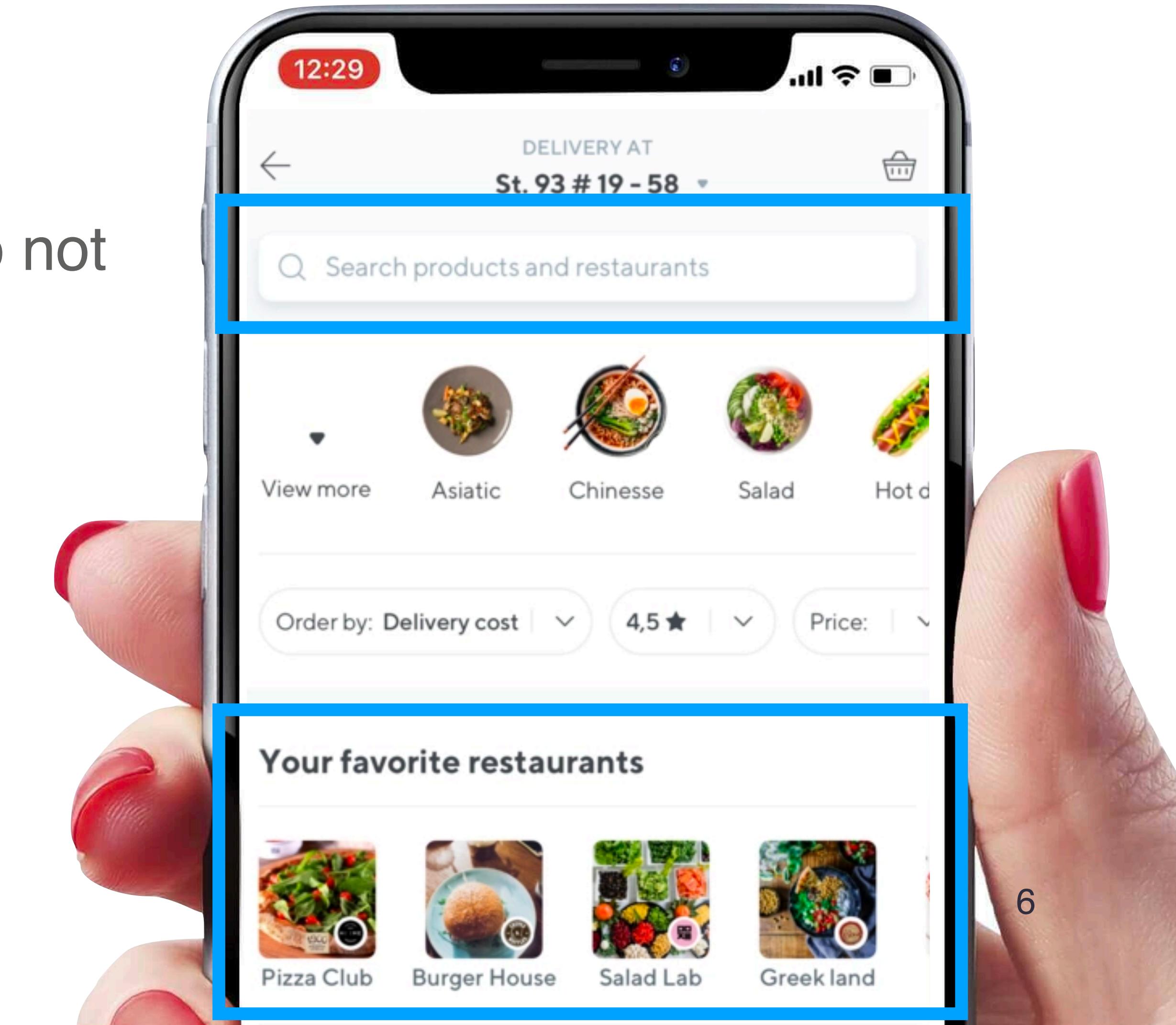
- A product feature that personalizes your experience
- A search bar that remembers your order
- A suggestion tab that can compare your tastes with thousands of similar customers
- An auto-checkout that know what you want after a crappy day



Data Products



- Good products AND good data products do not appear out of thin air
- Many start as Minimum Viable Products

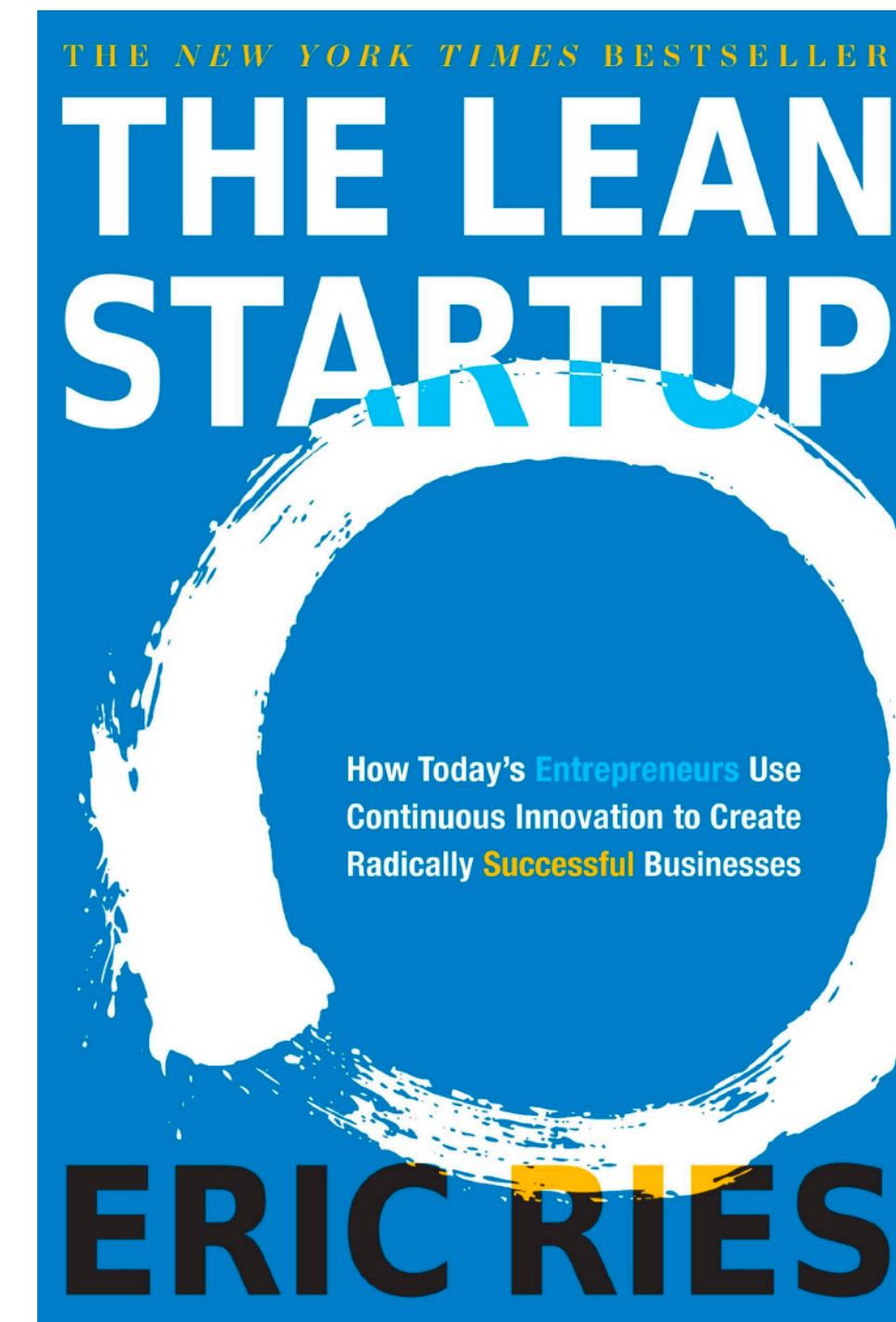


Minimum Viable Product



Eric Ries - The Lean Startup

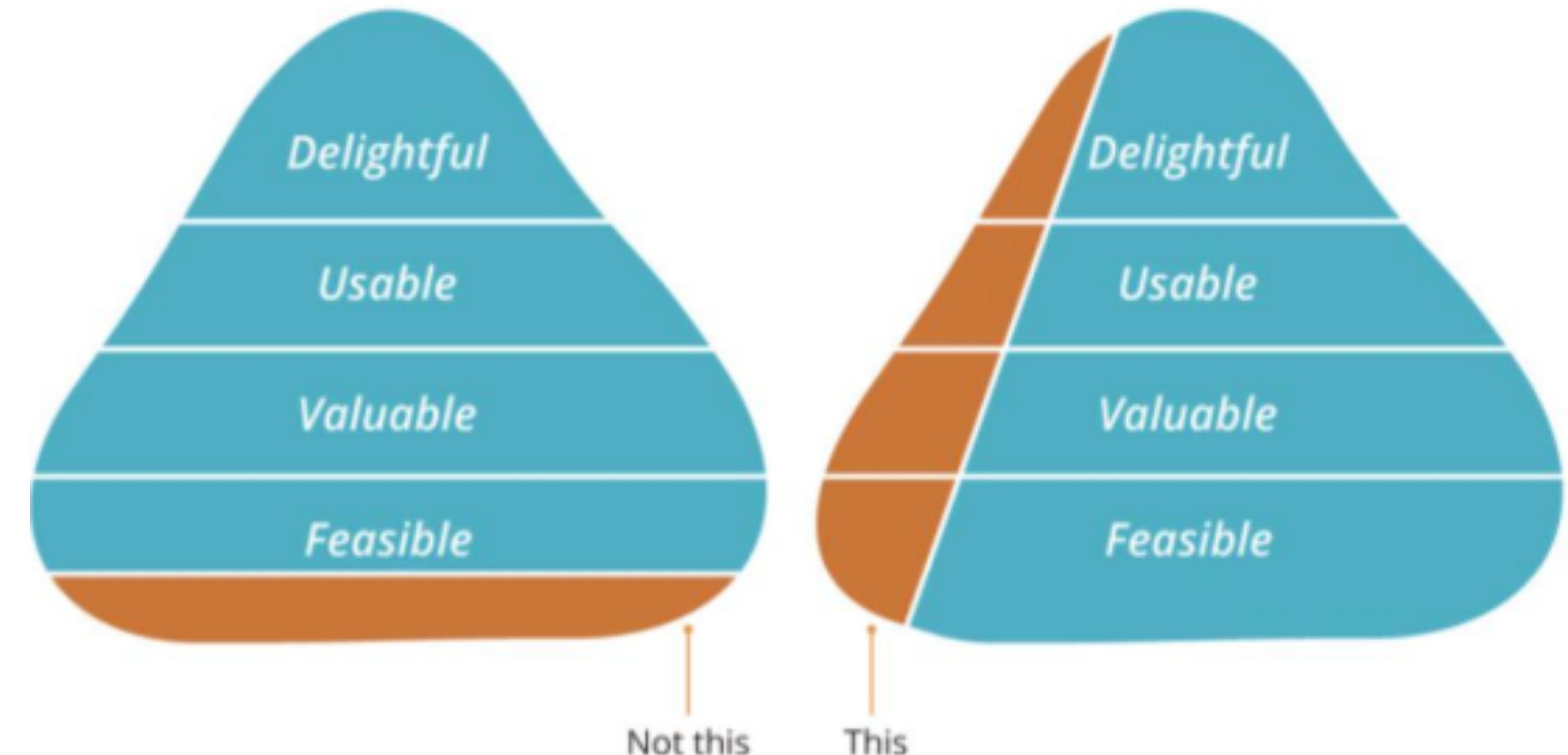
- The minimum feature set to express what the product CAN be
- An easy iteration to incorporate client feedback
- Constant validation of product/market fit at every step



Minimum Viable Data Product

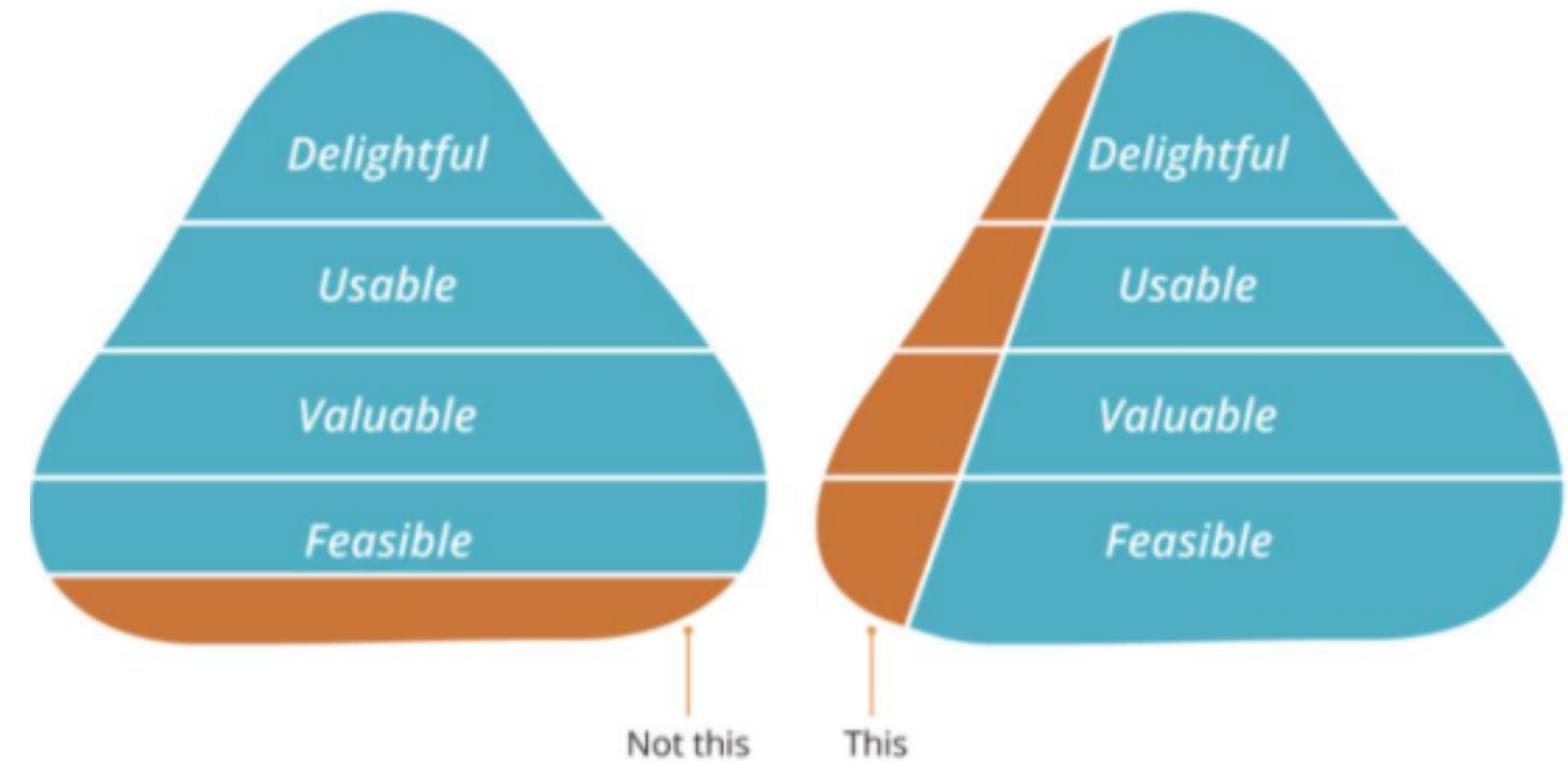


- Does this product introduce information the customer cares about?
- Is it doing so with enough accuracy to be valuable?
- Valuable for getting internal buy-in from management



credits @boagworld

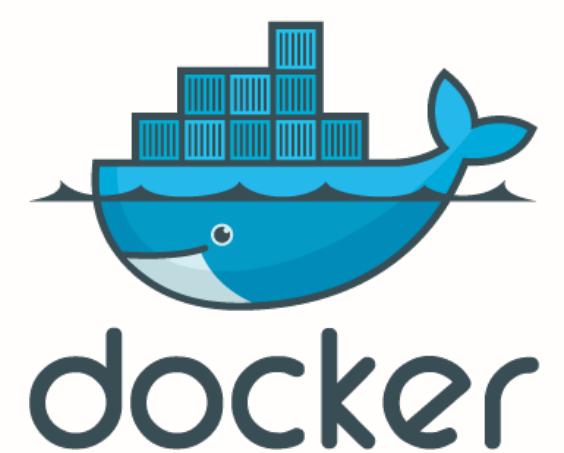
- With both MVPs and MVDPs there is a single factor that is consistent for success



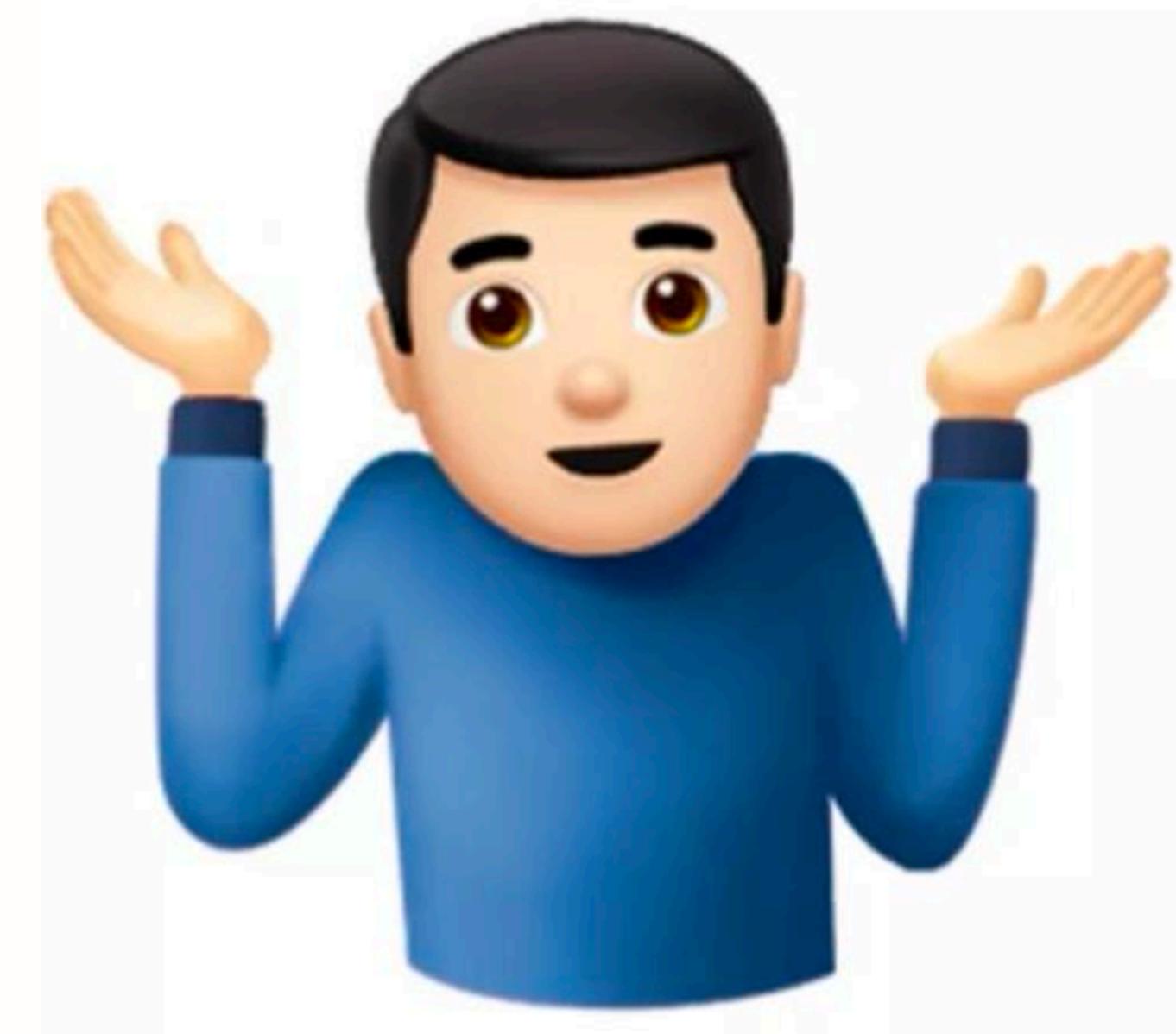
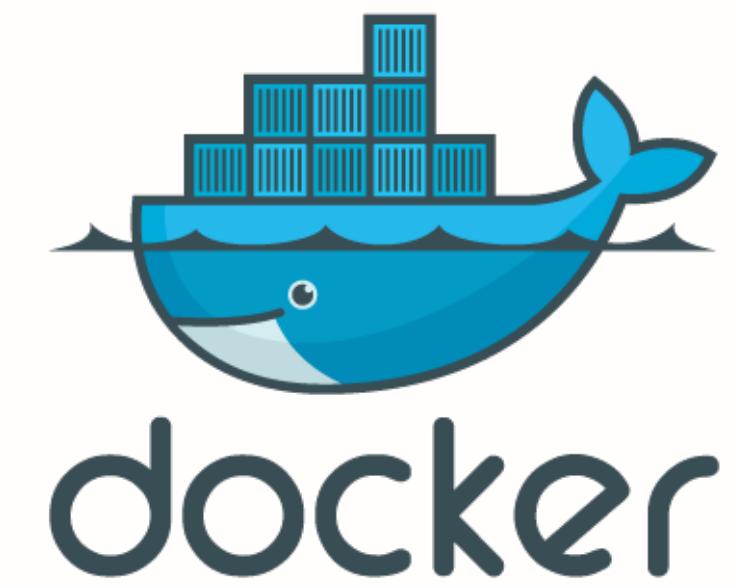
credits @boagworld

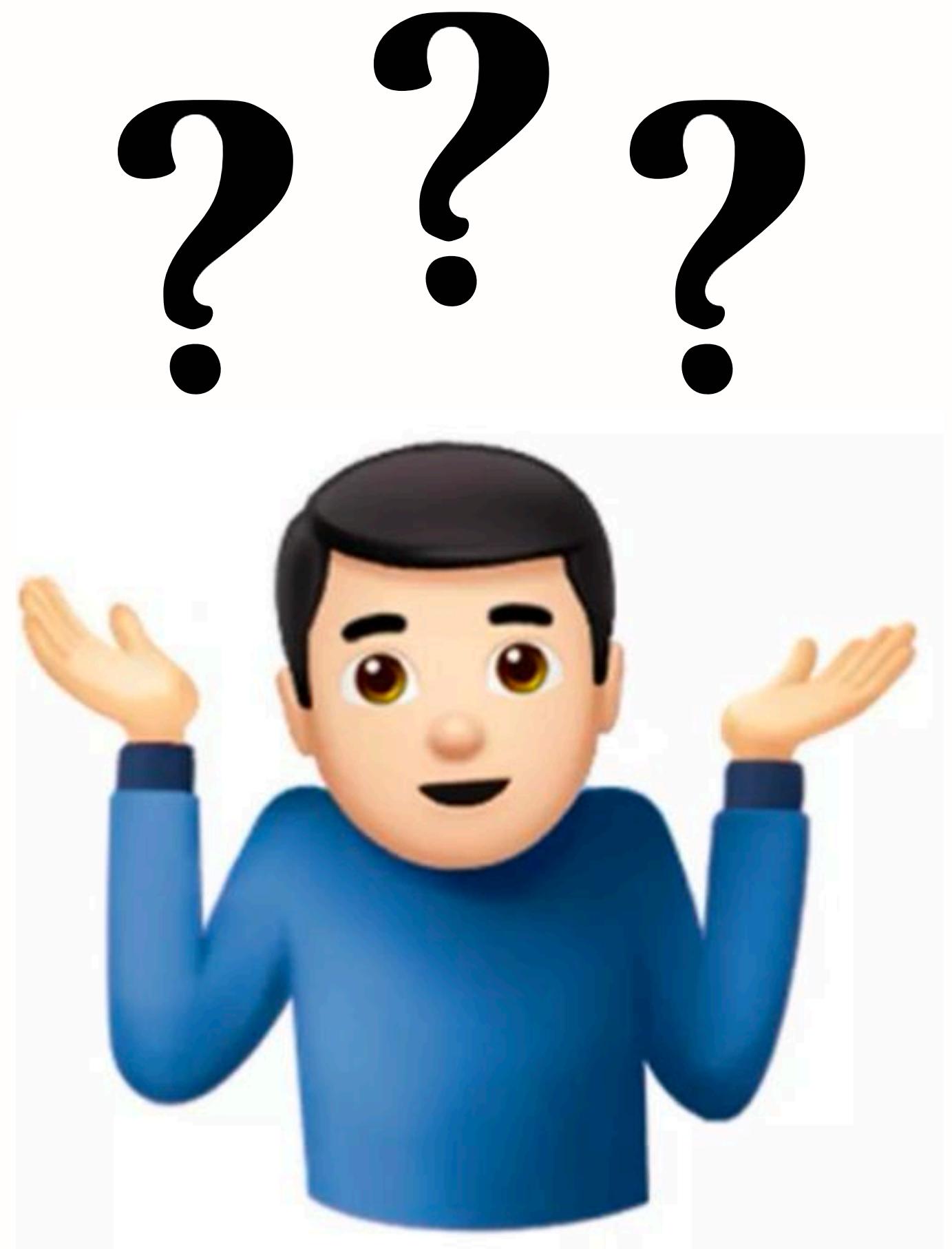
The Most Important Part of an MVP is The Ability to Iterate Quickly

From 0 to Data Science Pipeline



<https://github.com/gonzalodiaz/pydata2019cba-mvdp-airflow-jupyter>







Jupyter Notebooks



Watch 308 Star 6,282 Fork 2,631

ASTRONOMER

- Notebooks are a great way to tell stories
- Jupyter Notebooks are interactive Python/R/Julia/JS/* interpreters
- Enables iterative experimentation with data
- Pandas, Scikit-learn, Numpy
- Spark, Tensorflow

The screenshot shows a Jupyter Notebook interface. The top menu bar includes File, Edit, View, Run, Kernel, Tabs, Settings, and Help. A toolbar below the menu contains icons for file operations like new, open, save, and run. The main area displays a notebook titled "first_notebook.ipynb". The first cell contains the code `[]: print("Hola World!!!")`. The output of this cell is "This is the result of a long experimentation process" followed by three ellipses and the calculation `5 + 5 = 10`. Below this, another cell shows the expression `5 + 5`, and a third cell is currently empty.

Jupyter Docker Stacks



- Easy Notebook setup
- Many base Jupyter images with popular data science libraries pre-built
- E.g.: tensorflow, pyspark, scikitlearn
- <https://hub.docker.com/u/jupyter>

A screenshot of a web page showing three search results for "Jupyter Docker Stacks" on the Docker Hub. Each result is a card with a gray cube icon, the image name, the author ("jupyter"), the update time ("Updated 15 days ago"), and a brief description.

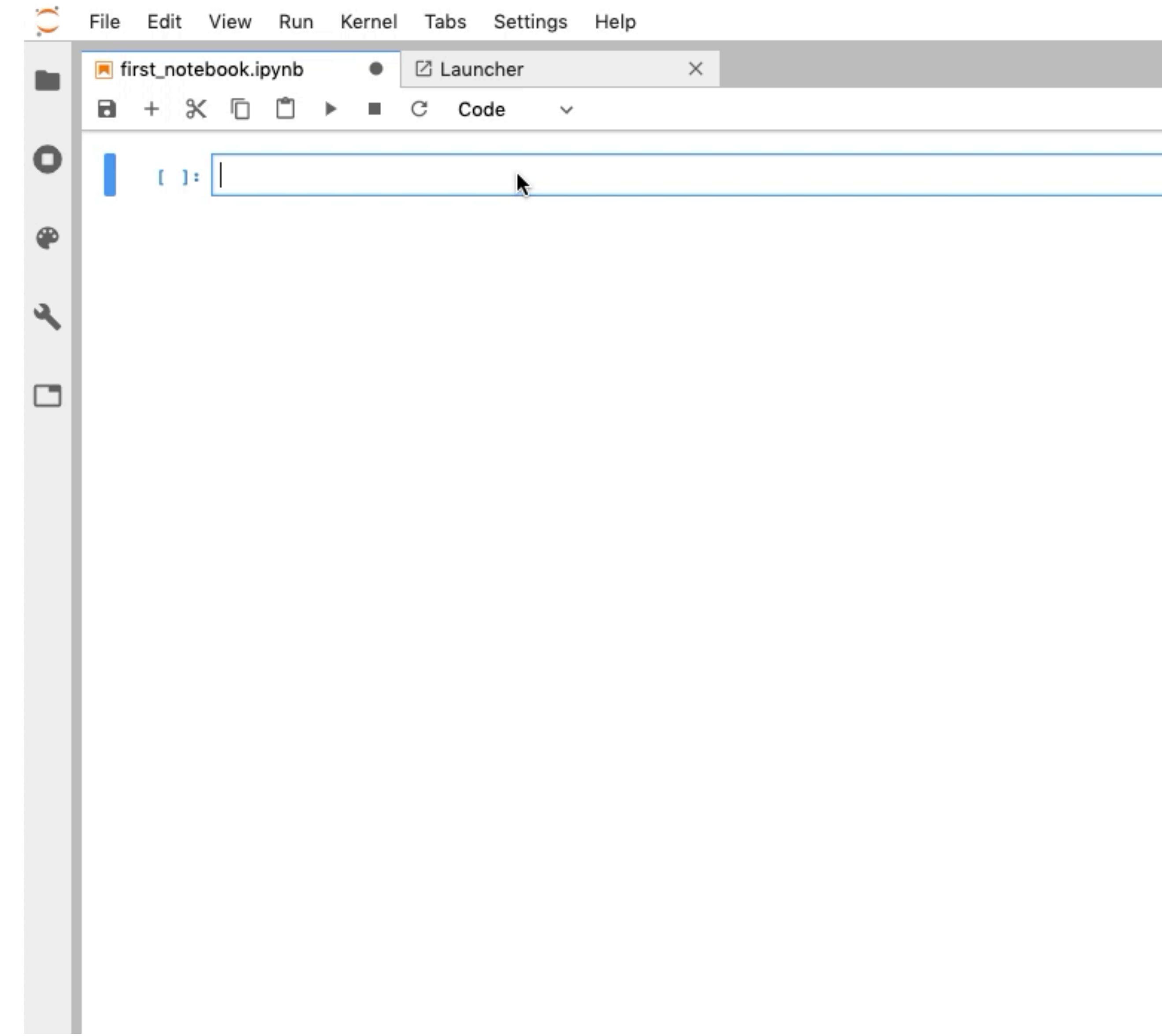
- jupyter/all-spark-notebook
By [jupyter](#) • Updated 15 days ago
Jupyter Notebook Python, Scala, R, Spark, Mesos Stack from <https://github.com/jupyter/docker-stacks/tree/master/all-spark-notebook>
Container
- jupyter/tensorflow-notebook
By [jupyter](#) • Updated 15 days ago
Jupyter Notebook Scientific Python Stack w/ Tensorflow from <https://github.com/jupyter/docker-stacks/tree/master/tensorflow-notebook>
Container
- jupyter/pyspark-notebook
By [jupyter](#) • Updated 15 days ago
Jupyter Notebook Python Stack w/ PySpark from <https://github.com/jupyter/docker-stacks/tree/master/pyspark-notebook>
Container

JupyterLab

[Watch](#) 347 [Star](#) 8,777 [Fork](#) 1,338

Rappi
ASTRONOMER

- Great integrated development environment for Jupyter Notebooks.
- More than just Notebooks
- Will replace Jupyter Notebooks soon





→ pydata2019cba-mvdp-airflow-jupyter git:(master) x └



JupyterLab

localhost:8888/lab?

File Edit View Run Kernel Tabs Settings Help

Untitled.ipynb KaggleTitanic.ipynb

Name / templates / nb_originals / Last Modified

KaggleTitanic.ipynb a minute ago

Parameters

```
[ ]: n_estimators = 50
datasets_root_path = '/home/jovyan/work/jupyter/templates/datasets'
```

Import Libraries

```
[ ]: # Ignore warnings
import warnings
warnings.filterwarnings('ignore')

# Handle table-like data and matrices
import numpy as np
import pandas as pd

# Modelling Algorithms
from sklearn.ensemble import RandomForestClassifier

# Modelling Helpers
from sklearn.model_selection import train_test_split
```

Data Loading

```
[ ]: train = pd.read_csv(f'{datasets_root_path}/titanic/train.csv')
test = pd.read_csv(f'{datasets_root_path}/titanic/test.csv')

full = train.append( test , ignore_index = True, sort = False )
titanic = full[ :891 ]

del train , test

print ('Datasets:' , 'full:' , full.shape , 'titanic:' , titanic.shape)
```

Data Preparations

```
[ ]: ## Data Preparation
sex = pd.Series( np.where( full.Sex == 'male' , 1 , 0 ) , name = 'Sex' )
embarked = pd.get_dummies( full.Embarked , prefix='Embarked' )
```

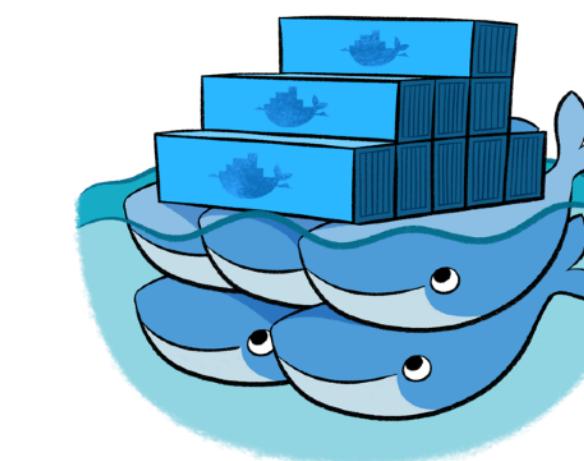
20

Distributed Jupyter Notebooks Computing

Rappi
ASTRONOMER



Enterprise Gateway



JupyterHub

Rappi

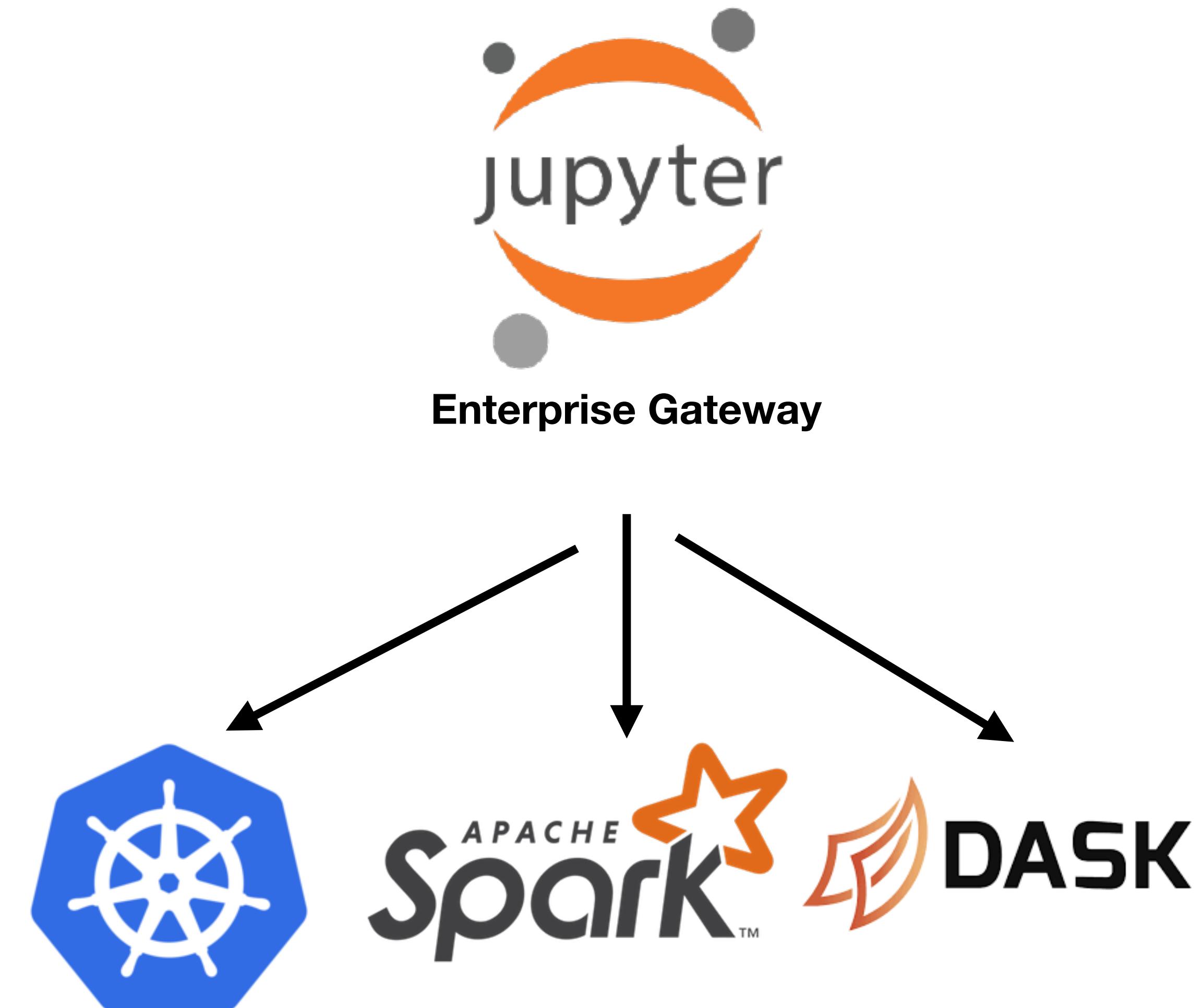
- Remote running of Jupyter Notebooks
 - Can run Multi-tenant on a single instance
 - Offers solution for Kubernetes



Jupyter Enterprise Gateway



- Runs on YARN, Kubernetes, etc.
- Easy Management of multi-tenant environments
- Integrates natively with Spark, Dask, and other parallel libraries



- Git integration for easy launching and sharing of Jupyter Notebooks
- Not for enterprise: Notebooks must be public



Papermill

[Watch](#) 84 [Star](#) 2,457 [Fork](#) 198



- Parameterize Notebooks through cell tagging
- Execute using Python API or CLI
- Store to S3, Azure



+ C

/ templates / papermill_tutorial_basic /

Name	Last Modified
input.ipynb	seconds ago
papermill_base.ipynb	seconds ago

input.ipynb papermill_base.ipynb C Code

Set cell metadata!

```
{  
    "tags": [  
        "parameters"  
    ]  
}
```

[]: #parameters
a = 'a'
b = 'b'

[]: print(f'a = {a}')

[]: print(f'b = {b}')

[]:

+ /lab? Tabs Settings Help

C Untitled.ipynb X KaggleTitanic.ipynb ●

Last Modified a minute ago

Parameters

```
[ ]: n_estimators = 50
datasets_root_path = '/home/jovyan/work/jupyter/templates/datasets'
```

Import Libraries

```
[ ]: # Ignore warnings
import warnings
warnings.filterwarnings('ignore')

# Handle table-like data and matrices
import numpy as np
import pandas as pd

# Modelling Algorithms
from sklearn.ensemble import RandomForestClassifier

# Modelling Helpers
from sklearn.model_selection import train_test_split
```

Data Loading

```
[ ]: train = pd.read_csv(f'{datasets_root_path}/titanic/train.csv')
test = pd.read_csv(f'{datasets_root_path}/titanic/test.csv')

full = train.append( test , ignore_index = True, sort = False )
titanic = full[ :891 ]
```

Scaling Jupyter Notebooks



- Have the algorithm
- Can parameterize notebook model builds using paper mill
- But how do we productionize?

Apache Airflow

[Watch](#) 663 [Star](#) 13,982 [Fork](#) 5,202



Apache Airflow



- Workflow scheduler with python DAGs
- Easy to schedule, scale, and monitor pipelines

Apache Airflow

Rappi

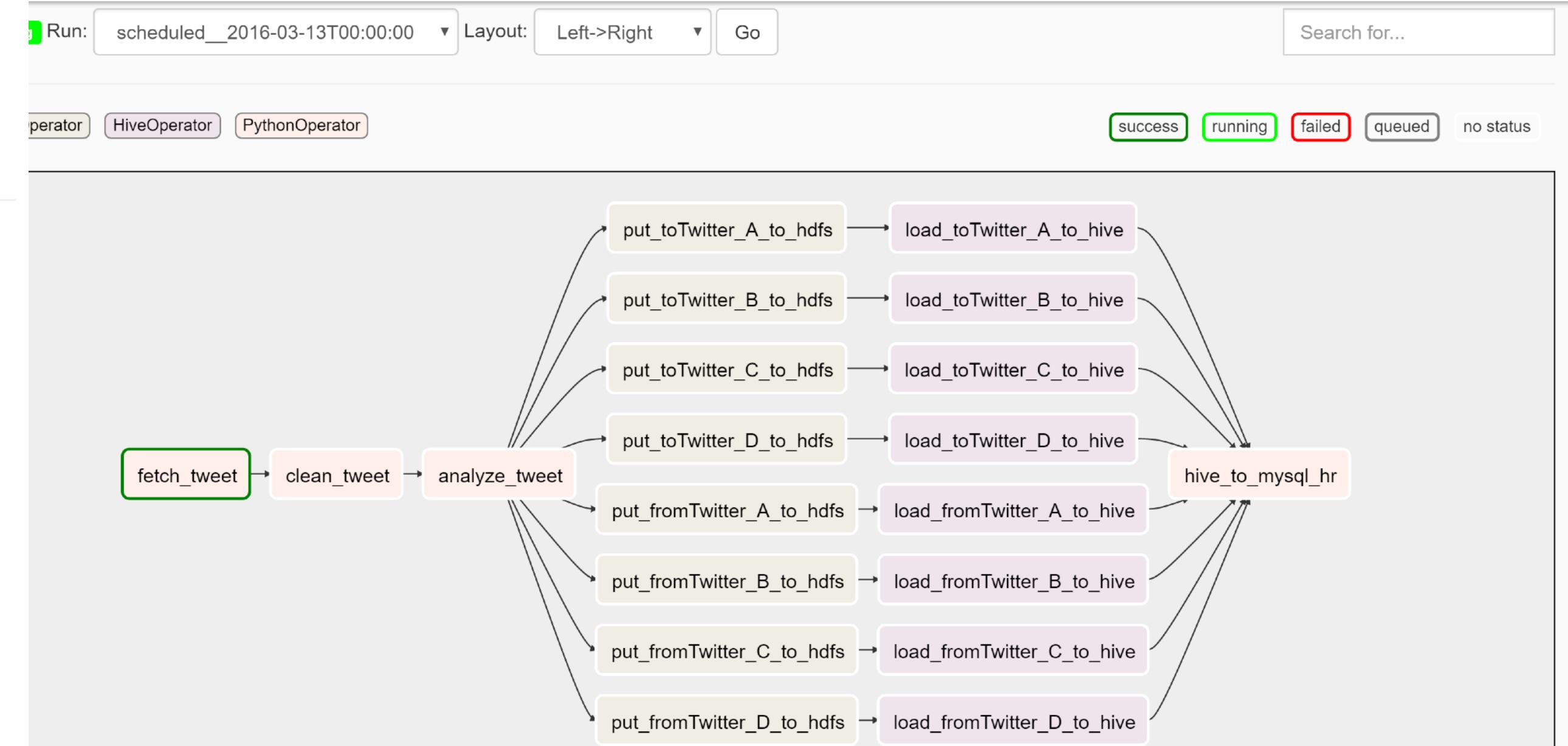
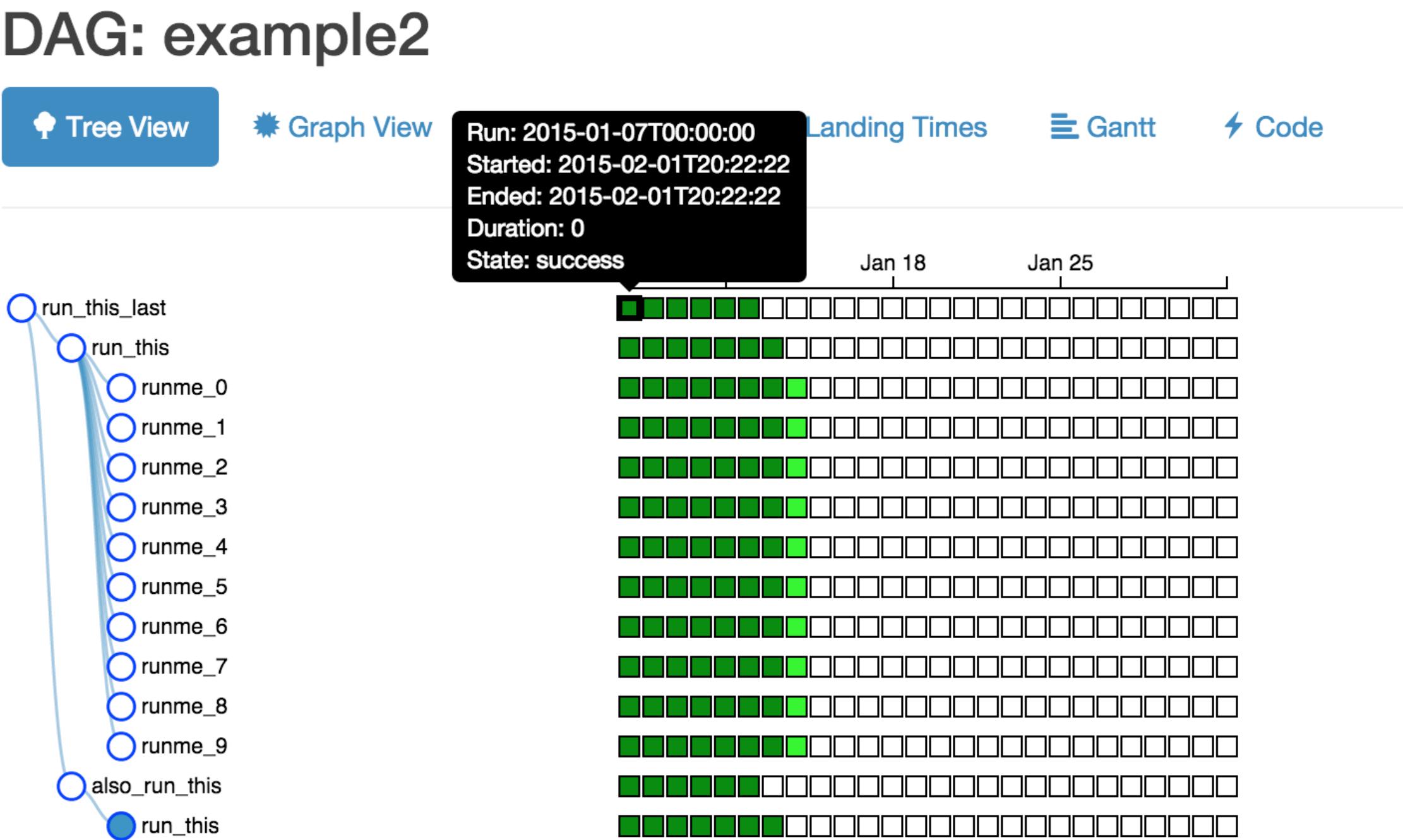
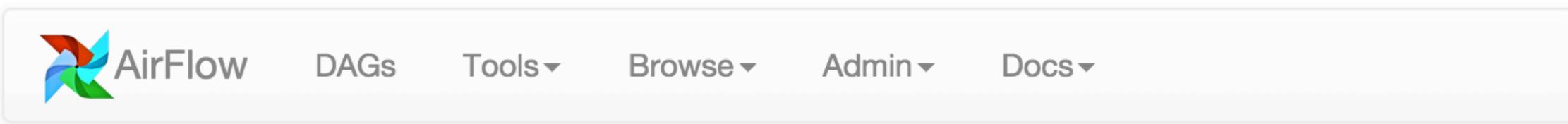
ASTRONOMER

The screenshot shows the Apache Airflow web interface with the title "DAGs". The top navigation bar includes links for "Airflow", "DAGs", "Data Profiling", "Browse", "Admin", "Docs", and "About". The timestamp "11:27 UTC" and a power icon are also present. The main content area displays a table of four DAG entries:

	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
<input checked="" type="checkbox"/>	batch_mysql_v2	01***	airflow	2 ○ ○ ○ ○ ○	2017-08-08 01:00 i	28 ○ ○ ○	View Edit Logs Metrics Task Instances DAG Runs
<input checked="" type="checkbox"/>	batch_postgresql_v1	30 0 ***	airflow	○ ○ ○ 1 ○ ○	2017-08-08 00:30 i	18 ○ ○ 1	View Edit Logs Metrics Task Instances DAG Runs
<input checked="" type="checkbox"/>	batch_sqlserver_v2	30 1 ***	airflow	2 ○ ○ ○ ○ ○	2017-08-08 01:30 i	28 ○ ○ ○	View Edit Logs Metrics Task Instances DAG Runs
<input checked="" type="checkbox"/>	sales_ftp_v1	0 1 ***	airflow	2 ○ ○ ○ ○ ○	2017-08-08 01:00 i	28 ○ ○ ○	View Edit Logs Metrics Task Instances DAG Runs

Below the table, a message says "Showing 1 to 4 of 4 entries" and a page navigation bar shows "Previous" (disabled), "1", and "Next". A link "Hide Paused DAGs" is also visible.

Apache Airflow



```
dag = DAG('tutorial', default_args=default_args)

# t1, t2 and t3 are examples of tasks created by instantiating operators
t1 = BashOperator(
    task_id='print_date',
    bash_command='date',
    dag=dag)

t2 = BashOperator(
    task_id='sleep',
    bash_command='sleep 5',
    retries=3,
    dag=dag)

templated_command = """
    {% for i in range(5) %}
        echo "{{ ds }}"
        echo "{{ macros.ds_add(ds, 7) }}"
        echo "{{ params.my_param }}"
    {% endfor %}
"""

t3 = BashOperator(
    task_id='templated',
    bash_command=templated_command,
    params={'my_param': 'Parameter I passed in'},
    dag=dag)

t2.set_upstream(t1)
t3.set_upstream(t1)
```

Define Operators

Set Dependencies

Data Science in Airflow



- Hyper-parameter Tuning
- Canary Testing
- A/B Testing
- Model Deployment

Hyper-parameter tuning



- Choosing a set of optimal hyperparameters for a learning algorithm.
- E.g.: Loss-rate, tolerance, etc.

Grid-search



- Essentially a for-loop
- Pythonic DAG creation shines here
- Can nest for-loops for multiple degrees



DAGs

Data Profiling ▾

Browse ▾

Admin ▾

Docs ▾

About ▾

2019-09-22 19:21:29 UTC

DAG [example_papermill_operator_hyperparam] is now fresh as a daisy × Off **DAG: example_papermill_operator_hyperparam**

schedule: None

 Graph View  Tree View  Task Duration  Task Tries  Landing Times  Gantt  Details  Code  Trigger DAG  Refresh  DeleteNone Base date: 2019-09-22 19:21:08 Number of runs: 25 Run: Layout: Left->Right Go

Search for...

PapermillOperator

success running failed skipped up_for_reschedule up_for_retry queued no_status

DAG [example_papermill_operator_hyperparam] is now fresh as a daisy

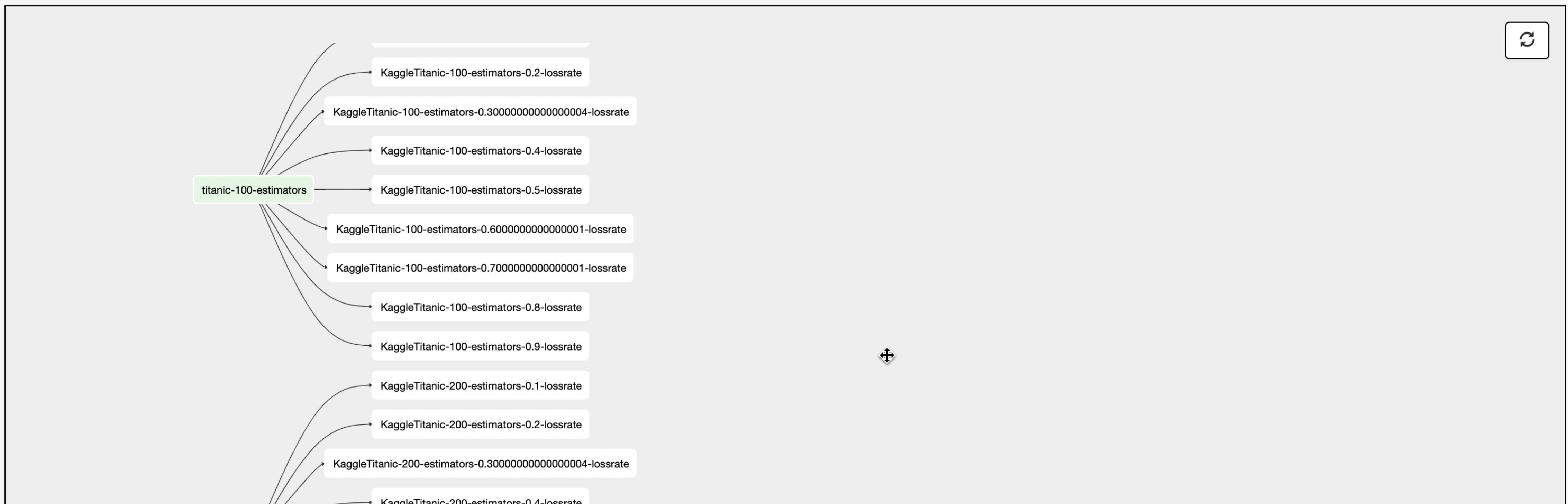
 Off **DAG: example_papermill_operator_hyperparam**

schedule: None

[Graph View](#) [Tree View](#) [Task Duration](#) [Task Tries](#) [Landing Times](#) [Gantt](#) [Details](#) [Code](#) [Trigger DAG](#) [Refresh](#) [Delete](#)

None Base date: 2019-09-22 19:22:29 Number of runs: 25 Run: Layout: Left->Right Go

Search for...

[DummyOperator](#) [PapermillOperator](#)[success](#) [running](#) [failed](#) [skipped](#) [up_for_reschedule](#) [up_for_retry](#) [queued](#) [no_status](#)

Grid-search

- Essentially a for-loop
- Pythonic DAG creation shines here
- Can nest for-loops for multiple degrees
- However not super efficient...

```
num_estimators = [100,200,500]
loss_rates = [.01, .04, .07]
with DAG(
    dag_id='example_papermill_operator', default_args=args,
    schedule_interval=None
) as dag:
    start = DummyOperator()
    dag >> start
    last = start
    for num_estimator in num_estimators:
        merge = DummyOperator()
        merge << last
        for loss_rate in loss_rates:
            task_name = 'KaggleTitanic-{}-{}-estimators'.format(num_estimator, loss_rate)
            output_name = "/nb_outputs/{}/ipynb".format(task_name)
            titanic_estimator = PapermillOperator(
                task_id=task_name,
                input_nb="/nb_originals/KaggleTitanic.ipynb",
                output_nb=output_name,
                parameters={"n_estimators": num_estimator, "loss_rate": loss_rate}
            )
            merge << titanic_estimator
```

Optimized Grid-search



- Set a “maximum” fit
- Maintain a “source of truth”

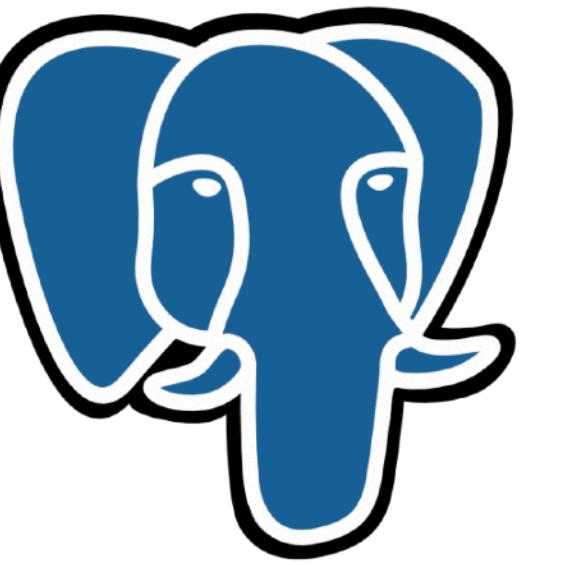
```
[15]: if get_current_max() > max_fit  
      model.fit( train_X , train_y )
```

Optimized Grid-search

Rappi
ASTRONOMER



Amazon S3



PostgreSQL



redis

Stochastic Grid Search

- Use ShortCircuitOperators to complete loop early if optimal answer found
- Use XCom to communicate the results of the previous round to the next round

```
num_estimators = [100,200,500]
loss_rates = [.01, .04, .07]
with DAG(
    dag_id='example_papermill_operator', default_args=args,
    schedule_interval=None
) as dag:
    start = DummyOperator()
    dag >> start
    last = start
    for num_estimator in num_estimators:
        merge = DummyOperator()
        merge << last
        for loss_rate in loss_rates:
            task_name = 'KaggleTitanic-{}-{}-estimators'.format(num_estimator, loss_rate)
            output_name = "/nb_outputs/{}/ipynb".format(task_name)
            titanic_estimator = PapermillOperator(
                task_id=task_name,
                input_nb="/nb_originals/KaggleTitanic.ipynb",
                output_nb=output_name,
                parameters={"n_estimators": num_estimator, "loss_rate": loss_rate}
            )
            merge << titanic_estimator
```

Canary Testing

- Method for ensuring that your new models are an improvement compared to old ones
- Can also be used to detect model errors (e.g. model isn't a huge deviation from previous models)

```
def compare_results(**context):
    new_model_result = context['task_instance'].xcom_pull(task_ids='new_model')
    old_results = []
    for model_name in models_to_compare:
        old_results.append(context['task_instance'].xcom_pull(task_ids='old-model-{}'.format(model_name)))
    if new_model_result > max(old_results):
        return "deploy_new_model"
    else:
        return "alert_of_failed_model"

with DAG(dag_id='example_canary_test', default_args=args) as dag:
    num_tasks = 5
    compare_all_results = BranchPythonOperator(
        task_id='compare_results',
        python_callable=compare_results,
        dag=dag,
        provide_context=True
    )

    deploy_new_model = DummyOperator(task_id="deploy_new_model")
    alert_of_failed_model = DummyOperator(task_id="alert_of_failed_model")
    compare_all_results >> deploy_new_model
    compare_all_results >> alert_of_failed_model

    new_model = PythonOperator(
        task_id="new_model",
        python_callable=lambda: 5
    )
    new_model >> compare_all_results
    i = 5

    for model in models_to_compare:
        attempt = PythonOperator(
            task_id="old-model-{}".format(model),
            python_callable=lambda: i
        )
        i += 1
        attempt >> compare_all_results
```

Project

1: Project

2: Favorites

3: Structure

4: pydata2019cba-mvdp-airflow-jupyter ~/code/pr

5: airflow_pydata

6: dags

7: operators

8: papermill_operator.py

9: example_papermill_dag.py

10: example_papermill_hyperparam_dag.py

11: example_shortcircuit_dag.py

12: pgdata

13: docker-compose.yml

14: Dockerfile

15: requirements.txt

16: astronomer_pydata

17: .astro

18: dags

19: operators

20: example_canary_test.py

21: example_papermill_dag.py

22: example_papermill_hyperparam_dag.py

23: example_shortcircuit_dag.py

24: datasets

25: output

26: titanic

27: test.csv

28: train.csv

29: include

30: nb_originals

31: KaggleTitanic.ipynb

32: plugins

33: .dockerignore

34: .env

35: .gitignore

36: airflow_settings.yaml

37: Dockerfile

38: packages.txt

39: requirements.txt

40: jupyter

41: templates

42: .ipynb_checkpoints

43: datasets

44: titanic

45: 31

46: 32

47: 33

48: 34

49: 35

50: 36

51: 37

52: 38

53: 39

54: 40

55: 41

56: 42

57: 43

58: 44

59: 45

60: 46

61: 47

62: 48

63: 49

64: 50

65: 51

66: 52

67: 53

68: 54

69: 55

70: 56

71: 57

72: 58

73: 59

74: 60

75: 61

76: 62

77: 63

```
models_to_compare = ["model-1", "model-bar", "model-foo"]

def compare_results(**context):
    new_model_result = context['task_instance'].xcom_pull(task_ids='new_model')
    old_results = []
    for model_name in models_to_compare:
        old_results.append(context['task_instance'].xcom_pull(task_ids='old-model-{}'.format(model_name)))
    if new_model_result > max(old_results):
        return "deploy_new_model"
    else:
        return "alert_of_failed_model"

with DAG(dag_id='example_canary_test', default_args=args) as dag:
    num_tasks = 5
    compare_all_results = BranchPythonOperator(
        task_id='compare_results',
        python_callable=compare_results,
        dag=dag,
        provide_context=True
    )

    deploy_new_model = DummyOperator(task_id="deploy_new_model")
    alert_of_failed_model = DummyOperator(task_id="alert_of_failed_model")
    compare_all_results >> deploy_new_model
    compare_all_results >> alert_of_failed_model

    new_model = PythonOperator(
        task_id="new_model",
        python_callable=lambda: 5
    )
    new_model >> compare_all_results
    i = 0

    for model in models_to_compare:
        attempt = PythonOperator(
            task_id="old-model-{}".format(model),
            python_callable=lambda: i
        )
        i += 1
        attempt >> compare_all_results
```

Airflow at Rappi

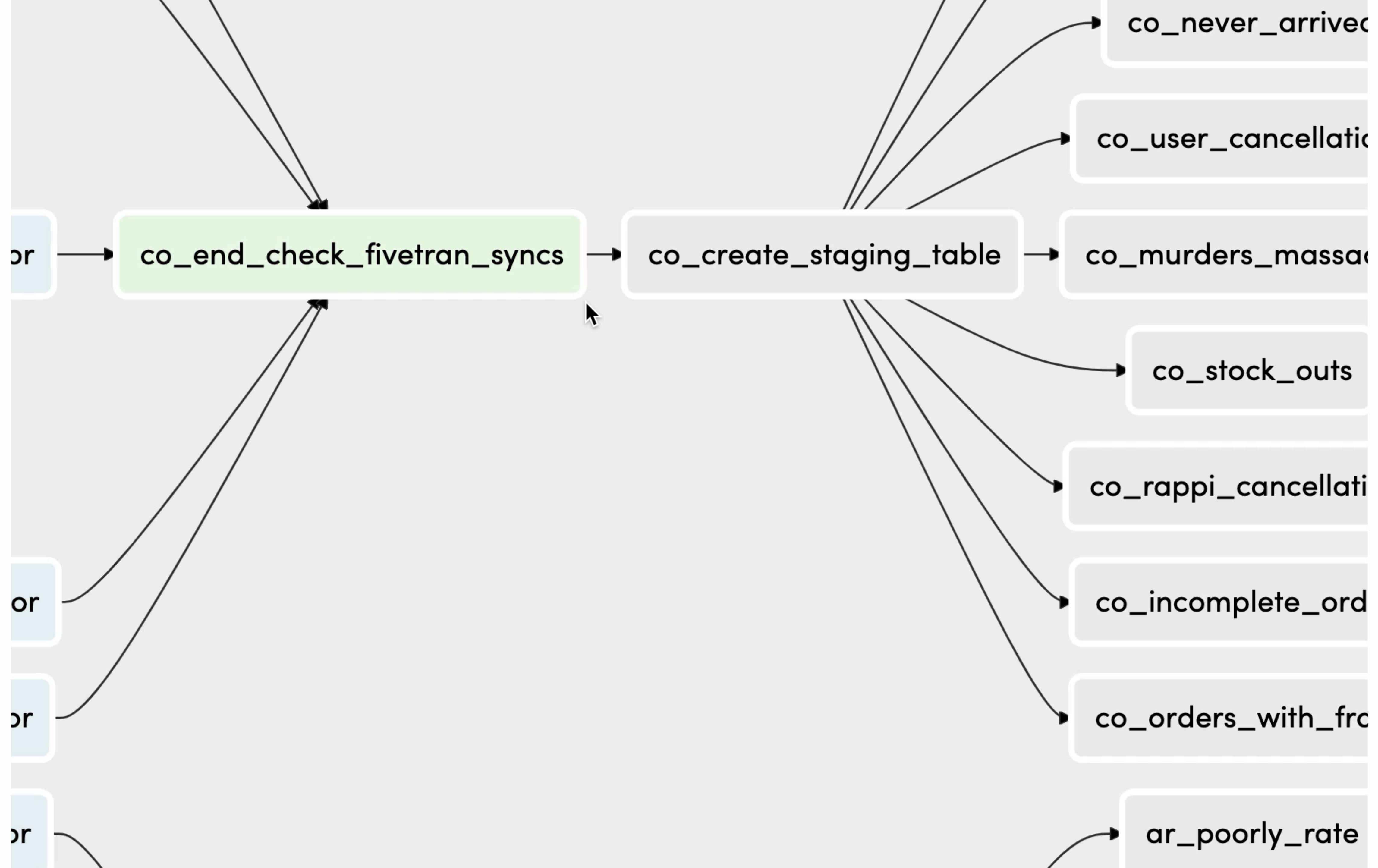


- Take advantage of Airflow extensibility
- RappiFlow: Custom plugins + Astronomer
- SQL + Pandas + Programmatic Pipelines to generalize the DAGs
- Docker Images for Training and Offline Prediction

Airflow at Rappi

Rappi





```

class TrainAndPredictCatalogOperator(KubernetesPodOperator):
    """Train and Predict Recommendation using the Docker Image."""

    @apply_defaults
    def __init__(self, training_params: dict, image_name: str, image_tag_name: str, image_pull_secret: str,
                 k8s_namespace: str, *args, **kwargs):
        """Class Initializer."""
        namespace = k8s_namespace
        image = f'{image_name}:{image_tag_name}'
        name = f'{training_params['country_code']}ml-train-cpgs-catalog'
        super().__init__(
            namespace=namespace, image=image, image_pull_secrets=image_pull_secret, name=name, *args, **kwargs
        )
        self.get_logs = True
        self.in_cluster = True
        self.is_delete_operator_pod = True
        self.env_vars = {
            'COUNTRY': training_params['country_code'],
            'DATE': '{{ ds_nodash }}',
            'AWS_BUCKET': settings.CPGS_DATA_SCIENCE_S3_BUCKET,
        }

        self.resources = Resources(
            request_memory=get_k8s_resources_mapping(training_params['node_size'])['memory'],
            limit_memory=get_k8s_resources_mapping(training_params['node_size'])['memory'],
            request_cpu=get_k8s_resources_mapping(training_params['node_size'])['cpu'],
            limit_cpu=get_k8s_resources_mapping(training_params['node_size'])['cpu']
        )

        self.pool = training_params['pool']

        sf_secrets = Secret(
            deploy_type='env',
            deploy_target=settings.SNOWFLAKE_CPGS_DS_WRITE_CONN_ID,
            secret=settings.SNOWFLAKE_CPGS_DS_WRITE_CONN_ID.replace('_', '-')
        )

        s3_secrets = Secret(
            deploy_type='env',
            deploy_target=settings.S3_CPGS_DS_CONN_ID,
            secret=settings.S3_CPGS_DS_CONN_ID.replace('_', '-')
        )
        self.secrets = [sf_secrets, s3_secrets]

```

Scaling Apache Airflow



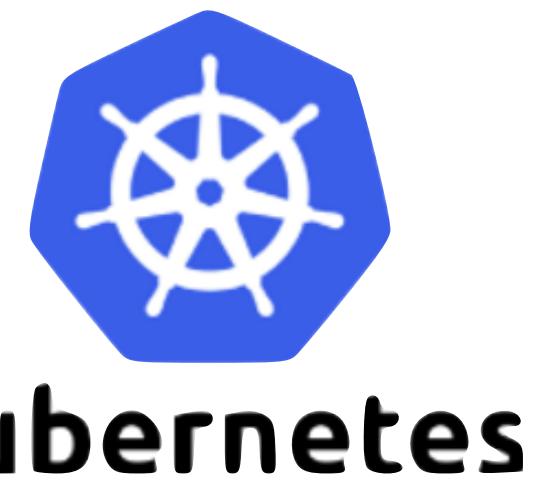
kubernetes

KubernetesExecutor

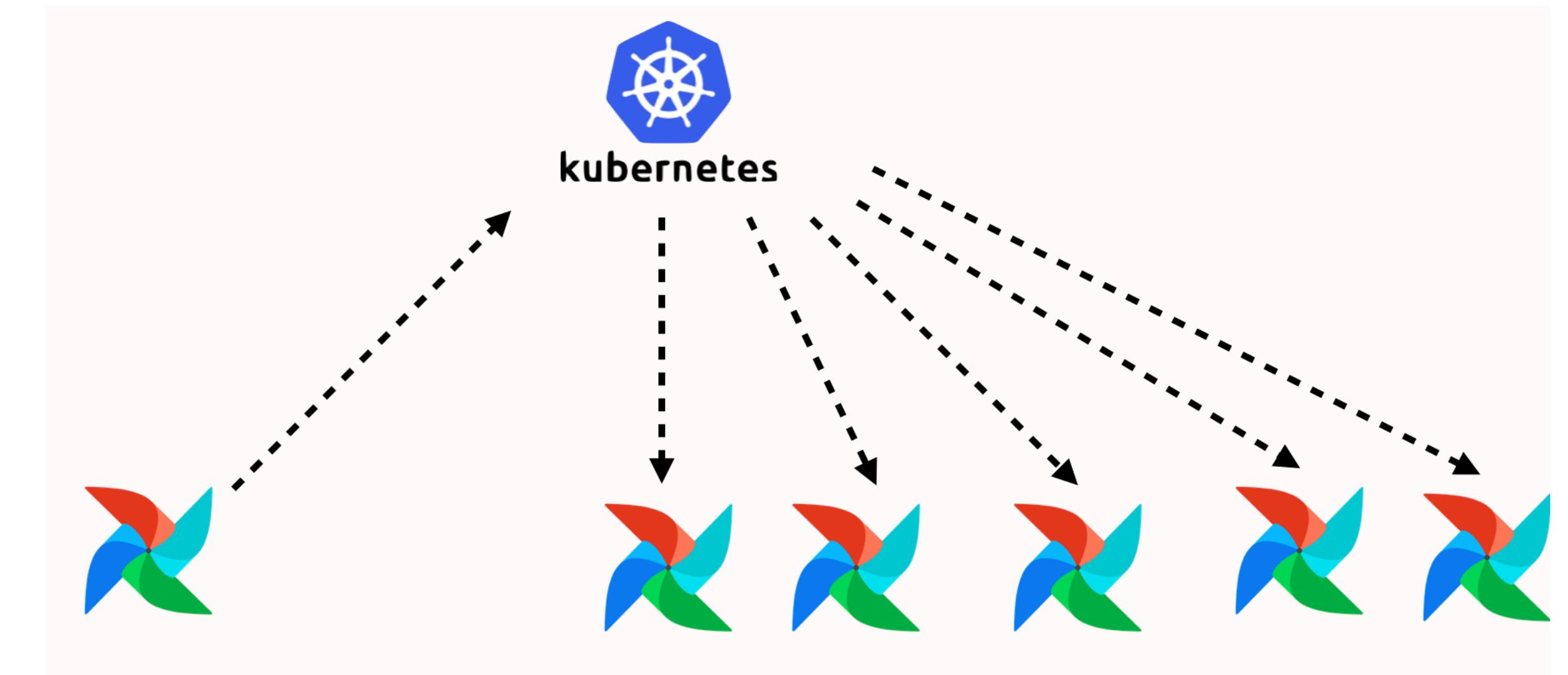


CeleryExecutor

KubernetesExecutor



- Dynamic Allocation
- Fault Tolerance
- Only Requires Kubernetes



CeleryExecutor



- Immediate SLAs
- Multiple tasks per-worker
- Queues for differing requirements

ExecutorConfig



kubernetes

- Users can define elements of tasks in the task description
- Set RAM/CPU requests per-task
- service-accounts for PII
- Node affinity to a GPU-machine

```
t = BashOperator(  
    task_id = 'account-test',  
    bash_command = 'gcloud auth application-default login',  
    dag = dag,  
  
    executor_config = {  
        'request_memory': '128Mi',  
        'limit_memory': '128Mi'  
        'image': 'airflow/scipy:1.1.5'  
        'gcp-service-account' : 'service-account@xxx.iam.gserviceaccount.com'  
    }  
)
```

Deploying Apache Airflow



kubernetes

KubernetesExecutor



CeleryExecutor

Deploying Apache Airflow



kubernetes

KubernetesExecutor



CeleryExecutor

Astronomer



- Disclosure: I work for Astronomer
- Apache Airflow as a Service on Kubernetes
- Easy deployment of K8sExecutor and CeleryExecutor
- Monitoring, support,



Astronomer Local Development



- CLI allows for local airflow server without signing up
- curl -sSL https://install.astronomer.io | sudo bash



```
dimberman astronomer (* | context | default) $
```

```
dimberman astronomer (* | context | default) $  
dimberman astronomer (* | context | default) $ mkdir test  
dimberman astronomer (* | context | default) $ cd test/  
dimberman test (* | context | default) $ █
```

dimberman astronomer (\$ | context | default) \$

```
berman test (⎈ | context default) $ ls
ockerfile      airflow_settings.yaml  dags      include      packages.txt      plugins      requirements.txt
berman test (⎈ | context default) $ vi dags/example-dag.py
berman test (⎈ | context default) $ astro dev start
iv file ".env" found. Loading...
ending build context to Docker daemon 10.75kB
Step 1/1 : FROM astronomerinc/ap-airflow:0.10.1-alpha.14-1.10.5-onbuild
 0.10.1-alpha.14-1.10.5-onbuild: Pulling from astronomerinc/ap-airflow
48c3bd43c5: Already exists
a4c1370ba7: Pull complete
06a6e31dc93: Pull complete
6afb80ad72: Pull complete
3fd51898e3: Pull complete
254078e654: Pull complete
3c2ca9ebe7: Pull complete
76233955f1d: Pull complete
Digest: sha256:5358698d4c0a0644c3fd5ba022f7aaa7793461cfeca49177639fdaf84f184b3e
Status: Downloaded newer image for astronomerinc/ap-airflow:0.10.1-alpha.14-1.10.5-onbuild
Executing 5 build triggers
--> Running in c3d8dbad470d
  fetch http://dl-cdn.alpinelinux.org/alpine/v3.10/main/x86_64/APKINDEX.tar.gz
  fetch http://dl-cdn.alpinelinux.org/alpine/v3.10/community/x86_64/APKINDEX.tar.gz
  fetch http://dl-cdn.alpinelinux.org/alpine/edge/community/x86_64/APKINDEX.tar.gz
  fetch http://dl-cdn.alpinelinux.org/alpine/edge/testing/x86_64/APKINDEX.tar.gz
  fetch http://dl-cdn.alpinelinux.org/alpine/edge/main/x86_64/APKINDEX.tar.gz
: 234 MiB in 69 packages
Removing intermediate container c3d8dbad470d
--> Running in 792f3e843298
WARNING: You are using pip version 19.2.1, however version 19.2.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
Removing intermediate container 792f3e843298
--> 45dc3804cc92
Successfully built 45dc3804cc92
Successfully tagged test_bf7430/airflow:latest
Creating network "testbf7430_airflow" with driver "bridge"
Creating volume "testbf7430_airflow_logs" with driver "local"
Creating volume "testbf7430_postgres_data" with driver "local"
[INFO][0089] [0/3] [postgres]: Starting
[INFO][0090] [1/3] [postgres]: Started
[INFO][0090] [1/3] [scheduler]: Starting
[INFO][0091] [2/3] [scheduler]: Started
[INFO][0091] [2/3] [webserver]: Starting
[INFO][0092] [3/3] [webserver]: Started
Workflow Webserver: http://localhost:8080
Postgres Database: localhost:5432/postgres
The default credentials are admin:admin
berman test (⎈ | context default) $ █
```

Testing Airflow Development

- Once you have a running local instance you can continually change your DAG and compare on WebServer
- Use local docker image to test images for easy iteration
- Treat airflow as an orchestrator and try to keep complex logic in the K8s pod/ SparkJob/Bigquery



Deploying Airflow

- Once the python files have been deployed to the airflow machines, Airflow will read the new values and update
- K8sOperator allows users to bake-in files in Docker image, deploy using git, and share files using a persistent-volume
- We recommend baking into docker images



New to Astronomer?



```
Building image...
Sending build context to Docker daemon 147.5kB
Step 1/5 : FROM astronomerinc/ap-airflow:0.10.1-alpha.15-1.10.5-onbuild
# Executing 5 build triggers
--> Using cache
--> Using cache
--> Using cache
--> Using cache
--> 5a9b7e448924
Step 2/5 : COPY nb_originals /nb_originals
--> 9fa83b3d609f
Step 3/5 : COPY datasets /datasets
--> 8371b7e26901
Step 4/5 : RUN chmod 777 -R /datasets
--> Running in d43360d4d79e
Removing intermediate container d43360d4d79e
--> 505a2a7c177c
Step 5/5 : RUN mkdir /nb_outputs && chmod 777 /nb_outputs
--> Running in ba4b33852642
Removing intermediate container ba4b33852642
--> dc8d45f4eb90
Successfully built dc8d45f4eb90
Successfully tagged exact-spacecraft-1900/airflow:latest
Pushing image to Astronomer registry
The push refers to repository [registry.datarouter.ai/exact-spacecraft-1900/airflow]
c04e606d1092: Pushed
9fe552deb343: Pushed
936003bf8301: Pushed
167af4c9608e: Pushed
3163524aae03: Pushed
403ab3f7e6a8: Layer already exists
96f7b8559712: Layer already exists
ad99845b229c: Layer already exists
19379c7d9566: Layer already exists
2ee04610c8cd: Layer already exists
2465aef6a67e: Layer already exists
d65d405770c8: Layer already exists
58b4b74cd080: Layer already exists
d9729030c916: Layer already exists
63ea16be0911: Layer already exists
d8145dd556fa: Layer already exists
03901b4a2ea8: Layer already exists
cli-7: digest: sha256:61385770922947e44dd2187666d7f84d00b61269be1a1c5665b04bc72810f34a size: 3865
Untagged: registry.datarouter.ai/exact-spacecraft-1900/airflow:cli-7
Untagged: registry.datarouter.ai/exact-spacecraft-1900/airflow@sha256:61385770922947e44dd2187666d7f84
d00b61269be1a1c5665b04bc72810f34a
Deploy succeeded!
dimberman astronomer_pydata (airflow-additions) (✿ |context|default) $ 
```

}

```
File "/usr/local/lib/python3.6/site-packages/airflow/models/__init__.py", line 189, in get_fernet
    for fernet_part in fernet_key.split(',')
File "/usr/local/lib/python3.6/site-packages/airflow/models/__init__.py", line 189, in <listcomp>
    for fernet_part in fernet_key.split(',')
File "/usr/local/lib/python3.6/site-packages/cryptography/fernet.py", line 35, in __init__
    key = base64.urlsafe_b64decode(key)
File "/usr/local/lib/python3.6/base64.py", line 133, in urlsafe_b64decode
    return b64decode(s)
File "/usr/local/lib/python3.6/base64.py", line 87, in b64decode
    return binascii.a2b_base64(s)
binascii.Error: Incorrect padding

During handling of the above exception, another exception occurred:

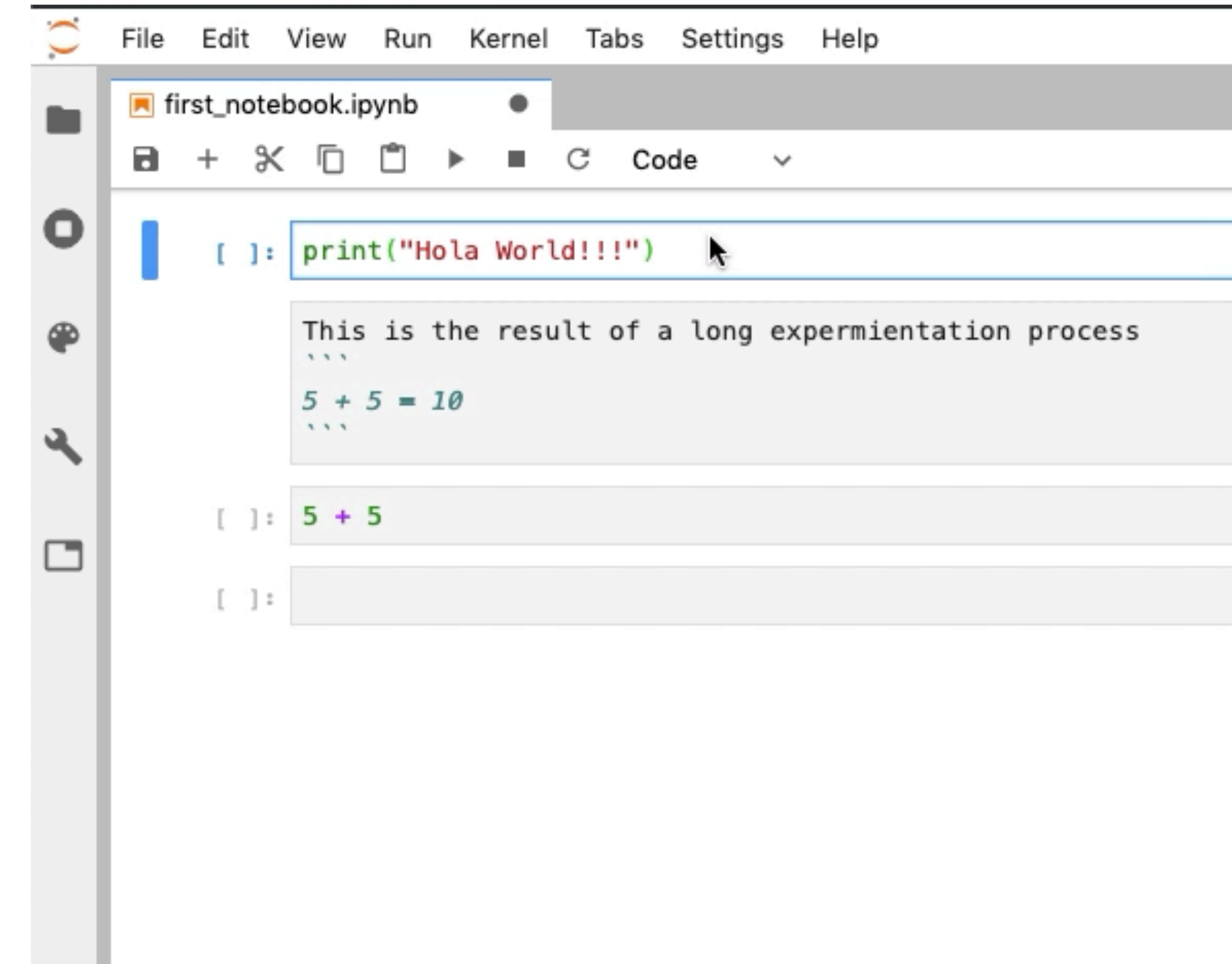
Traceback (most recent call last):
  File "/usr/local/bin/airflow", line 32, in <module>
    args.func(args)
  File "/usr/local/lib/python3.6/site-packages/airflow/bin/cli.py", line 1105, in resetdb
    db.resetdb(settings.RBAC)
  File "/usr/local/lib/python3.6/site-packages/airflow/utils/db.py", line 383, in resetdb
    initdb(rbac)
  File "/usr/local/lib/python3.6/site-packages/airflow/utils/db.py", line 102, in initdb
    schema='default')
  File "<string>", line 4, in __init__
  File "/usr/local/lib/python3.6/site-packages/sqlalchemy/orm/state.py", line 428, in _initialize_instance
    manager.dispatch.init_failure(self, args, kwargs)
  File "/usr/local/lib/python3.6/site-packages/sqlalchemy/util/langhelpers.py", line 67, in __exit__
    compat.reraise(exc_type, exc_value, exc_tb)
  File "/usr/local/lib/python3.6/site-packages/sqlalchemy/util/compat.py", line 277, in reraise
    raise value
  File "/usr/local/lib/python3.6/site-packages/sqlalchemy/orm/state.py", line 425, in _initialize_instance
    return manager.original_init(*mixed[1:], **kwargs)
  File "/usr/local/lib/python3.6/site-packages/airflow/models/connection.py", line 125, in __init__
    self.extra = extra
  File "<string>", line 1, in __set__
  File "/usr/local/lib/python3.6/site-packages/airflow/models/connection.py", line 181, in set_extra
    fernet = get_fernet()
  File "/usr/local/lib/python3.6/site-packages/airflow/models/__init__.py", line 193, in get_fernet
    raise AirflowException("Could not create Fernet object: {}".format(ve))
airflow.exceptions.AirflowException: Could not create Fernet object: Incorrect padding
$ (pydata) dimberman ~ (✿ |context|default) $ 
```

Session Contents Restored

```
Last login: Mon Sep 23 08:14:03 on ttys000
dimberman ~ (✿ |context|default) $ 
```

TL;DR...

- Iteratively built a working model using JupyterNotebooks
- Parameterized our notebooks using Papermill
- Tuned, tested, compared, and productionized our model using Apache Airflow
- Scaled our model building using Kubernetes or Celery



The screenshot shows a Jupyter Notebook interface with the following details:

- File Menu:** File, Edit, View, Run, Kernel, Tabs, Settings, Help.
- Toolbar:** Includes icons for file operations like Open, Save, and Close, along with a "Code" dropdown.
- Code Cell:** The first cell contains the Python code: `[]: print("Hola World!!!")`. The output of this cell is displayed below it: "This is the result of a long experimentation process" followed by three ellipses and the calculation `5 + 5 = 10`.
- Output Cells:** Below the first cell, there are two more cells. The second cell contains the expression `5 + 5`, and the third cell is currently empty.

Questions?

Thank you!



Daniel Imberman
Platform Engineer
[@danimberman](https://twitter.com/danimberman)
[@d_imberz](https://www.instagram.com/d_imberz)
www.astronomer.io



Gonzalo Diaz
Senior Backend Data Engineer
[@gdcor](https://twitter.com/gdcor)
[@gondcor](https://www.instagram.com/gondcor)
www.rappi.com.ar
Join us! gonzalo.diaz@rappi.com