# - Data Mining Project -

ಠ‿ಠ
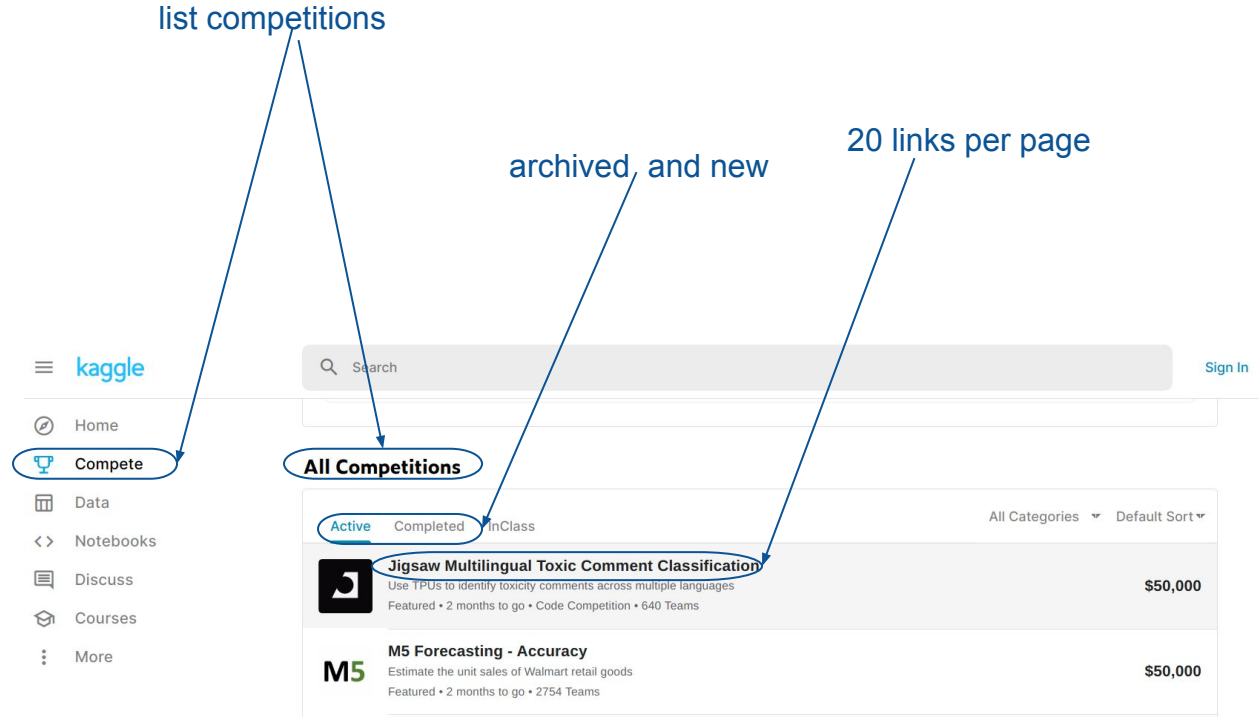
# The data

**Kaggle**, a subsidiary of Google LLC, is an online community of data scientists and machine learning practitioners.

**Kaggle** allows users to find and publish data sets, explore and build models in a web-based data-science environment, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges.

list competitions

archived, and new

20 links per page

kaggle

Search

Sign In

Home

Compete

**All Competitions**

Data

Notebooks

Active  Completed  InClass

All Categories ▾  Default Sort ▾

Discuss

Courses

More

**Jigsaw Multilingual Toxic Comment Classification**
Use TPUs to identify toxicity comments across multiple languages
Featured • 2 months to go • Code Competition • 640 Teams

$50,000

**M5 Forecasting - Accuracy**
Estimate the unit sales of Walmart retail goods
Featured • 2 months to go • 2754 Teams

$50,000

ರ.ರ

# The data

(~˘▾˘)~



second stage

topic

title

prize

header

end date

organizator

description

notebooks

score

entries

teams count

leaderboard

team name

start date

# Kaggle changed competition page view:



**Chess ratings - Elo versus the Rest of the World**

This competition aims to discover whether other approaches can predict the outcome of chess games more accurately than the workhorse Elo rating system.

$617 · 252 teams · 9 years ago

Before:

Featured Prediction Competition

**Zillow Prize: Zillow's Home Value Prediction (Zestimate)**

Can you improve the algorithm that changed the world of real estate?
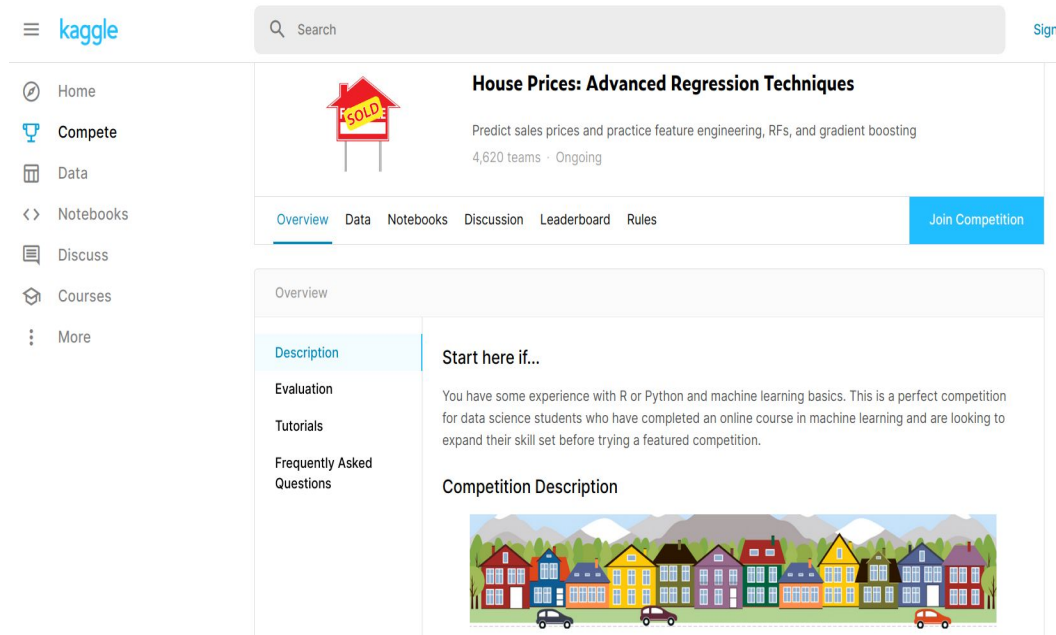
$1,200,000

Prize Money

Zillow · 3,775 teams · 2 years ago

: After

# Some numbers of how much data we scraped

- 5,000 competitions

- 15 000 scrapped pages

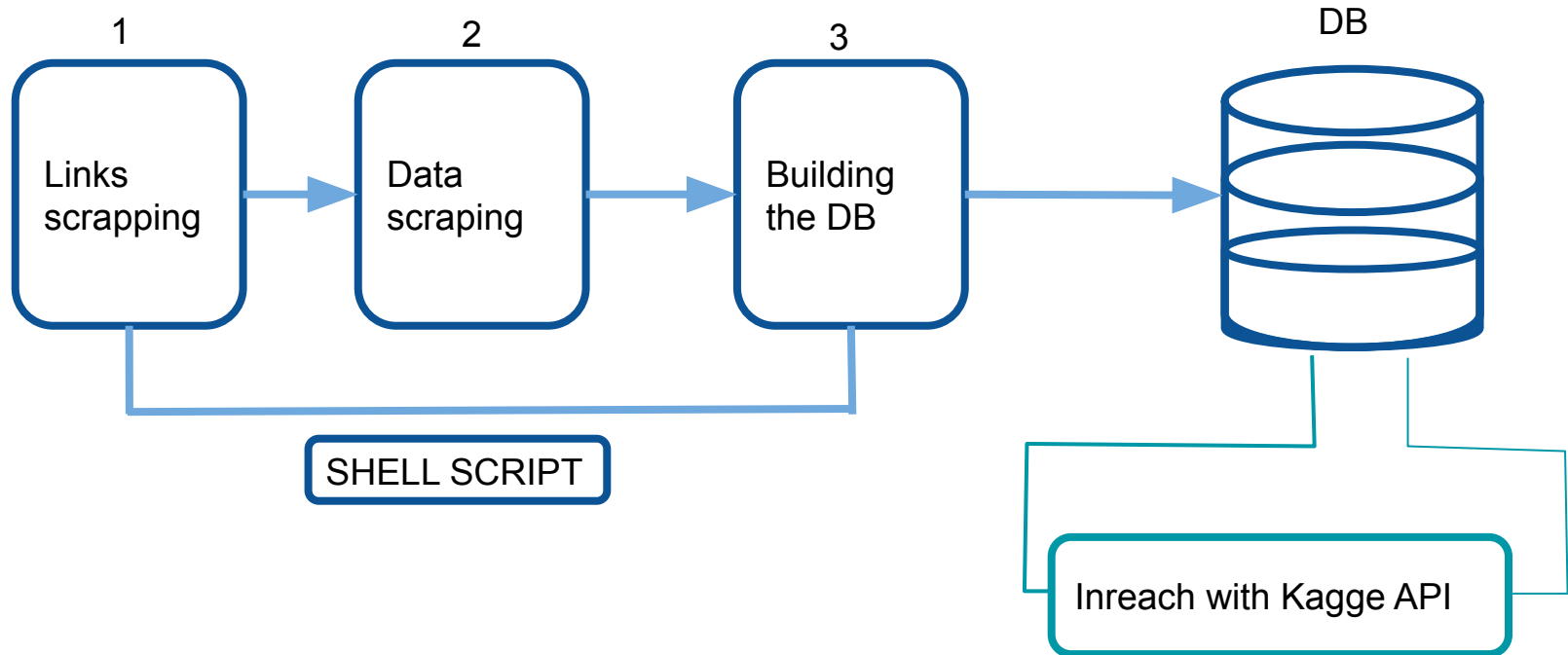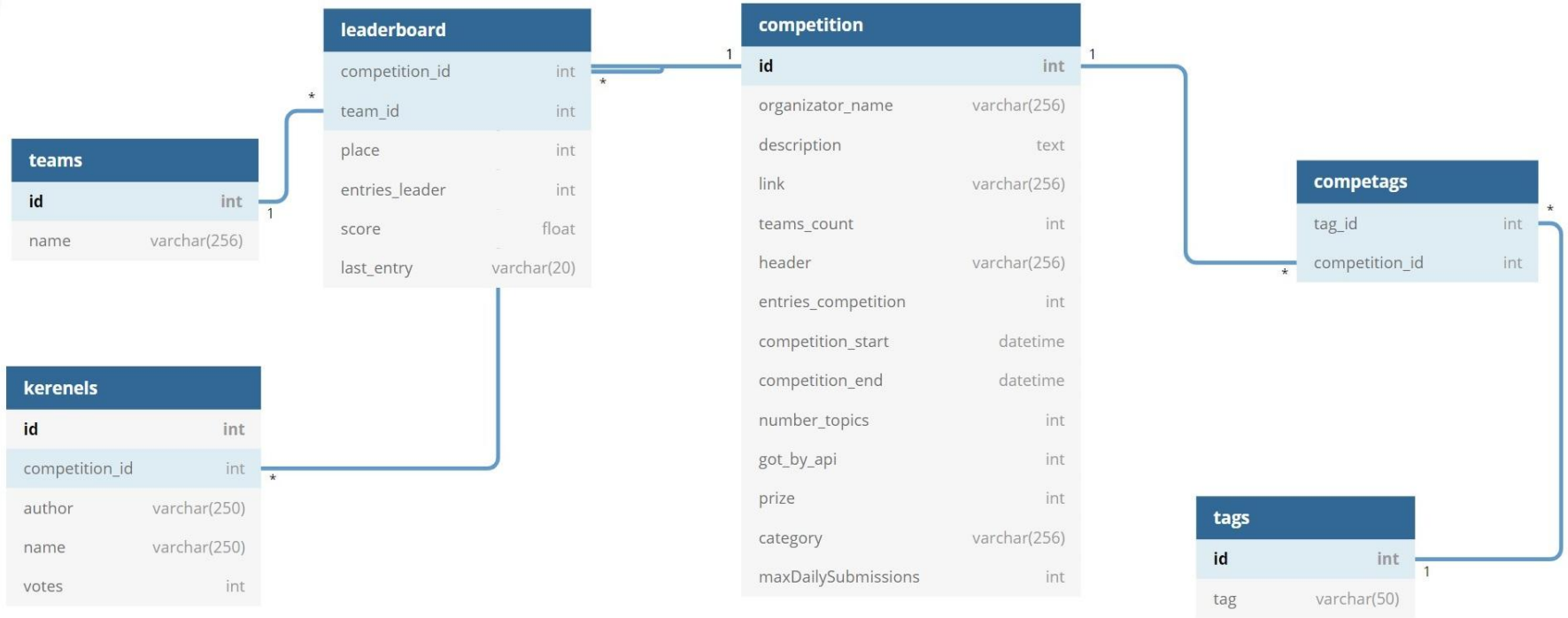- 90 000 teams

- Data from 2010 – 2020

⊙_⊙

# Project structure

(┌■_■)

# Database structure



*why we did our structure in this way?

( •_•)>⌐■-■

# Kaggle api

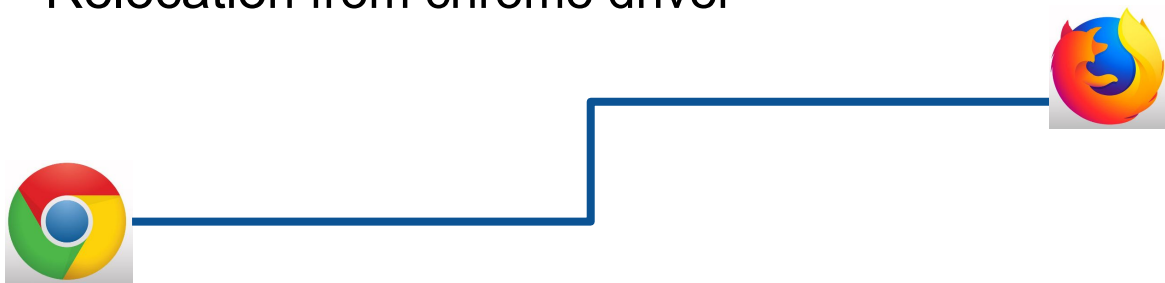|  | pro |  | contra |
|---|---|---|---|
| **+** | quick work | **-** | Scrap information only for 400 competitions of 6500 |
| **+** | fast connection | **-** | need to create account on kaggle to get a private key |

kaggle api

^و)ᵔᴥᵔ)

# Challenges we experienced and how we faced them

- At first the feeling that "I didn't understand at all what was wanted from me"　　　　ಠ_ಠ

- Kaggle is pretty old and changed the page structure several times

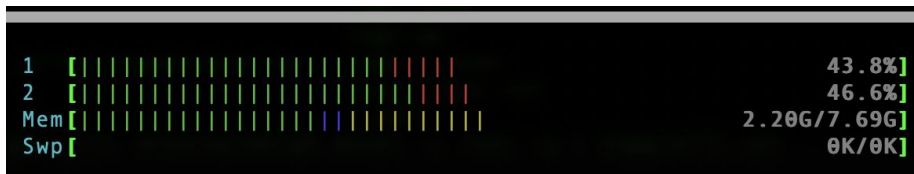- A lot of data: downloading takes a lot of time

- Relocation from chrome driver

# Challenges we experienced and how we faced them

- Troubles on Amazon server with connections ("!!!!!a server has to get rest. And what about me?!")

- Memory issue:                                  ಠ_ಥ

  the driver get all the memory
  and AWS freezes us.

  *Solution*: recreate the driver
  instance every 25 links

```
1   [||||||||||||||||||||||||            43.8%]
2   [||||||||||||||||||||||||||          46.6%]
Mem [||||||||||||||||||||          2.20G/7.69G]
Swp [                                     0K/0K]
```

Never forget about logging!!!

```
📄 DB_structure.png         204   2020-04-12 18:01:55,013 - __main__ - INFO - Extracting tags for https://www.kaggle.com/c/fia-machine-learning-t7
📄 download.log             205   2020-04-12 18:01:57,141 - __main__ - INFO - Collected data for link https://www.kaggle.com/c/fia-machine-learning-t7
📄 download_one.py          206   2020-04-12 18:01:57,141 - __main__ - INFO - Collected data for competitions.
📄 get_links.py             207   2020-04-12 18:01:57,142 - __main__ - INFO - the dictionary/json file saved into csv competition.csv
📄 geter_num_topics_prize_organiz 208   2020-04-12 18:01:57,142 - leaderboard - INFO - Extracting leader board data...
📄 insert_to_db.py          209   2020-04-12 18:02:13,307 - leaderboard - INFO - Extracted leader board data for link https://www.kaggle.com/c/passenger
📄 kaggle_api.py
```

# A short demo from redash.io

°₊✧ ੭(⁰▽⁰)੭ ✧₊°

Not a good platform for making visualisations!!!!

## Competitions by teams count grouped by category



## Teams count grouped by category



©http://ec2-54-183-241-235.us-west-1.compute.amazonaws.com/queries/3#6

# Number of topics in discussion

Without prize

With prize

People discuss more the competitions with prizes



6.83%

93.2%

19.6%

80.4%

☆(❀‿❀)☆

# Data Insights

- Very little of Kaggle competitions have prizes

- Competitions popularity doesn't depend on size of a prize

- Most of the competitions are some university assignments

- There are no teams that won more than 1 big prize.

Good luck on Kaggle!

f(ಠ‿↼)z