

Data Related Roles in today's Businesses

What possibilities are there after your PhD!

@Piotr Laczkowski

piotr.laczkowski@gmail.com



Polska
(Poznan)



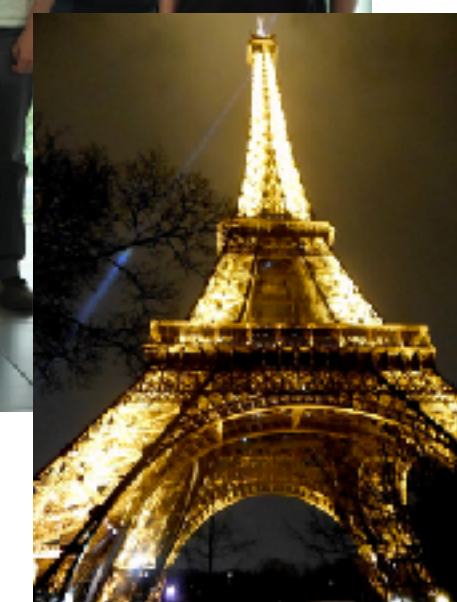
Grenoble
2005



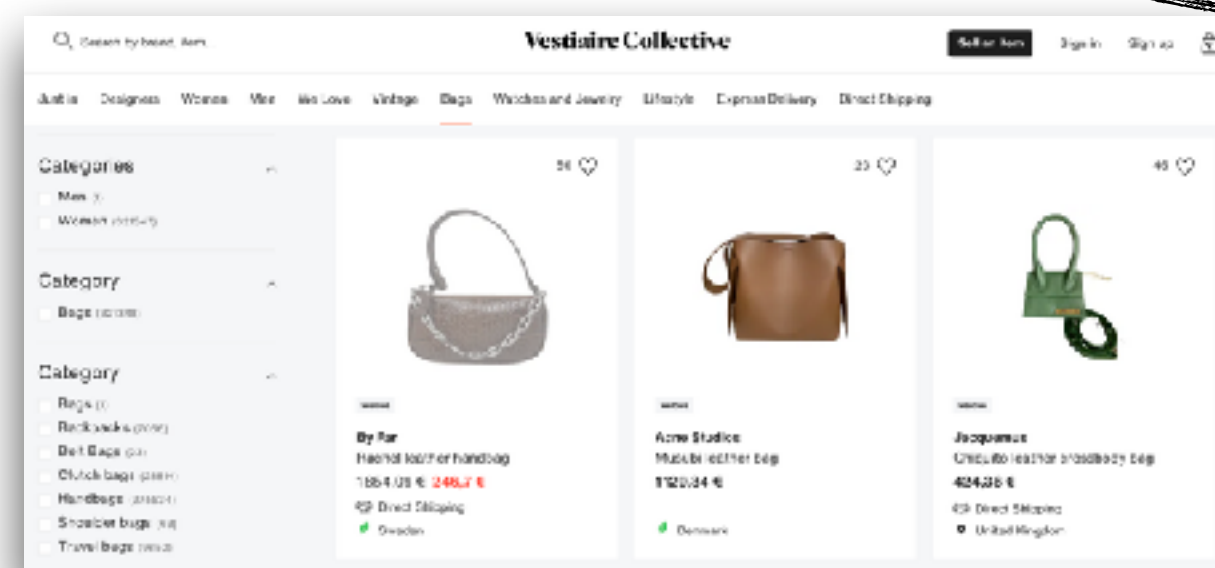
PhD (CEA)
2009-2012



PARIS
2012



Vestiaire Collective
2019



backmarket
2016



(1) You are THE ONLY ONE responsible for your path !!!

-> think about it already (plan it)

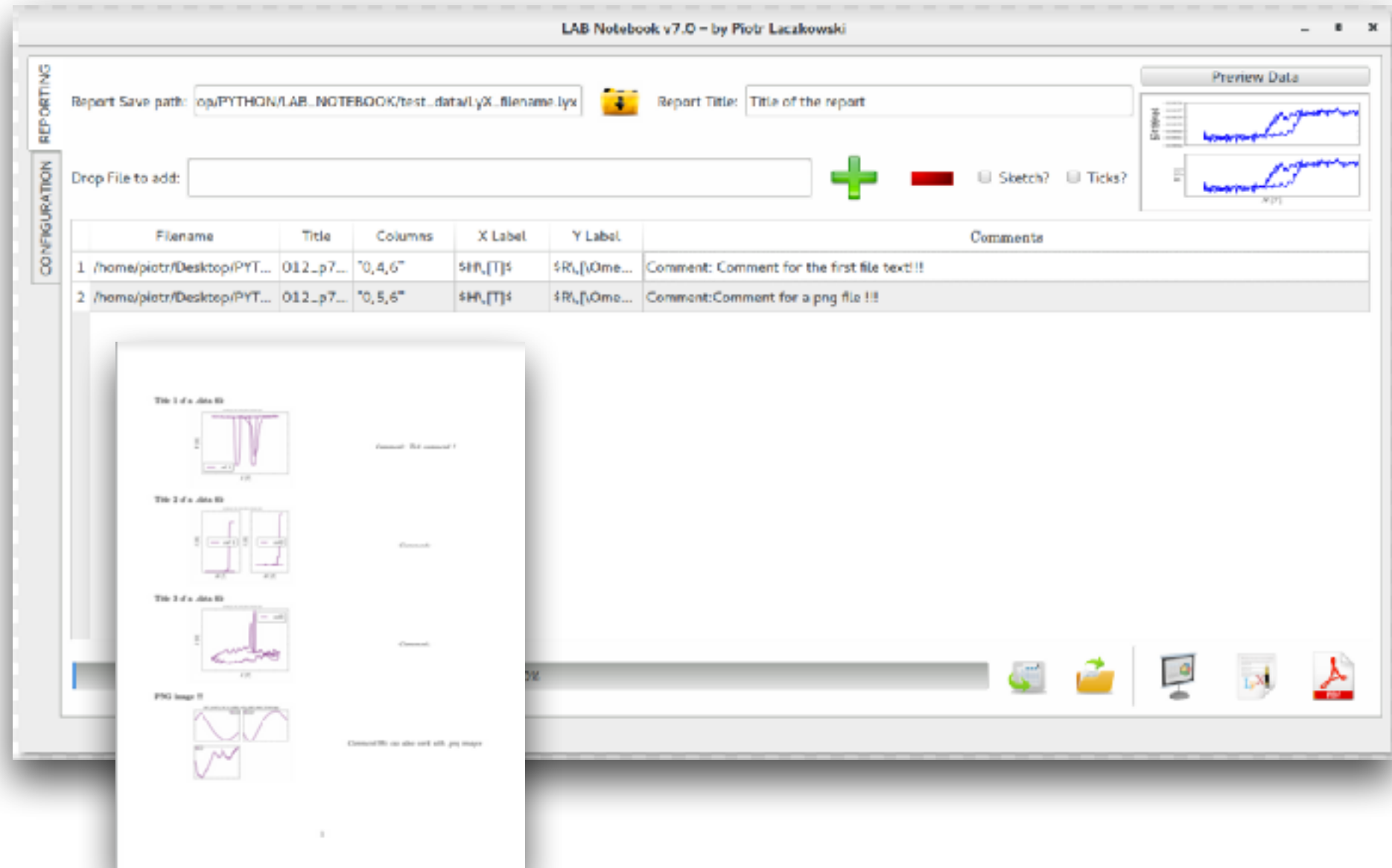
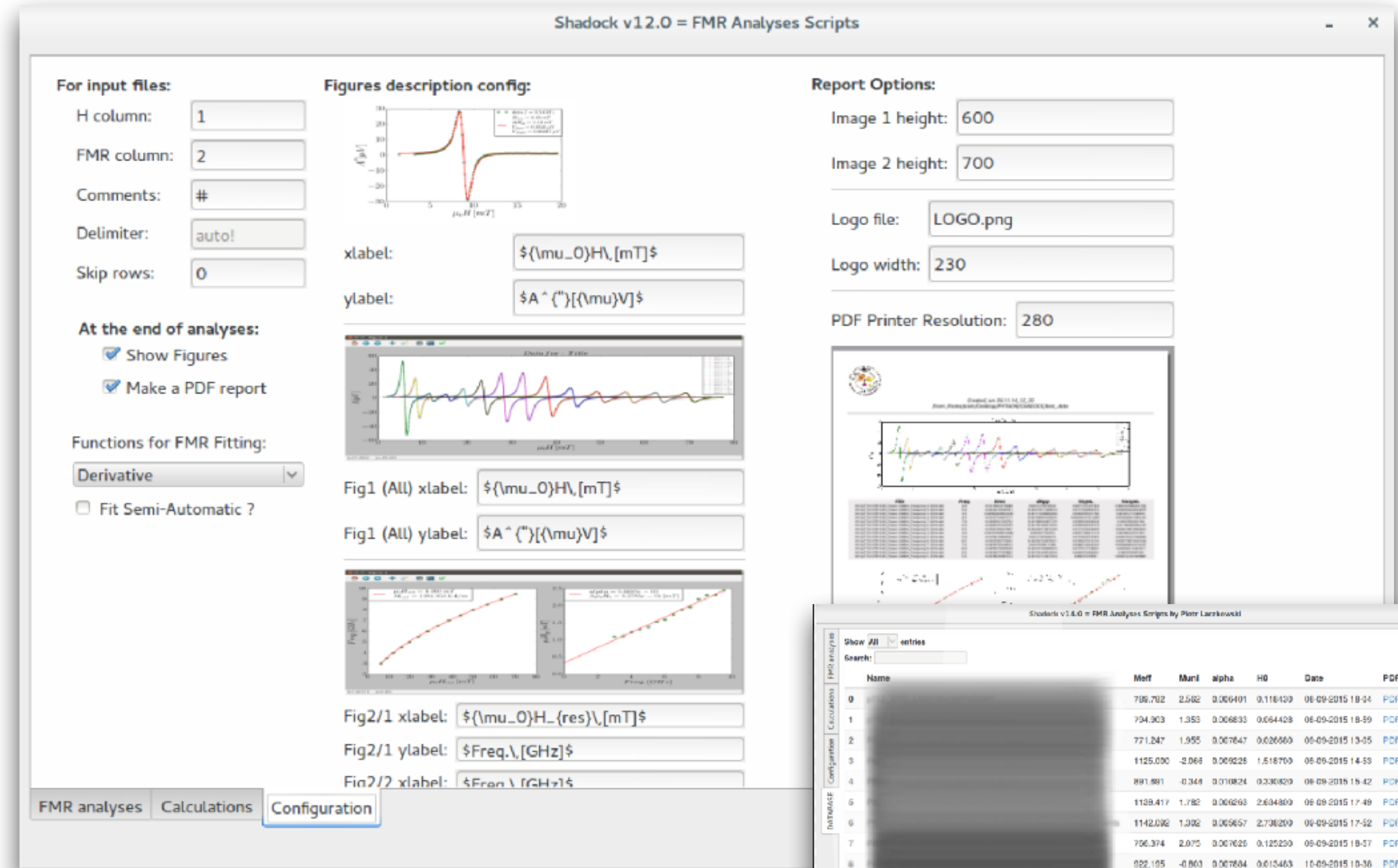
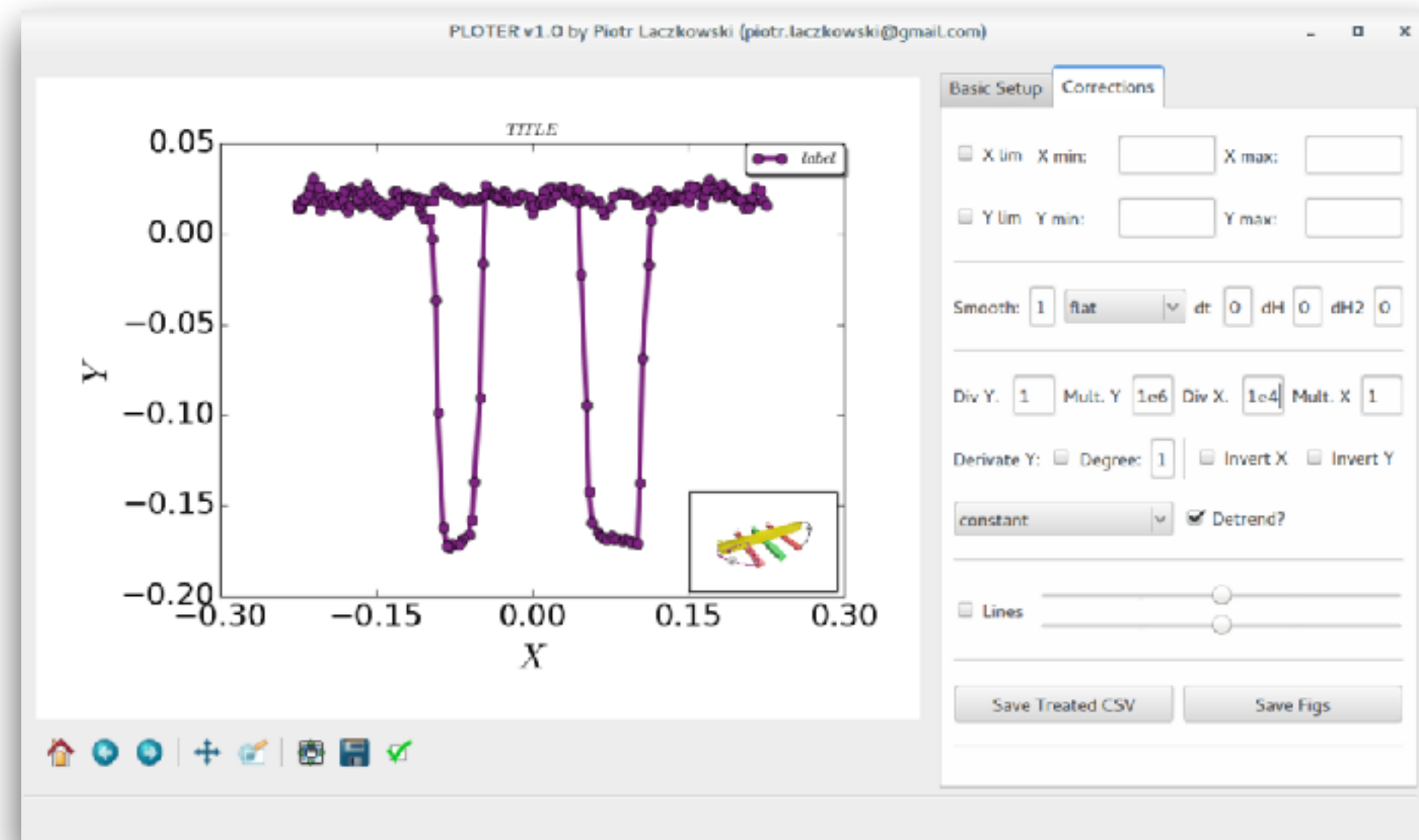
(2) use this time (PhD, ...) to discover what do you like to do

-> HAVE FUN doing it and gain skills !

**(3) If you have skills -> you DO NOT NEED any network not luck to
find the job, it can find you instead!**

PhD + Post-Doc Side Projects

Examples of Python side projects



Shaddock v12.0 = FMR Analyses Scripts by Piotr Laczowski

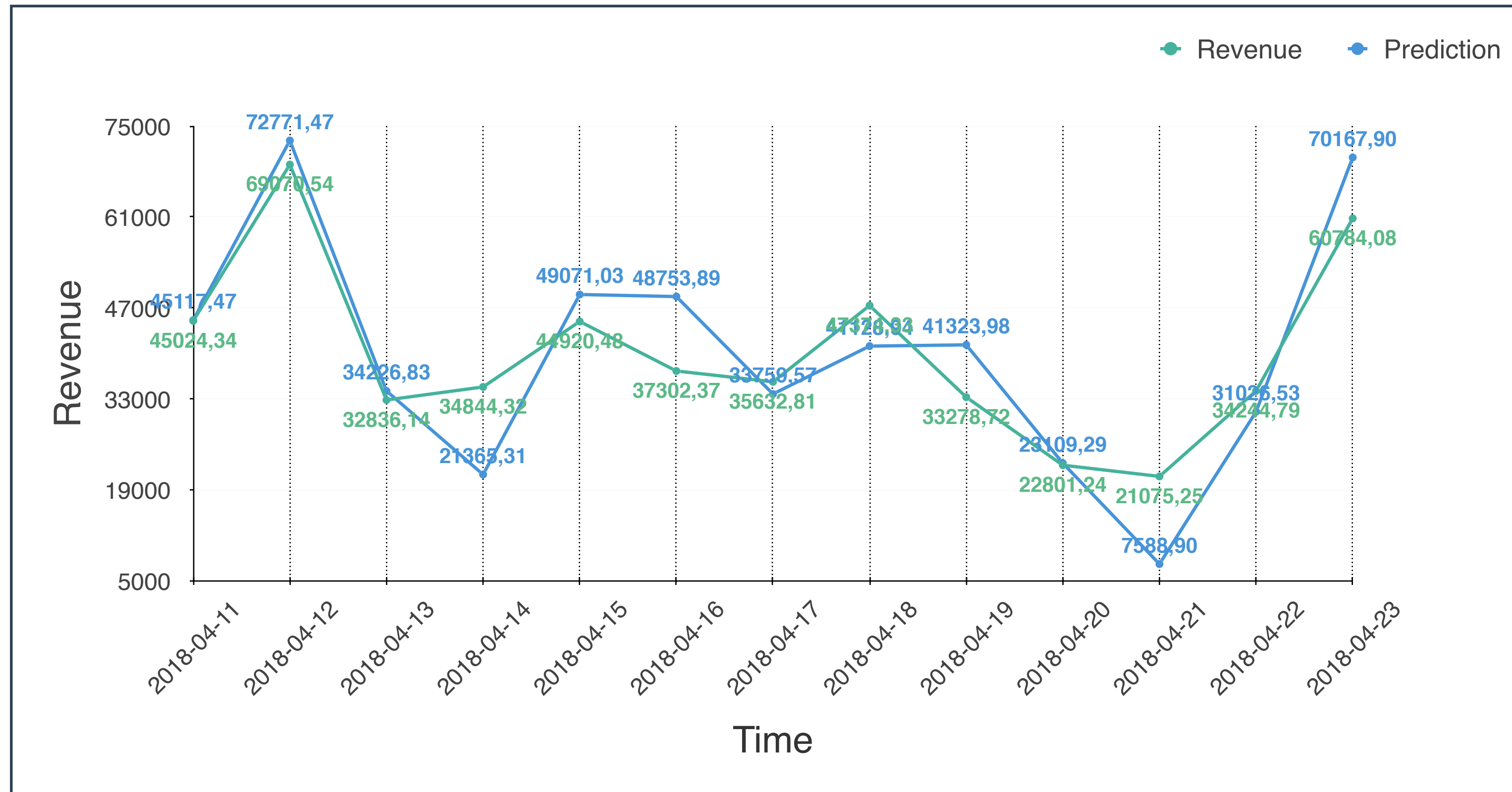
Show: All entries

Search:

Name	Meff	Munl	alpha	H0	Date	PDF_Report
0	786.702	2.562	0.006401	0.116400	06-09-2015 18:34	PDF
1	794.303	1.353	0.006833	0.064428	06-09-2015 18:59	PDF
2	771.247	1.955	0.007647	0.026560	06-09-2015 13:05	PDF
3	1125.000	-2.068	0.009228	1.518700	06-09-2015 14:53	PDF
4	881.881	-0.348	0.010824	0.330820	06-09-2015 15:42	PDF
5	1139.417	1.782	0.006265	2.634300	06-09-2015 17:49	PDF
6	1142.092	1.302	0.005657	2.736200	06-09-2015 17:52	PDF
7	796.374	2.075	0.007625	0.125250	06-09-2015 18:57	PDF
8	922.195	-0.863	0.007684	0.015463	10-09-2015 10:36	PDF
9	1083.638	2.997	0.006437	0.943680	10-09-2015 10:41	PDF
10	1146.088	0.889	0.005853	1.624900	10-09-2015 11:59	PDF
11	1147.435	1.590	0.005231	0.826590	10-09-2015 13:38	PDF
12	782.335	2.450	0.007420	0.172090	11-09-2015 10:33	PDF
13	796.303	2.610	0.007521	0.116370	11-09-2015 10:36	PDF
14	784.778	1.583	0.007760	-0.003475	11-09-2015 11:16	PDF
15	777.363	2.525	0.007795	0.164350	11-09-2015 14:50	PDF
16	1160.718	0.534	0.003871	2.810300	11-09-2015 15:44	PDF

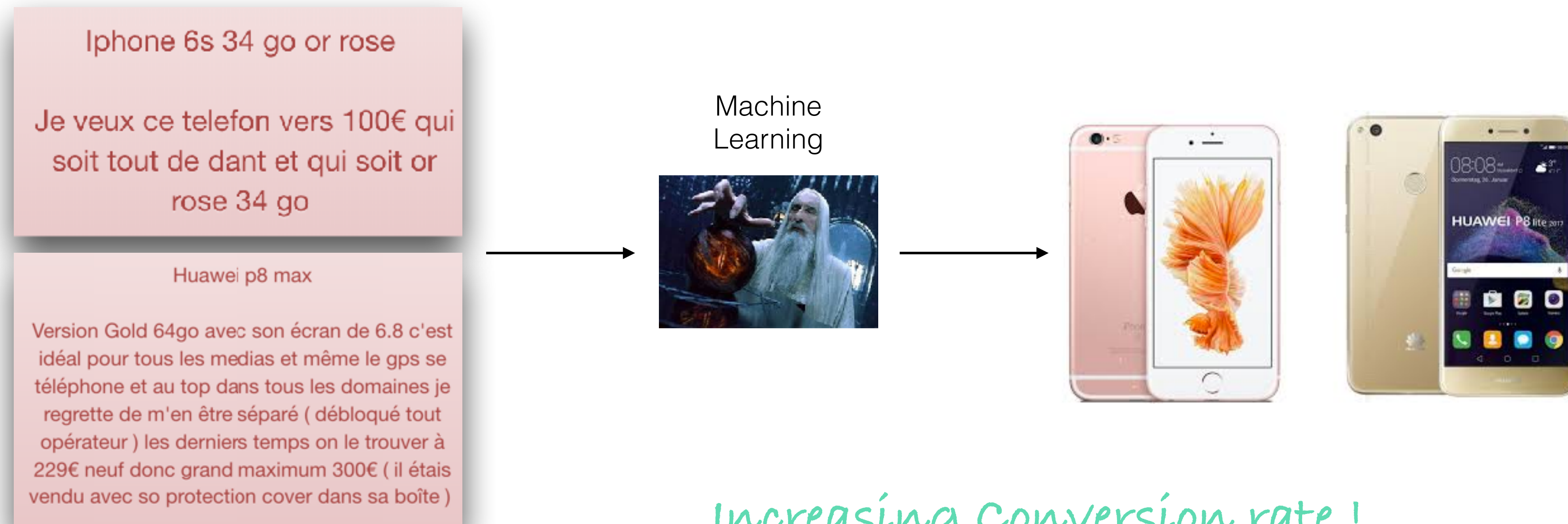
Back Market Projects

Predicting revenues (GMV) based on weather

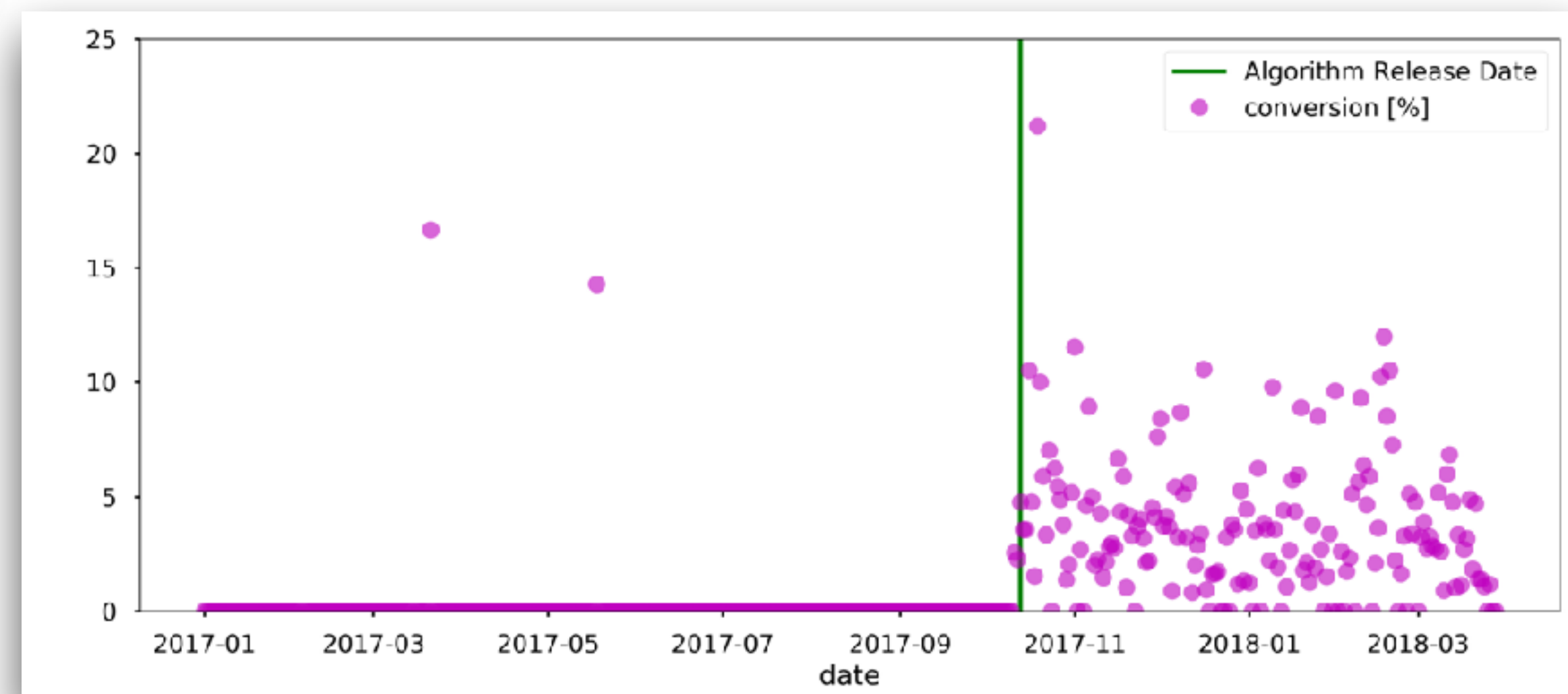


Predicting GMV while looking if its rains or shines with accuracy > 80%

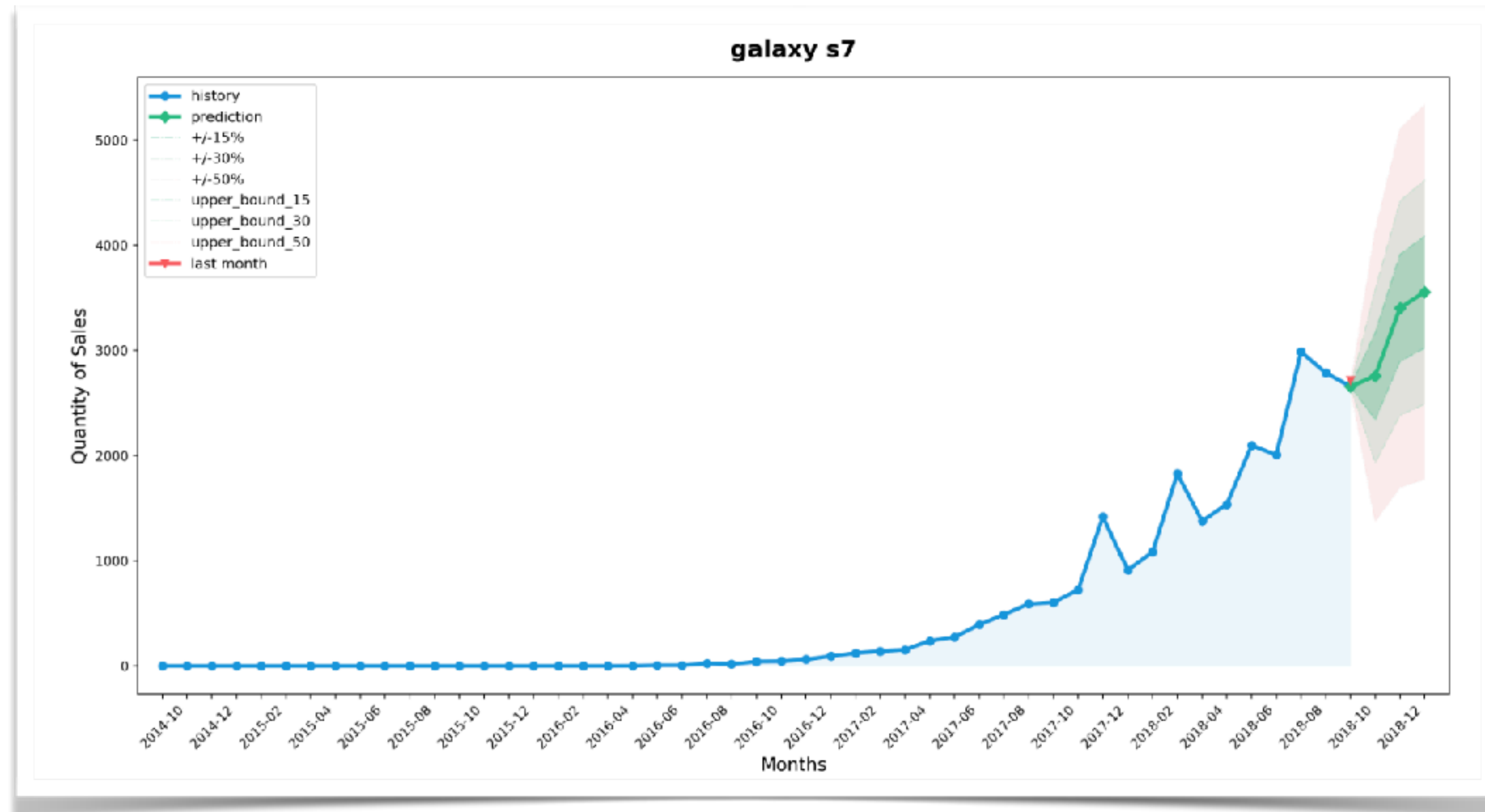
Clients Alerts Matching system



Increasing Conversion rate !



Products sales perditions - supply strategy planning!



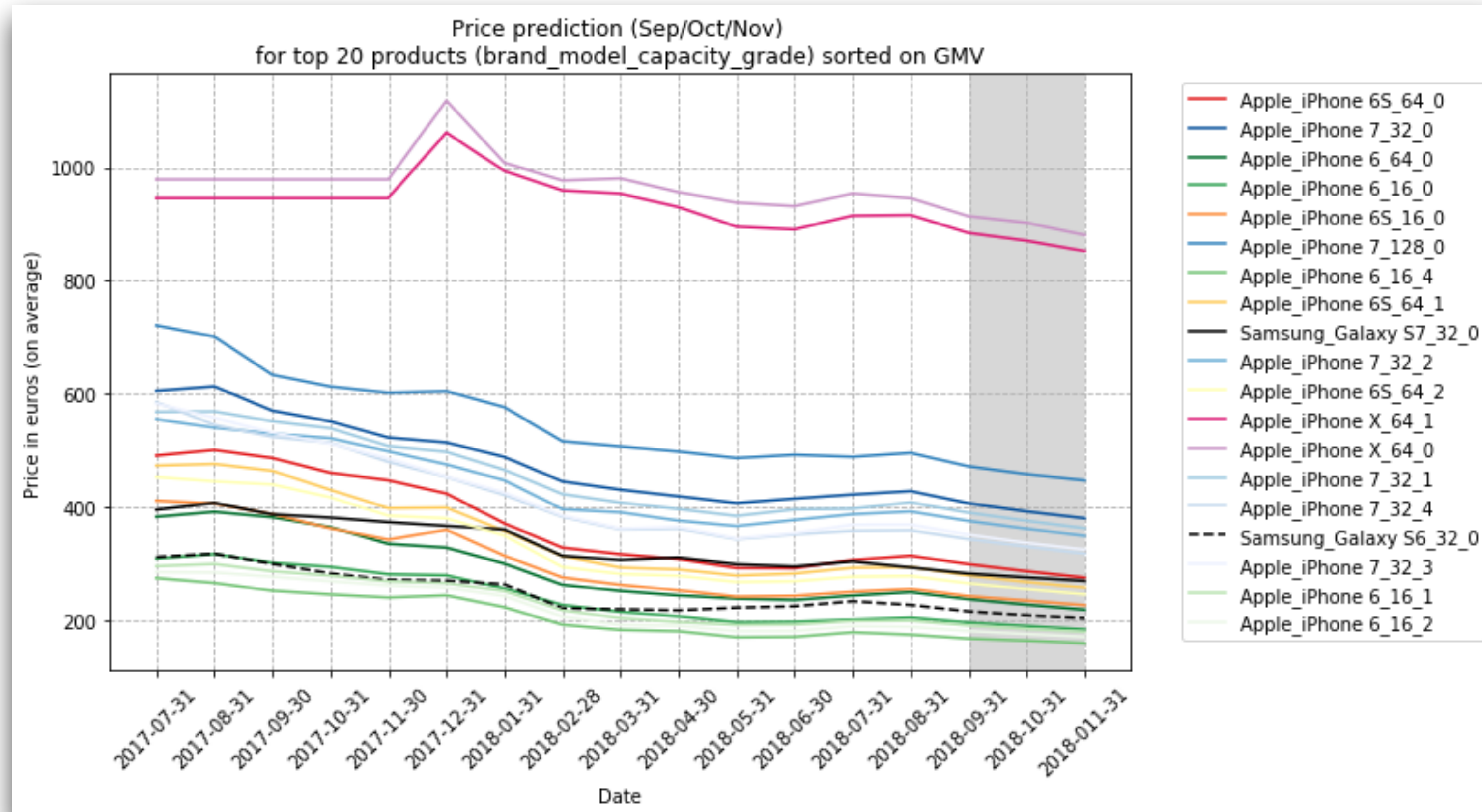
85% of the GMV covered with error < 5%

95% trends
perfectly predicted

+

85%
GMV Covered

Prices Predictions



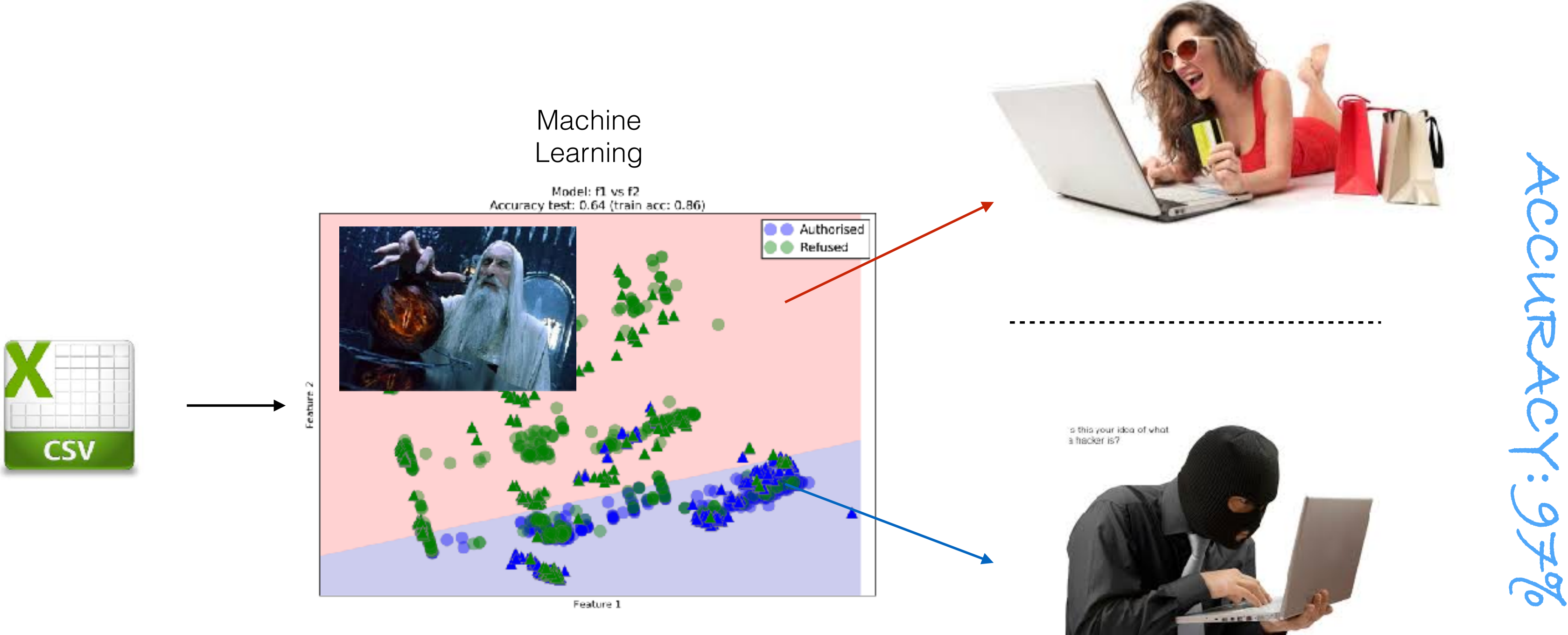
85% of the GMV covered with error < 5%

95.58% GMV trends
correctly predicted

+

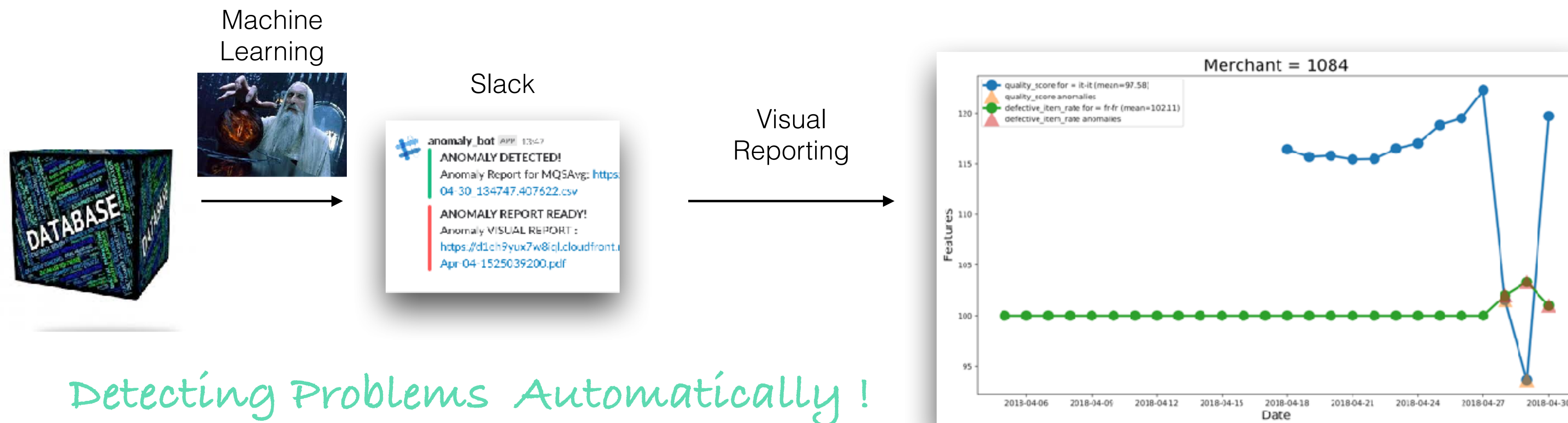
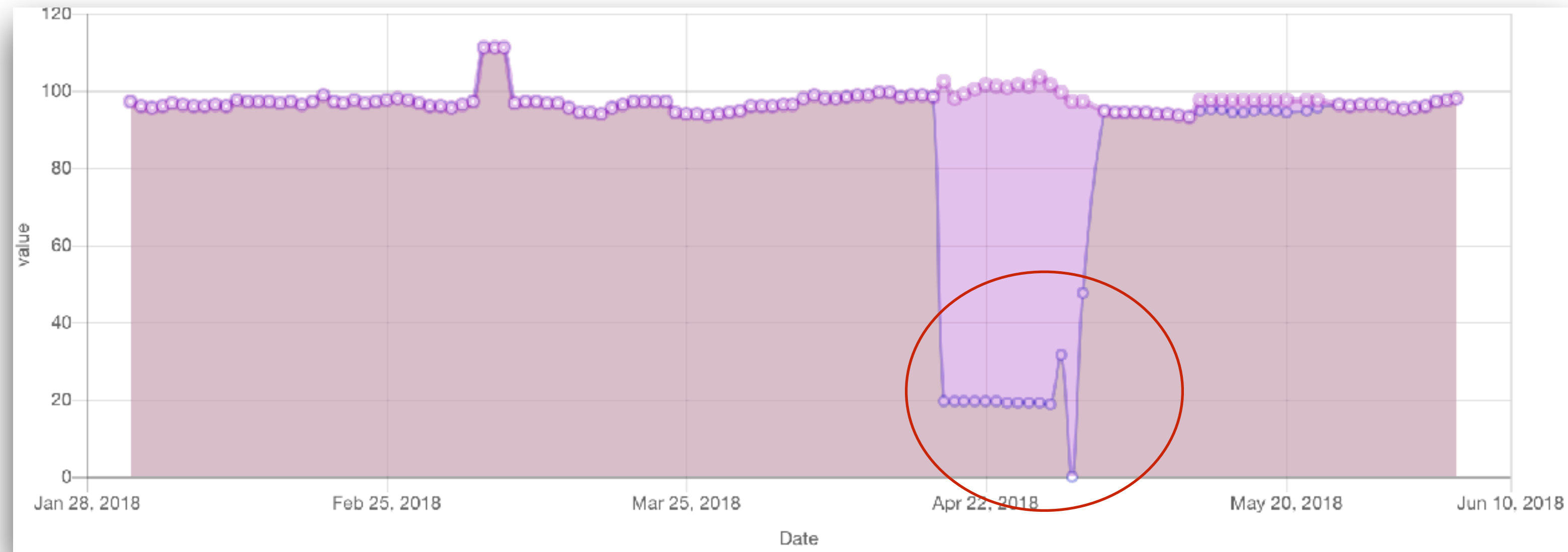
80%
GMV covered

Fraud detection

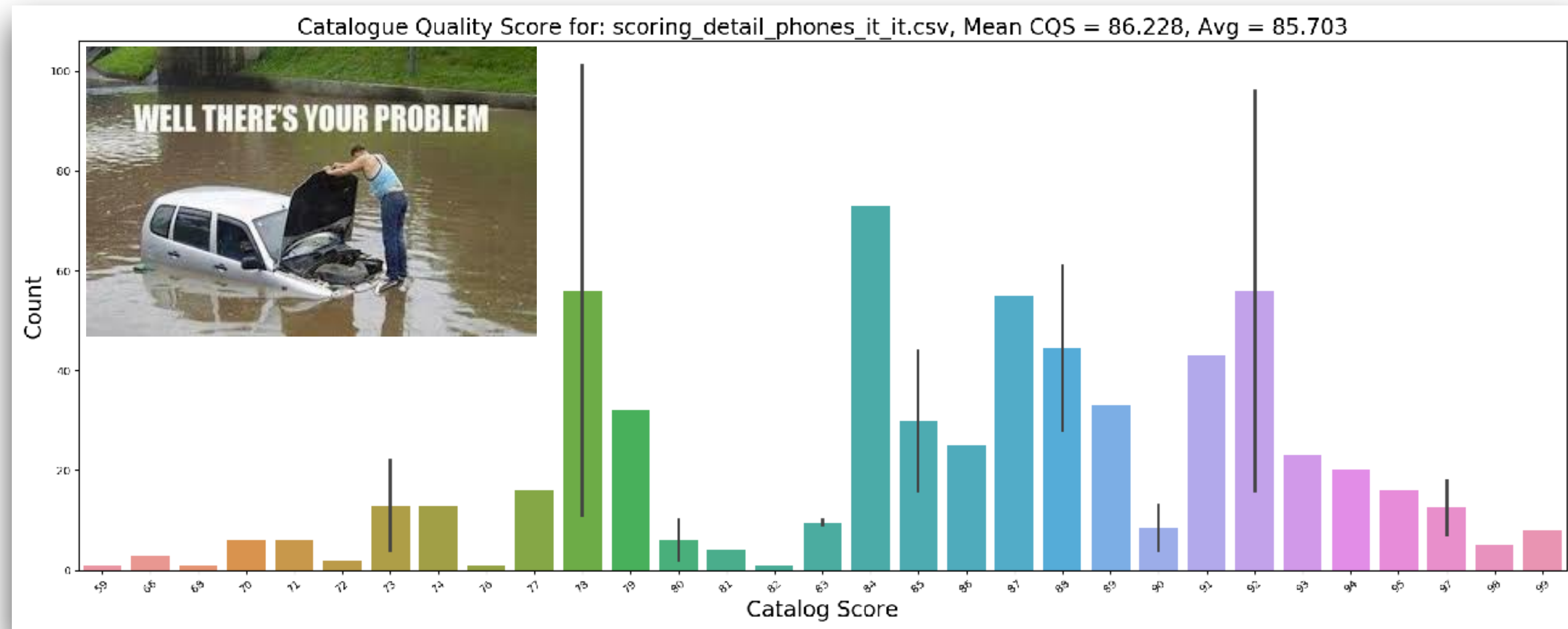


Preventing Fraud !

Anomalies Detection System



Catalogue Quality Score

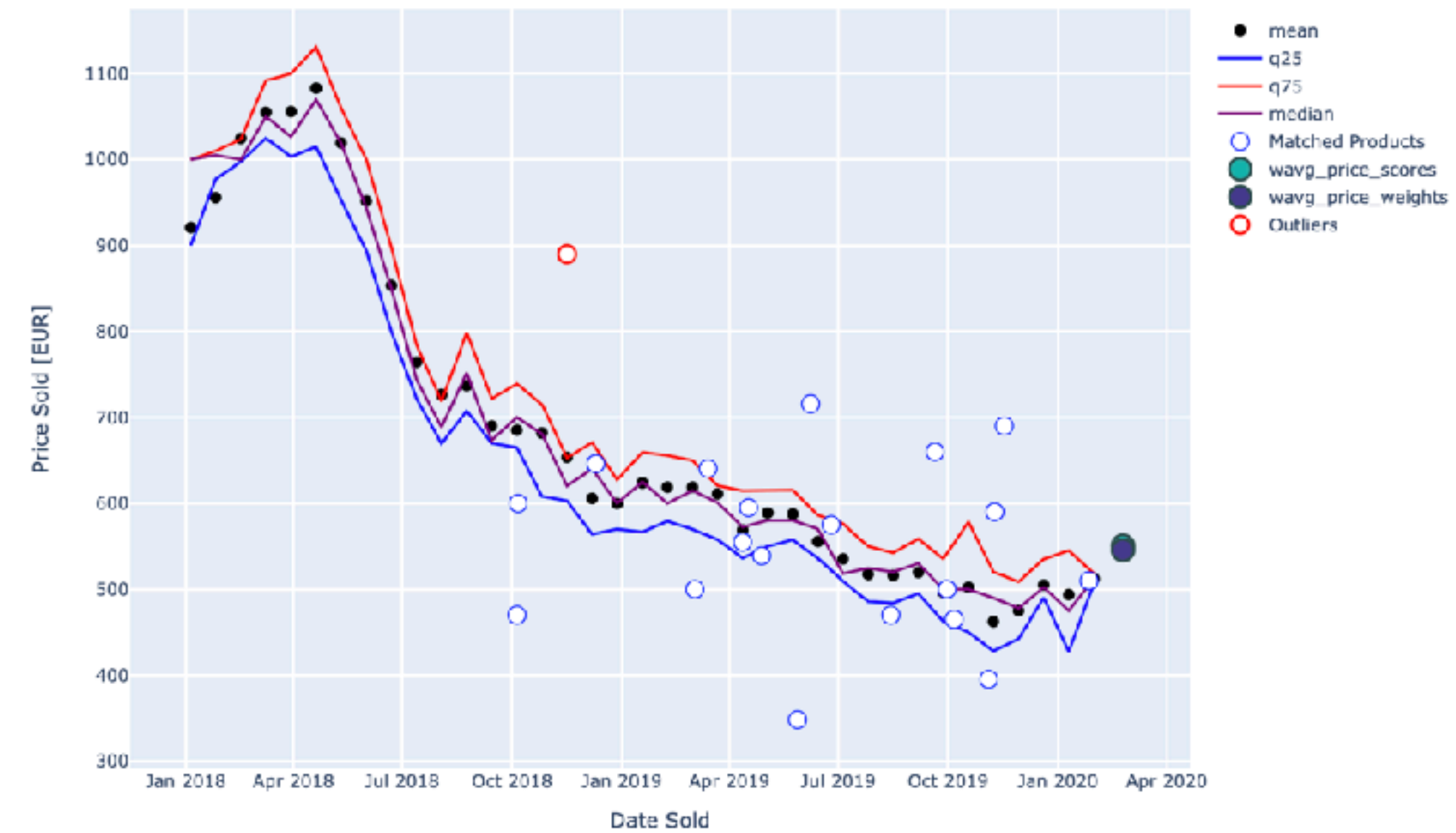
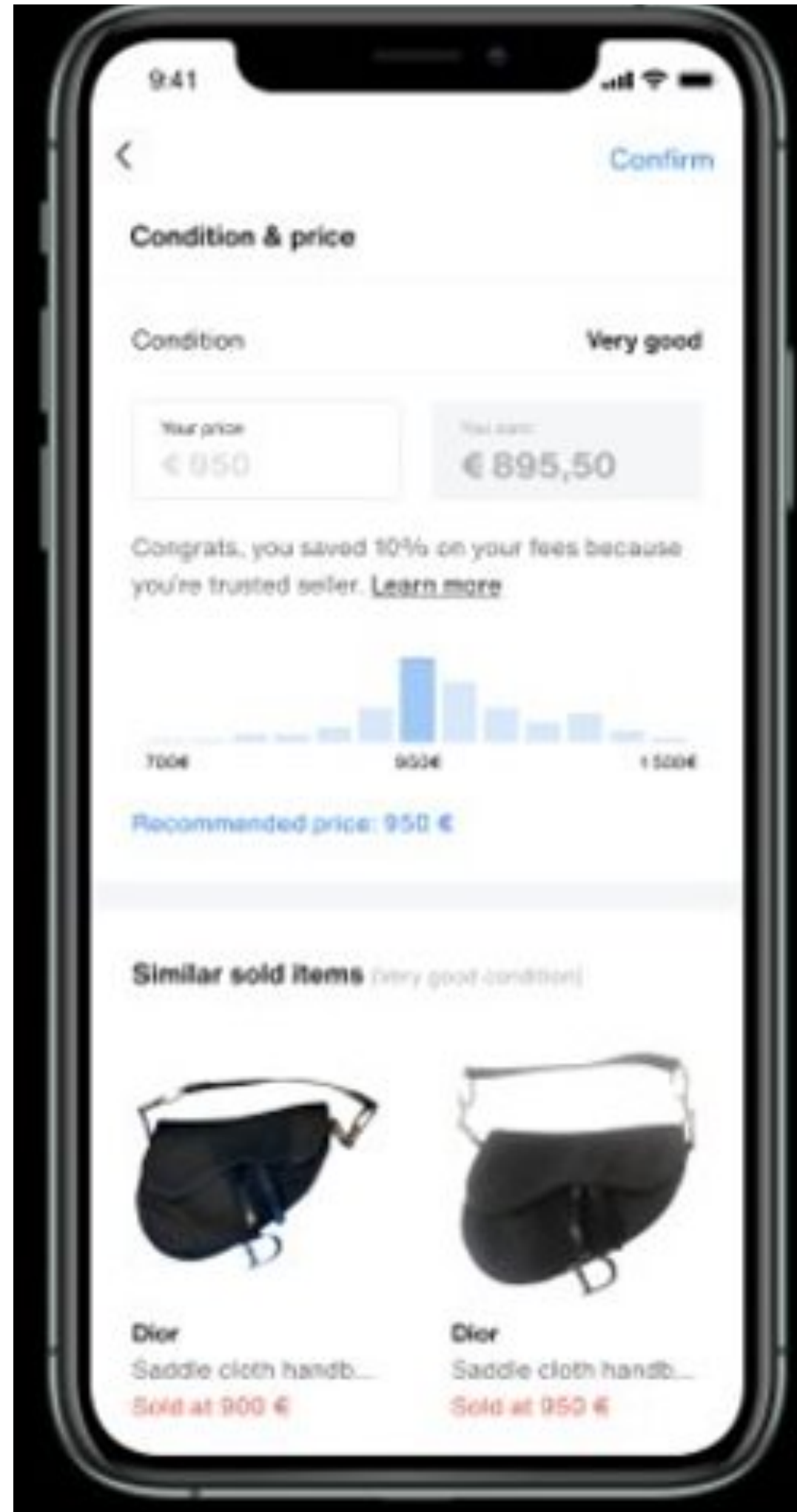


product_id	priority	score	title	color	connector	contrast	date_release	double_sim	image1_url	image2_url	image3_url	image4_url
3207	6.999140589454527	92.25352112676056	iPhone 6S 64 GB - Oro Rosa - sbloccato da tutti gli operatori	9.0	1.0	1.0	8.0	0.0	10.0	8.0	6.0	4.0
3206	6.657680361015061	79.5774647887324	iPhone 6S 16 GB - Oro Rosa - sbloccato da tutti gli operatori	9.0	1.0	1.0	8.0	0.0	10.0	0.0	0.0	0.0
1859	5.025024012941712	92.25352112676056	iPhone 6 64 GB - Grigio Siderale - sbloccato da tutti gli operatori	9.0	1.0	1.0	8.0	0.0	10.0	8.0	6.0	4.0
15023	4.13774834437086	88.02816901408451	iPhone 7 128 GB - Nero - sbloccato da tutti gli operatori	9.0	1.0	1.0	8.0	0.0	10.0	8.0	0.0	4.0
2225	2.8071562716220226	94.36619718309859	Samsung Galaxy Note 4 32 GB - nero - sbloccato da tutti gli operatori	9.0	1.0	1.0	8.0	0.0	10.0	8.0	6.0	4.0
3215	2.512512006470856	92.25352112676056	iPhone 6S 64 GB - Grigio Siderale - sbloccato da tutti gli operatori	9.0	1.0	1.0	8.0	0.0	10.0	8.0	6.0	4.0

Controlling/Increasing Data Quality !

Vestiaire Collective Projects

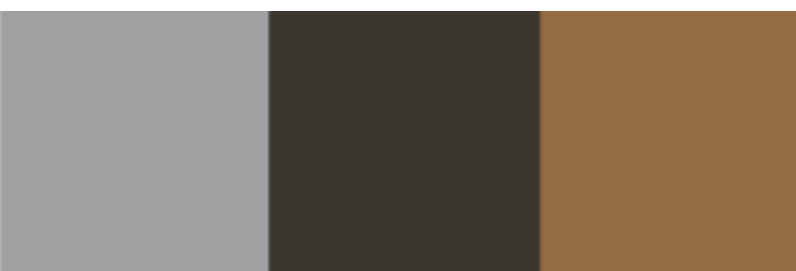
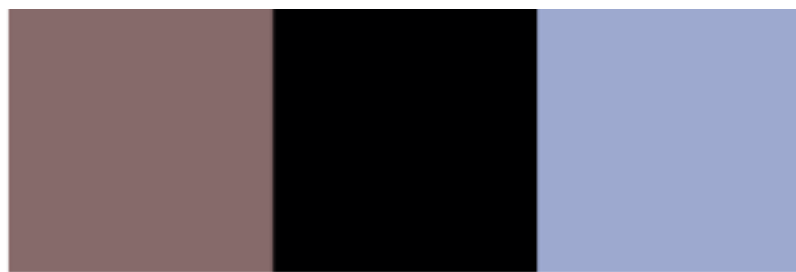
Pricing “unique” items



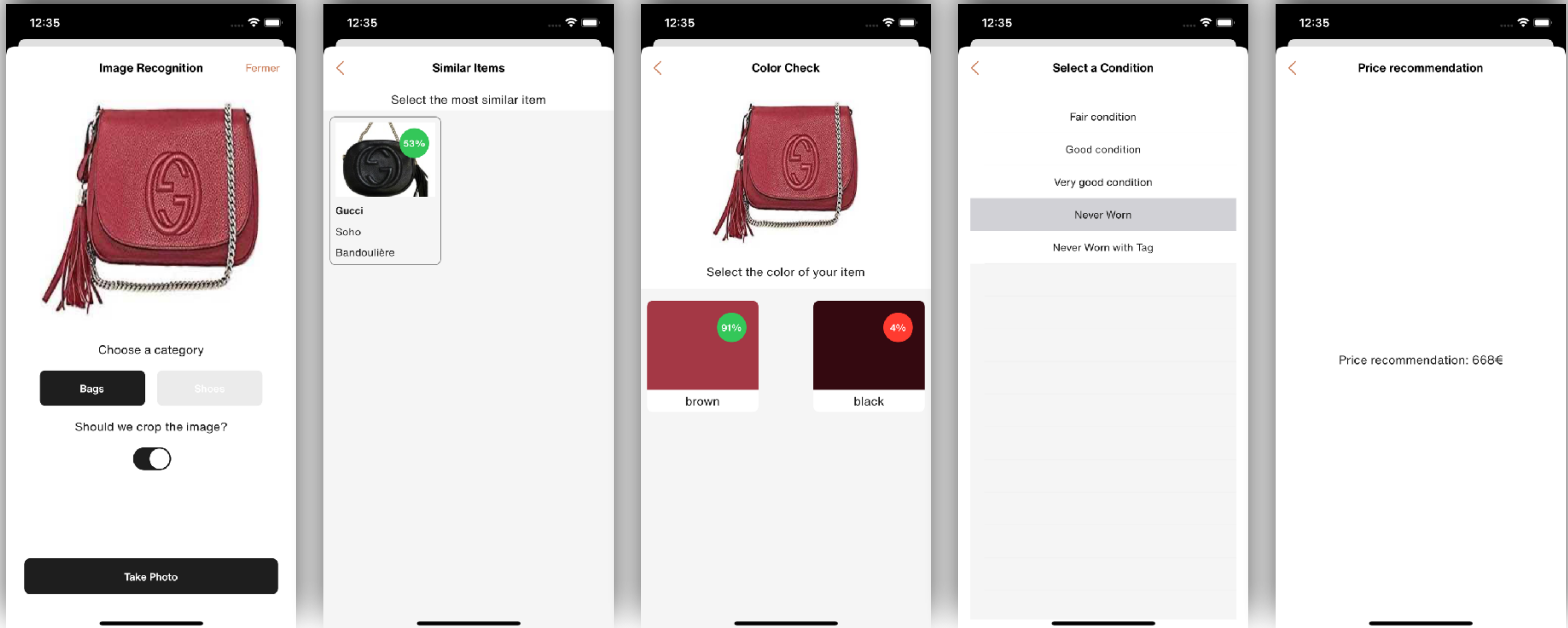
Semantic Image segmentation -> extracting objects!



Extracting Colours from the image

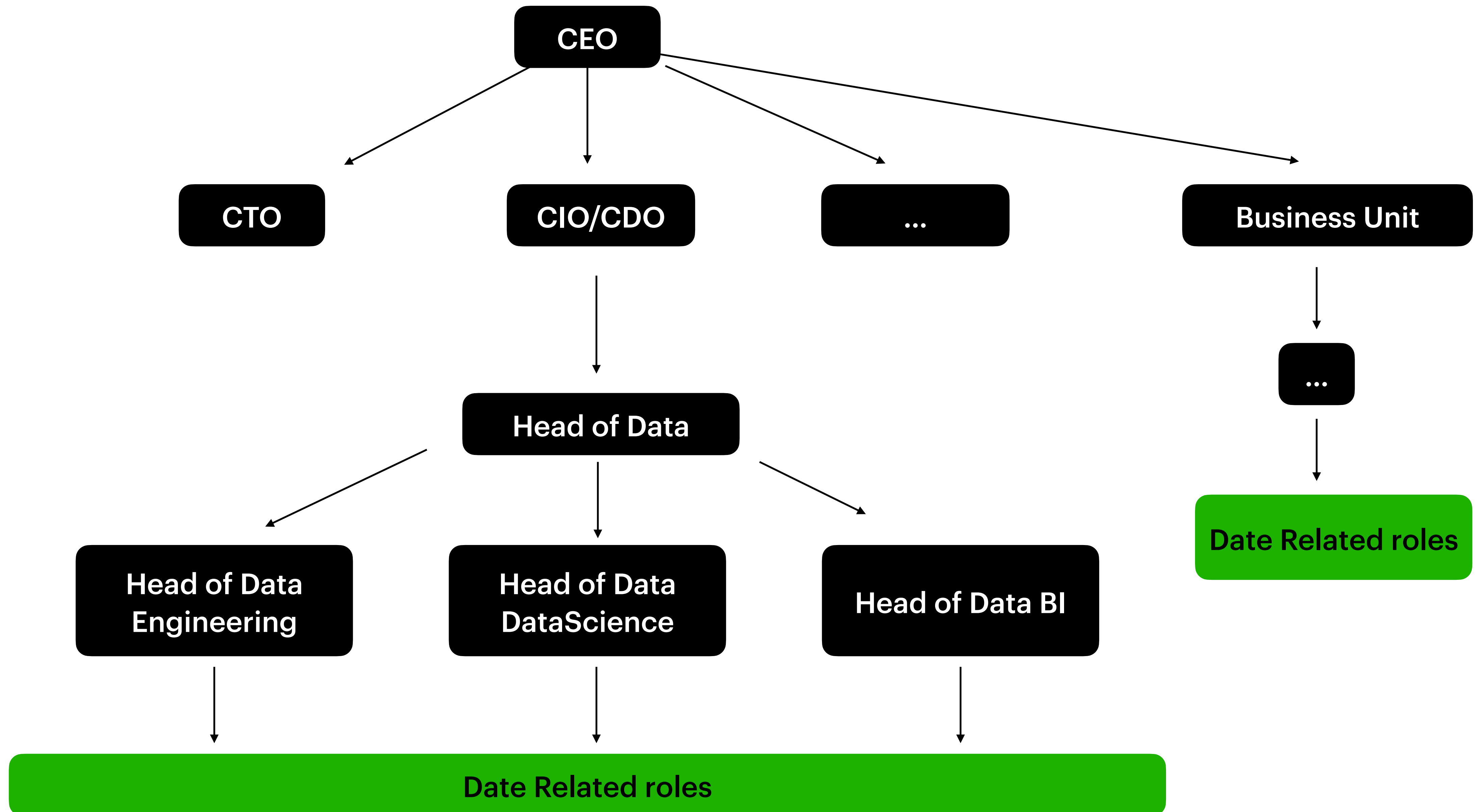


Reverse image search engine

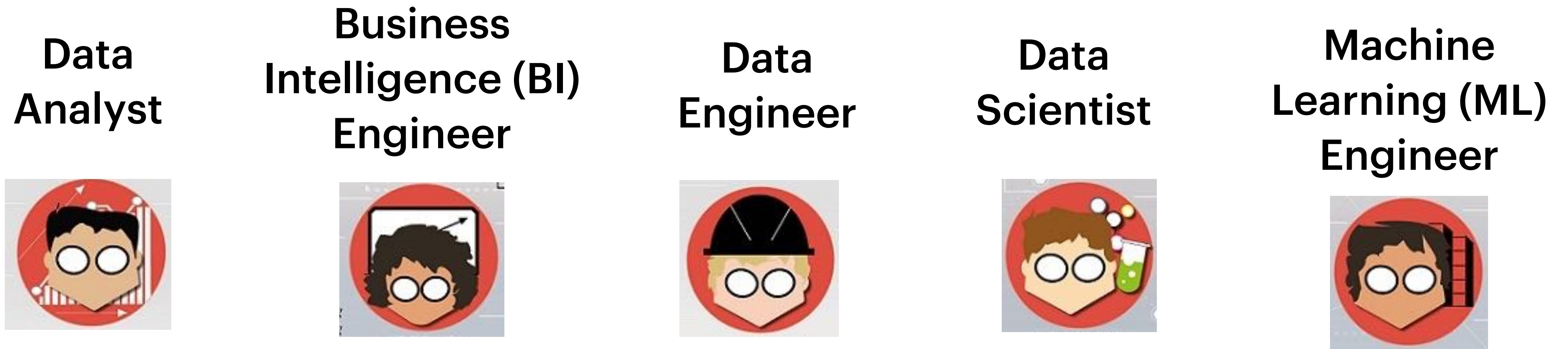


Data Related Roles

Data Oriented Roles in Business - typical tech organisation



Data Oriented Roles in Business



Typical "Data" career progression

Data Analyst



Role:

Analyze the data to respond to business related questions
Preparing reports and visualise data (storytelling)

Questions:

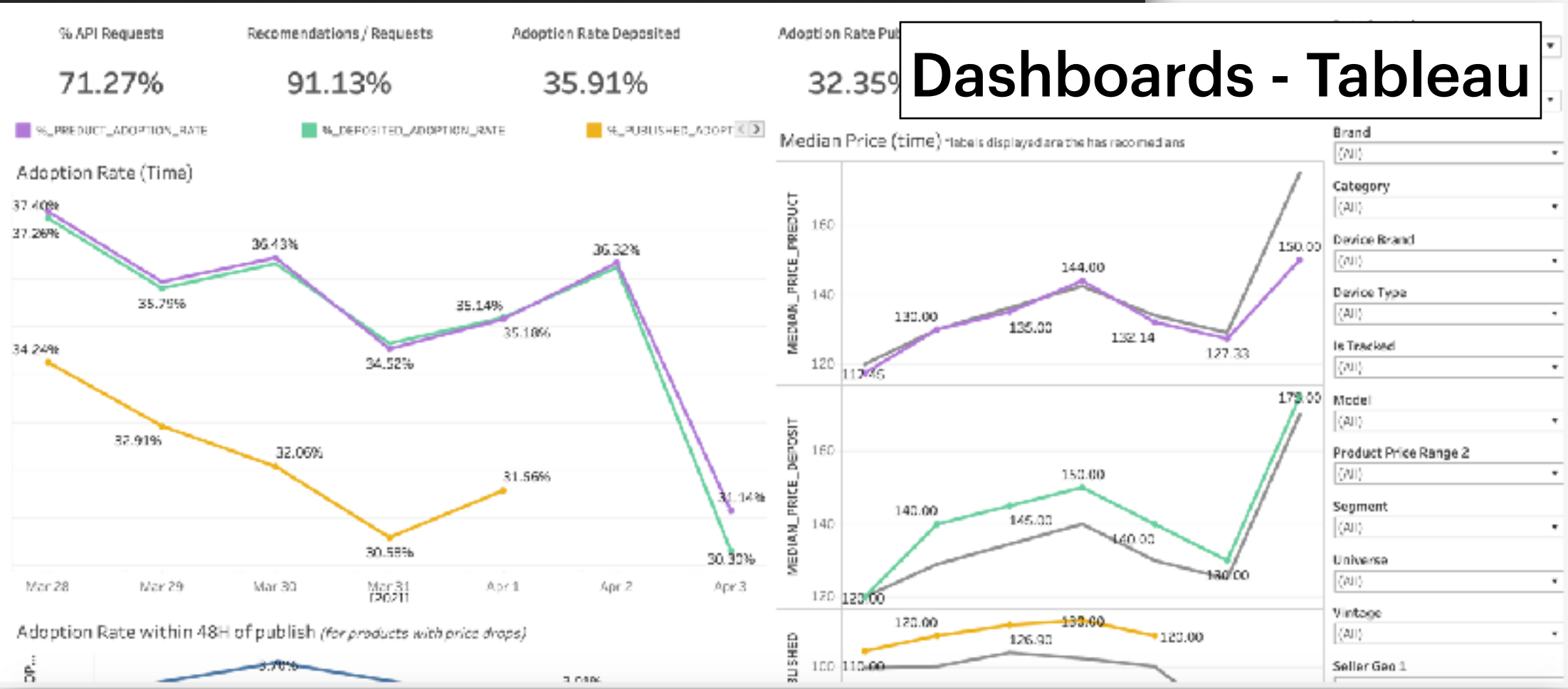
- What is the most common client path on the website?
- What was the impact of last promotion on sales?
- Do clients adopt our price recommendations ?

What do they use (day-to-day):

Postgres in Docker

SQL

```
select * from (select * from actor
JOIN film_actor fa on actor.actor_id = fa.actor_id
JOIN film f on fa.film_id = f.film_id JOIN film_ca
where fa.actor_id > 12) x
order by x.film_id
```



Tools:





Role:

- Combine various sources of data into aggregated information
- Create and provide new data based on existing data sources
- Preparing scalable dashboards for business

Questions:

- Can you add this information to the product table?
- what does this filed in this table means?
- How do we calculate the output price in this table ?
- why does my dashboard has old data ?

What do they use (day-to-day):

Postgres in Docker

SQL

```
select * from (select * from actor
JOIN film_actor fa on actor.actor_id = fa.actor_id
JOIN film f on fa.film_id = f.film_id
where fa.actor_id > 12) x
order by x.film_id
```

ETL

Extract
Transform
Load

DATA WAREHOUSE
DATA MART
HADOOP
FLAT FILES
XML

ETL - Airflow

Dashboards - Tableau

Tools:



Data Engineer



Role:

- Connect new sourced of data
- Manage DataLake (databases, accesses, roles, warehouses)
- Prepare data sources interfaces: APIs

Questions:


- Can you scrape our competitors data?
- we need this data to be available for mobile on the API
- we need access to this data table

What do they use (day-to-day):

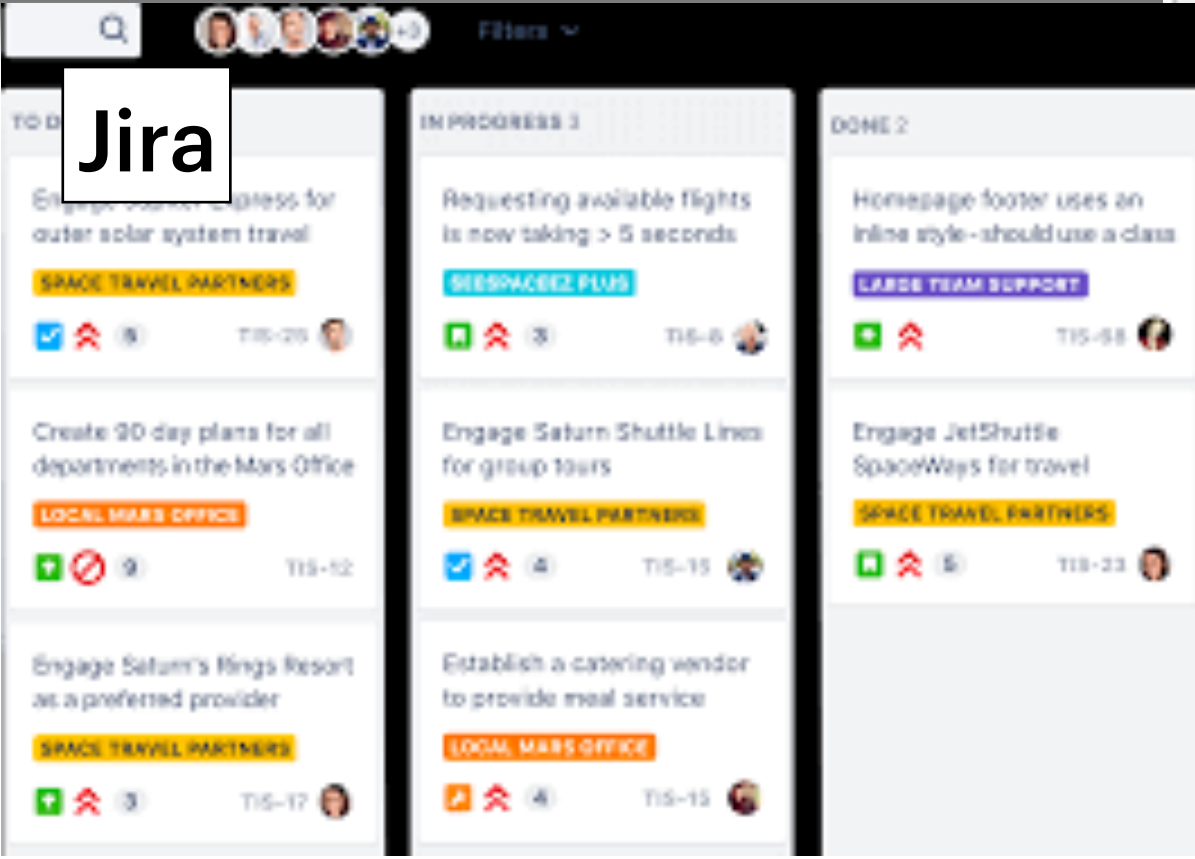
SQL

```
select * from (select * from actor
JOIN film_actor fa on actor.actor_id = fa.actor_id
JOIN film f on fa.film_id = f.film_id
where fa.actor_id > 12) x
order by x.film_id
```

ETL - Airflow



Jira



Python

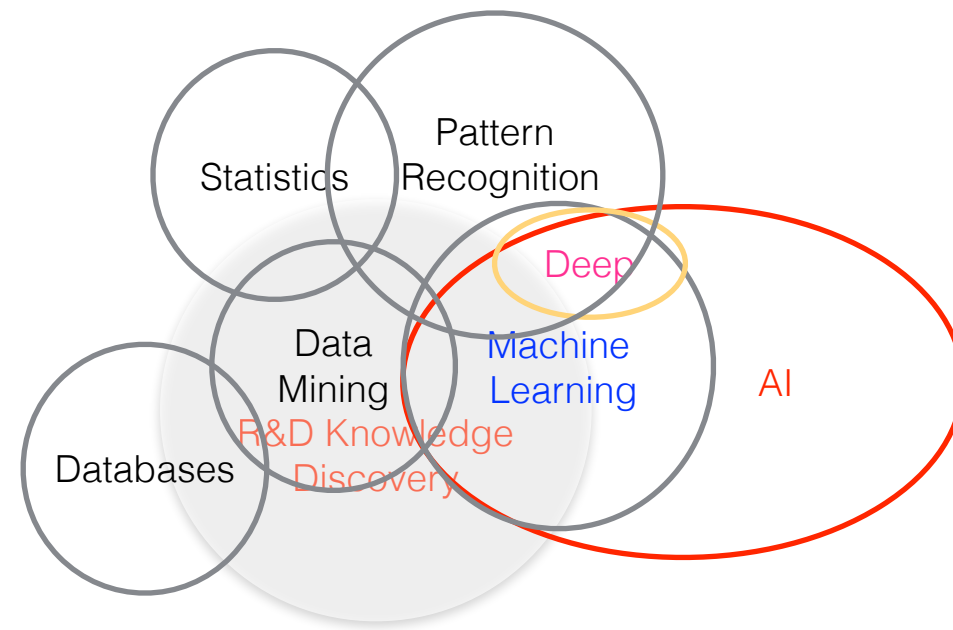
```
from annoy import AnnoyIndex
import joblib
import json
import pandas as pd
from loguru import logger

def build_annoy_index(feats_vecs, nr_trees=100, save_path="annoy_index.ann"):
    logger.info("building index")
    vector_len = feats_vecs[0].shape[0]
    logger.info(f"detected vector features: {vector_len}")
    ai = AnnoyIndex(vector_len, 'euclidean')
    logger.info("adding items")
    for i, v in enumerate(feats_vecs):
        ai.add_item(i, v)
    logger.info(f"building with: {nr_trees}")
    ai.build(nr_trees)
    logger.info(f"saving to: {save_path}")
    ai.save(save_path)
    logger.info("All Done - ✅")
    return ai
```

Tools:



Data Scientist



Role:

- Build statistical models to solve a business problem
- analyse and clean the data to select right approach for a problem
- manage a dedicated data projects: end-to-end

Questions:

can we predict orders/prices of all products per country?

can we fill-in missing values in our attributes?

can we test which feature is better for the user experience statistically?

which products should we show to specific users?

What do they use (day-to-day):

ETL - Airflow

Jira

Python

notebooks

Tools:






What do they use (day-to-day):

Role:


- deploy/scale all data science models technical architecture
- provide technical tools for scalability
- assure best practices in software development

Questions:

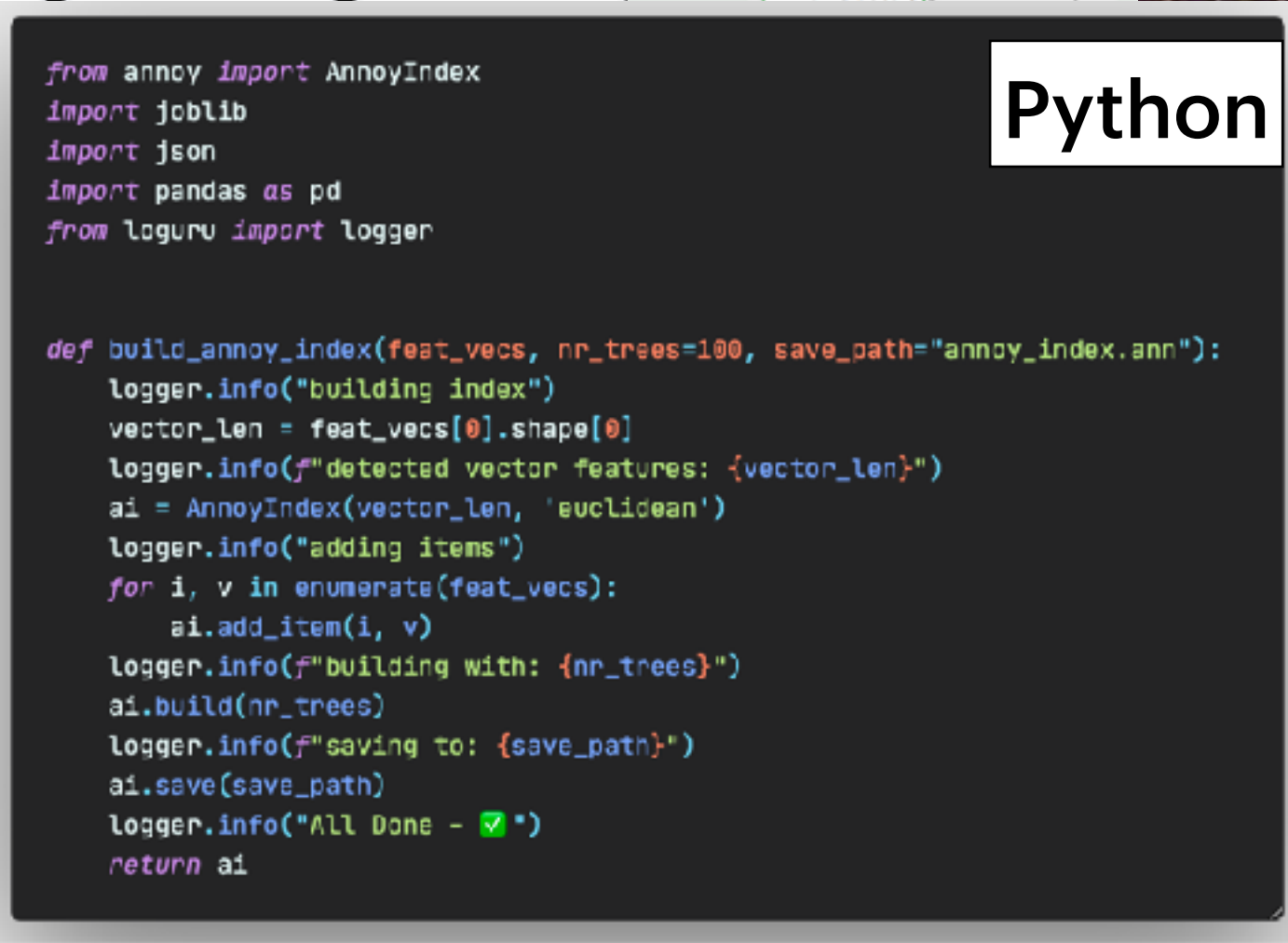
- how can I deploy this API on Kubernetes?
- can you deploy this model to production?



ETL - Airflow




Jira



Python

```
from annoy import AnnoyIndex
import joblib
import json
import pandas as pd
from loguru import logger

def build_annoy_index(feats_vecs, nr_trees=100, save_path="annoy_index.ann"):
    logger.info("building index")
    vector_len = feats_vecs[0].shape[0]
    logger.info(f"detected vector features: {vector_len}")
    ai = AnnoyIndex(vector_len, 'euclidean')
    logger.info("adding items")
    for i, v in enumerate(feats_vecs):
        ai.add_item(i, v)
    logger.info(f"building with: {nr_trees}")
    ai.build(nr_trees)
    logger.info(f"saving to: {save_path}")
    ai.save(save_path)
    logger.info("All Done - ✅")
    return ai
```



terminal

```
vivek@nixcraft:/tmp$ vi hello.sh
vivek@nixcraft:/tmp$
vivek@nixcraft:/tmp$ chmod +x hello.sh
vivek@nixcraft:/tmp$
vivek@nixcraft:/tmp$ ls -l hello.sh
-rwx 1 vivek vivek 31 Jan 21 15:08 hello.sh
vivek@nixcraft:/tmp$
vivek@nixcraft:/tmp$ ./hello.sh
Hello World
vivek@nixcraft:/tmp$
```

Tools:



Personal Recommendations

Recommendations:

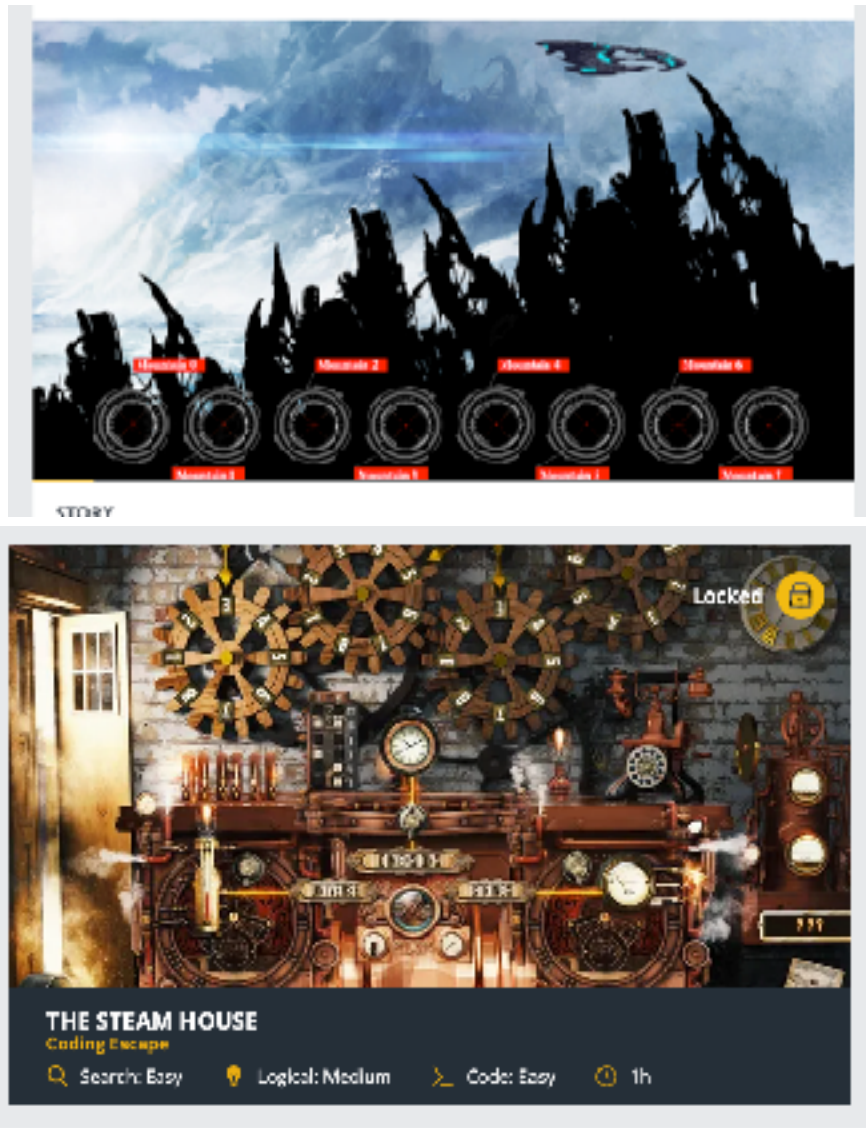
- you are the only one responsible for your path/future so **think about it NOW!**
- you do not need network nor luck -> **you need skills/knowledge** (SQL, Python, ...)!
- **plan** your project, **measure impact** and always keep the **bigger view** in mind
- propose change and **challenge** -> do it **by example**
- **find out what do you like to do** and go this path!
- start small and build on top of that -> do not over-engineer (not right away)
- **learn Git and Python** and SQL - by your own projects ! (this gives flexibility for **any carer path in tech**)
- remember **PhD is a real** experience (you learn lots of things) -> business needs to accept it!
- when doing something be able to **estimate what/when it is a success and when failure**
- please don't prepare you CV in TeX :)
- always think bigger -> do not depend on technology, network, luck
- 1 job offer = 1 customised CV
- **train** the recruitment process !

THANK YOU!

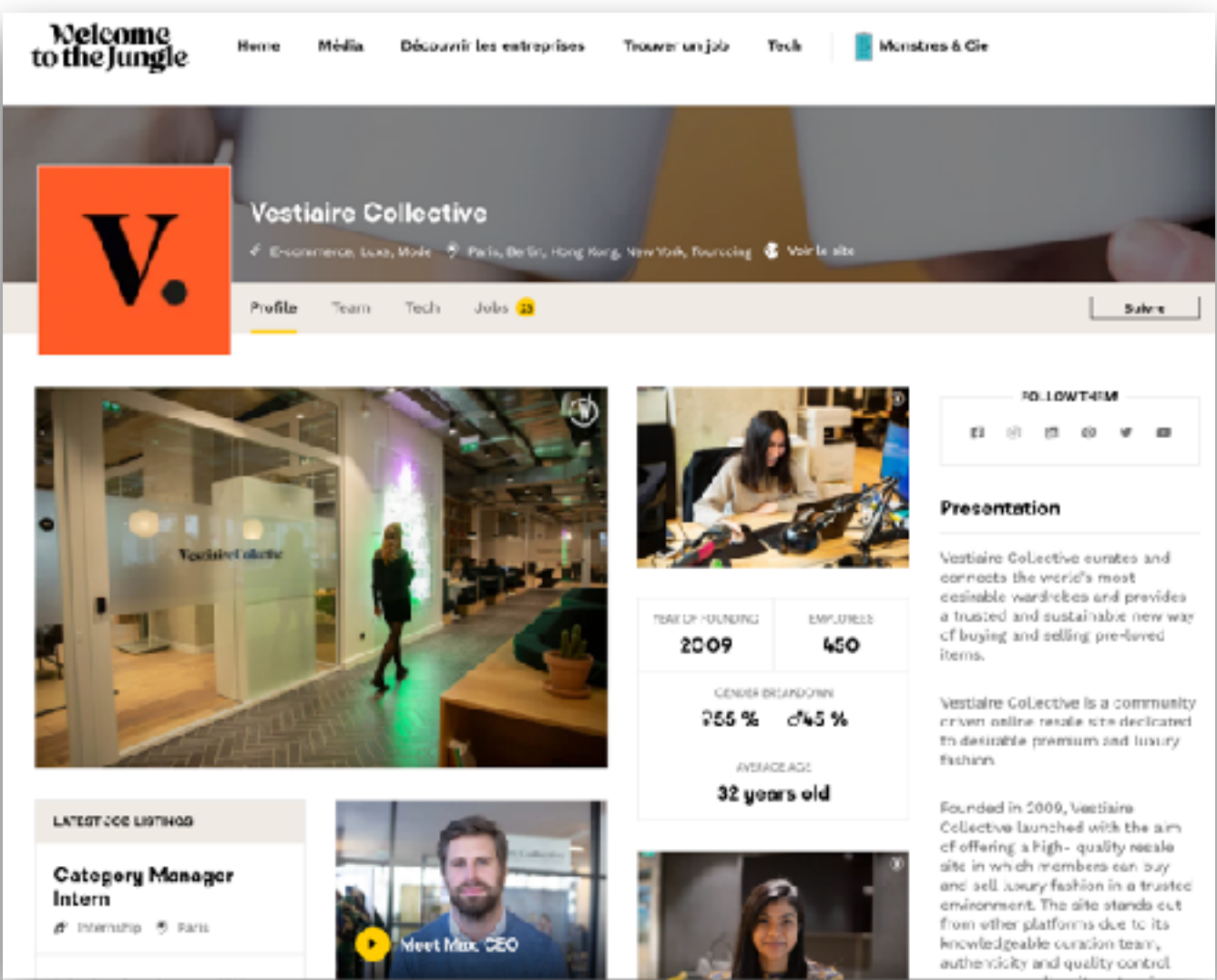
Recommendations:



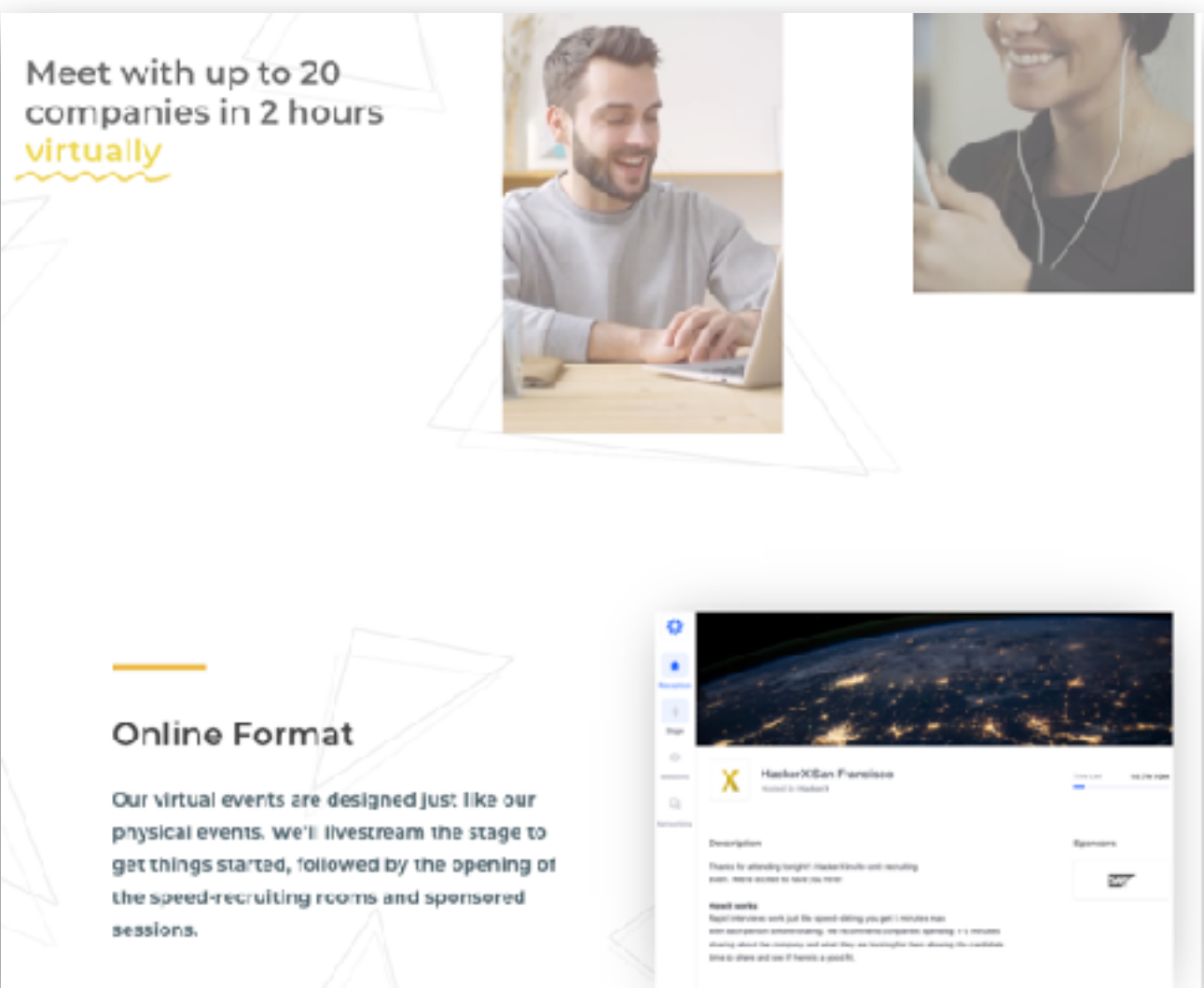
www.codinggame.com



www.welcometothejungle.com



<https://hackerx.org/>



THANK YOU!