

Machine learning *in the wild*

Tales from machine learning after college

2nd Edition

DIS 11/06/2023

Duarte O.Carmo

duarteocarmo.com - @duarteocarmo

Hello! I'm Duarte.

/du-art/ - it's Portuguese

ML/Software Engineer & contractor

From Portugal, based in Copenhagen, Denmark

I like **running**, and **writing** on my blog

Past: Strategy, Product Mgmt., New Ventures, Mgmt. Consulting

Now: I help companies solve **tough** problems end-to-end



Wequity

amplemarket

TALKATIVE



DIS

Today, we'll talk about machine learning from what I've seen out there

- How (I think) ML engineers should work
 - 3 example problems from the wild
 - “MLOps”
 - Learning
-
- *Opinions*
 - *Experiences*

MAGAZINE SPRING 2021 ISSUE / RESEARCH FEATURE

Why So Many Data Science Projects Fail to Deliver

Organizations can gain more business value from advanced analytics by recognizing and overcoming five common obstacles.

Mayur P. Joshi, Ning Su, Robert D. Austin, and Anand K. Sundaram • March 02, 2021

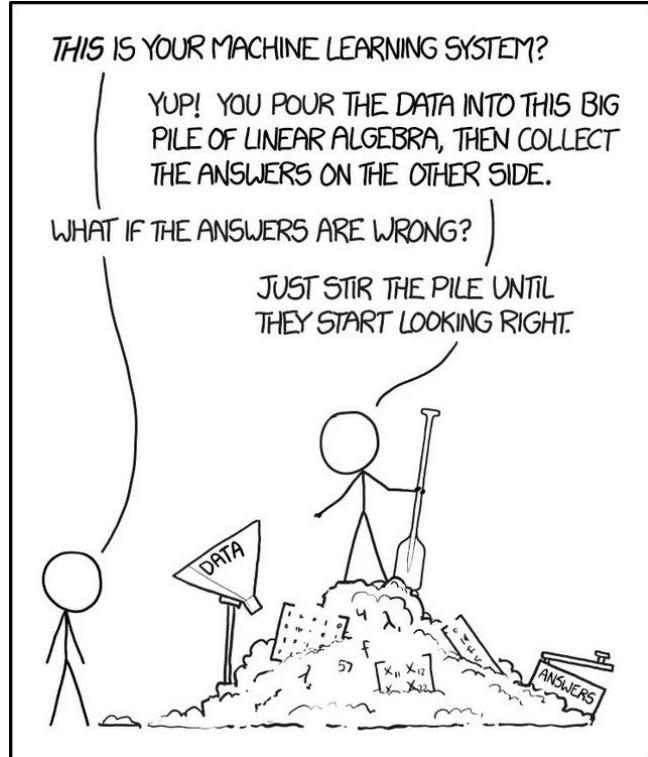
Reading Time: 14 min

1 | How I work

“We need a model”

(you probably don't)

Don't start with models, start with people



- Define the business goal, and the success metric
- This is real world (bad) data - not Kaggle: cr*p in, cr*p out
- Start with heuristics, and increase complexity as needed
- Put it out there as fast as possible, then iterate

Go end-to-end, early!



Don't build in the basement

You are makers at heart – and should treat your schedules like it



- Minimize time in meetings and double down on communication
- Fridays = no meetings
- We are on an emerging tech field, studying is important
- We are builders of things, disruptions are not welcome

Tools are irrelevant, until you have to use them every day



```
train_segmentation.yaml
...
kriging.py
...
EXPLORER
```

A screenshot of a VSCode workspace. On the left, there are two code editors: one showing `train_segmentation.yaml` and another showing `kriging.py`. The `train_segmentation.yaml` file contains configuration for a dataset loader. The `kriging.py` file is a Python script that performs kriging on a field cube. On the right, the VSCode Explorer sidebar shows the project structure for a DeepOCModel, including folders for DEEPOCMODEL, configs, datasets, models, notebooks, pipelines, steps, and tests, along with various configuration files like `objective_col.yaml` and `train_segmentation.yaml`.

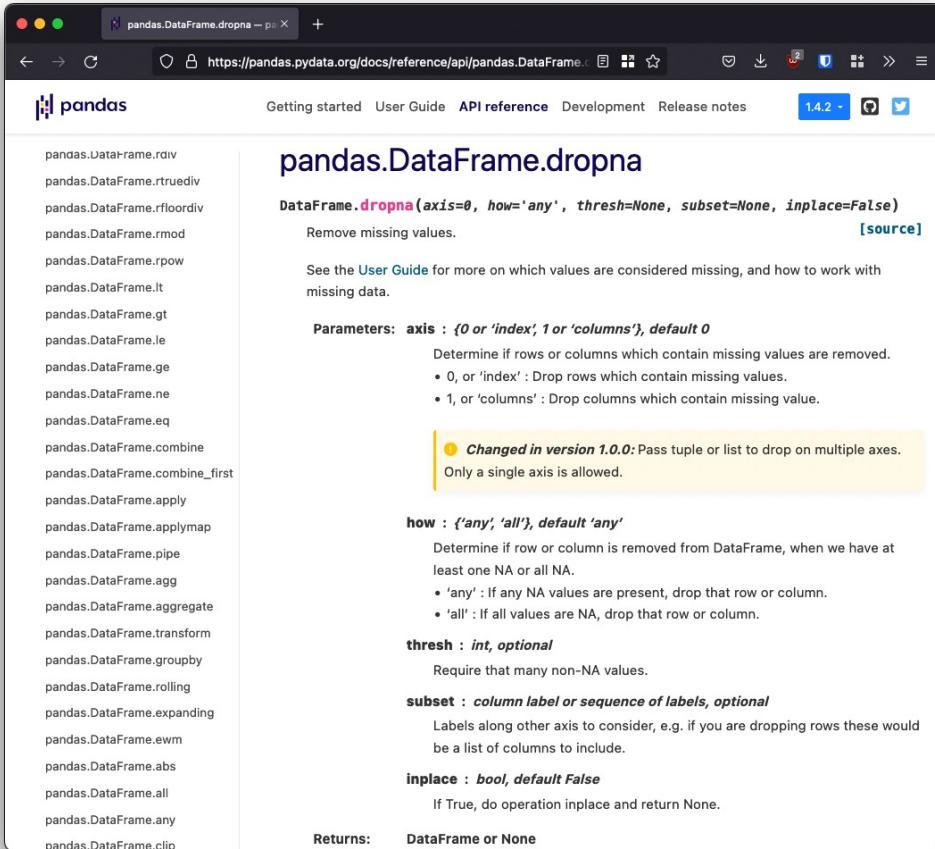
VSCode - (I actually use NeoVim)



A screenshot of a GitHub repository page for `DeepOCModel`. The page shows the repository structure, including branches, tags, and commits. A prominent message says "Your main branch isn't protected". Below this, the commit history lists several changes made by `duarteocarmo`, such as adding workers to config, closing a wandb loop, and updating the train test dataset. On the right side, there are sections for "Releases" (no releases published) and "Languages" (Jupyter Notebook 88.5%, Python 1.4%, Makefile 0.1%). At the bottom, there are "Suggested Workflows" for Actions Importer, SLSA Generic generator, and Python application.

GitHub, some people prefer GitLab

But there's quite nothing like reading

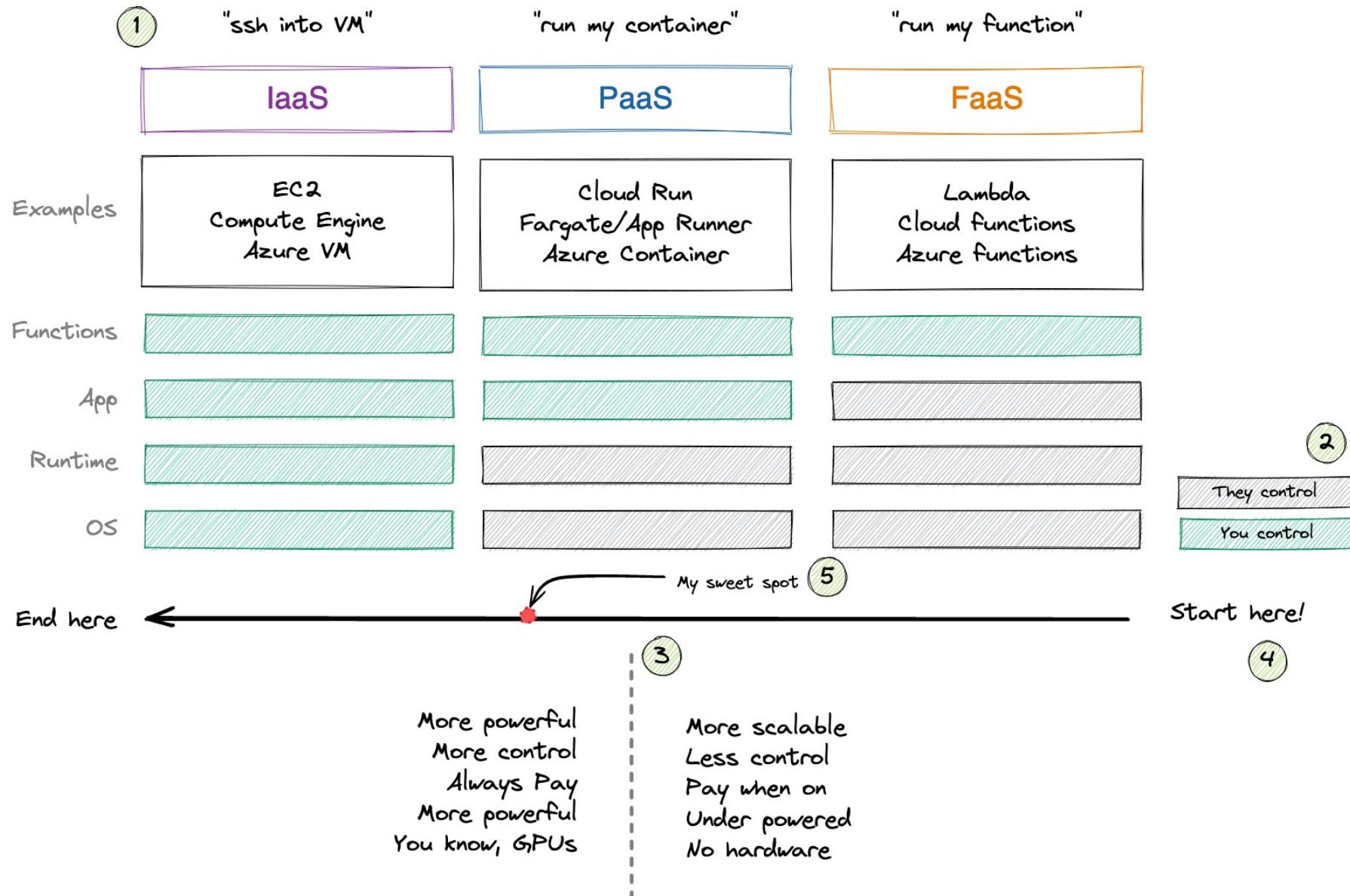


The screenshot shows a browser window displaying the pandas API documentation for the `DataFrame.dropna` method. The page is titled "pandas.DataFrame.dropna". It includes a sidebar with a list of other pandas DataFrame methods such as `raiv`, `rtruediv`, `rfloordiv`, `rmod`, `rpow`, `it`, `gt`, `le`, `ge`, `ne`, `eq`, `combine`, `combine_first`, `apply`, `applymap`, `pipe`, `agg`, `aggregate`, `transform`, `groupby`, `rolling`, `expanding`, `ewm`, `abs`, `all`, `any`, and `clip`. The main content area describes the `dropna` method, which removes missing values from a DataFrame. It details parameters for axis (0 or index, 1 or columns), how (any or all), thresh (number of non-NA values required), subset (columns to consider for dropping), and inplace (whether to do the operation in place). A note indicates that the method was changed in version 1.0.0 to accept tuples or lists for multiple axes, but only a single axis is allowed. The "Returns" section specifies that it returns a DataFrame or None.

- What does it do?
- Options?
- Default behaviors
- Maybe I can re-use this
- **It *actually* sticks**

Use AI to code

(But not at the expense of your brain)



2 | Problems

2.1 | Machine learning in Space



**Geospatial
data:
All data
generated from
observing our
planet in space**

Images: RGB bands from the Sentinel-2

Temperature: From multi/hyper spectral satellites

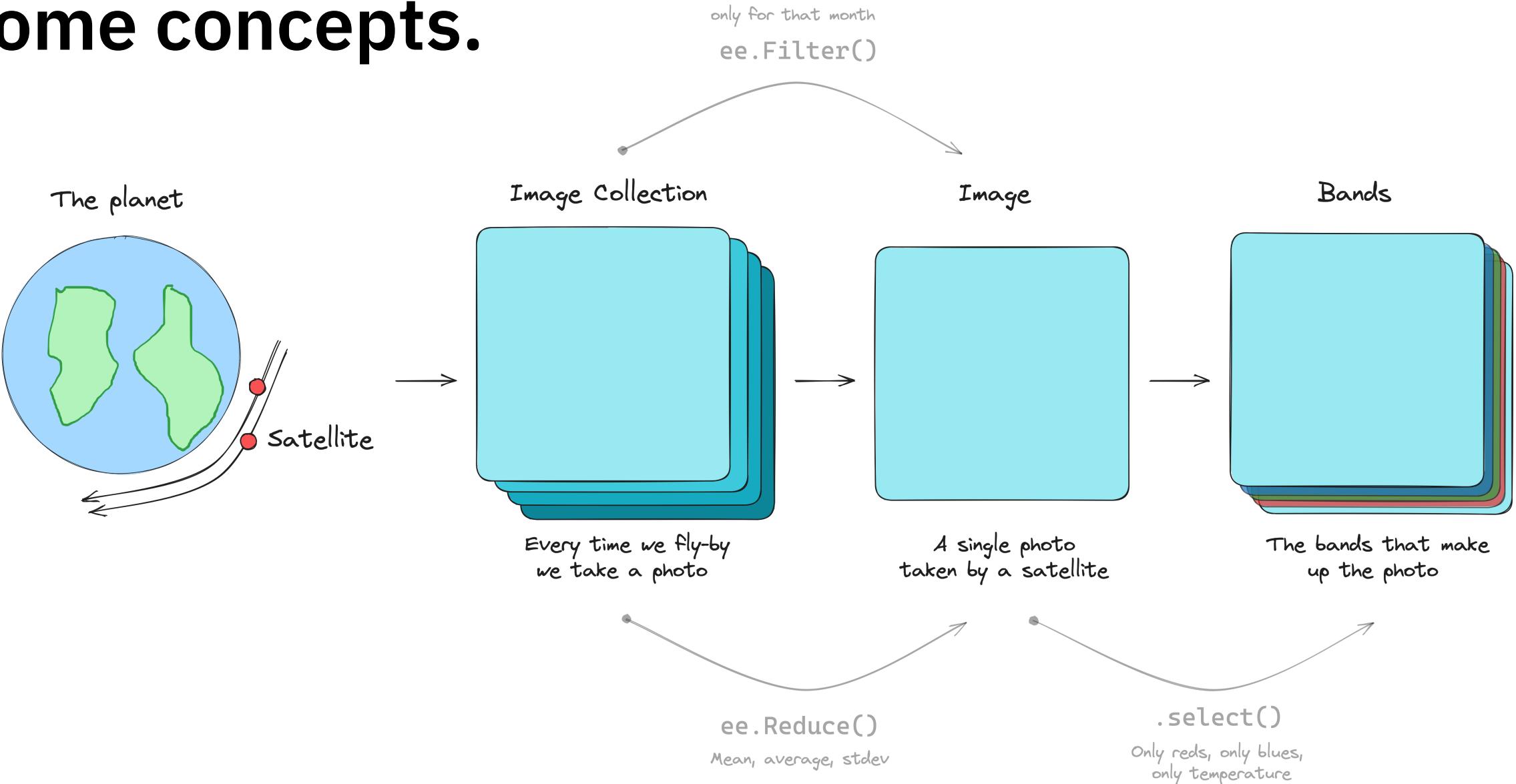
Vegetation: NDVI (vegetation index)

Landcover: Is this a city? Farm?

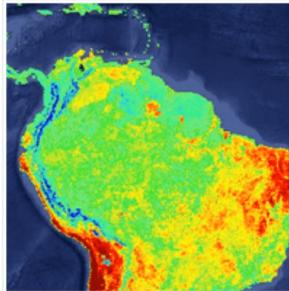
Precipitation: How much rainfall?

...

Some concepts.



MOD11A2.061 Terra Land Surface Temperature and Emissivity 8-Day Global 1km



Dataset Availability

2000-02-18T00:00:00Z–2023-07-04T00:00:00Z

For when is this data available?

Dataset Provider

NASA LP DAAC at the USGS EROS Center

Image Collection or Image?

Earth Engine Snippet

```
ee.ImageCollection("MODIS/061/MOD11A2")
```

Tags

8-day emissivity global lst mod11a2 modis nasa surface-temperature
terra usgs

Description of all bands

Description	Bands	Terms of Use	Citations	DOIs	
Resolution 1000 meters					What is the resolution?
Bands					
	Name	Units	Min	Max	Scale
	LST_Day_1km	K	7500	65535	0.02
	QC_Day				Offset
					Description
					Day land surface temperature
					Daytime LST quality indicators
	+ Bitmask for QC_Day				
	Day_view_time	h	0	240	0.1
	Day_view_angl	deg	0	130	-65
	LST_Night_1km	K	7500	65635	0.02
	QC_Night				Night land surface temperature
					Nighttime LST quality indicators
	+ Bitmask for QC_Night				
	Night_view_time	h	0	240	0.1
					Local time of night observation

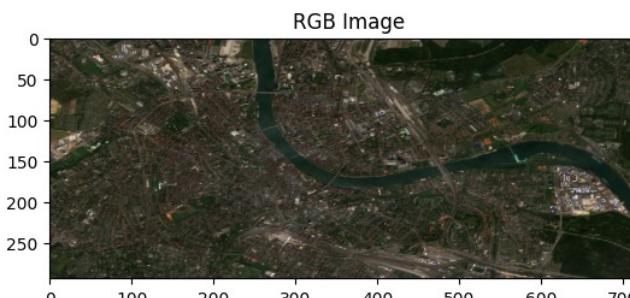
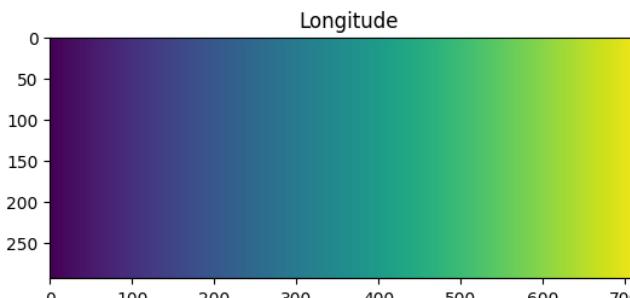
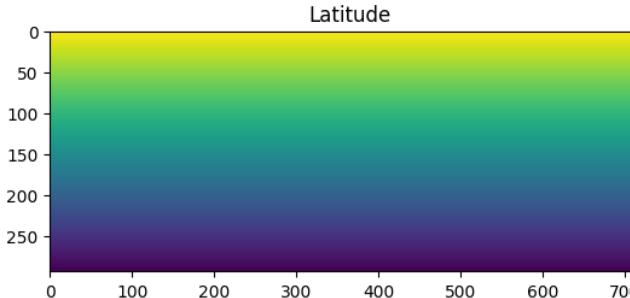
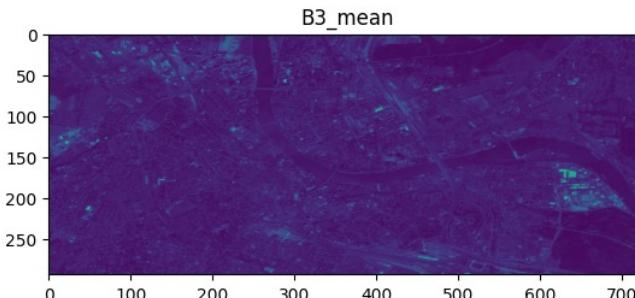
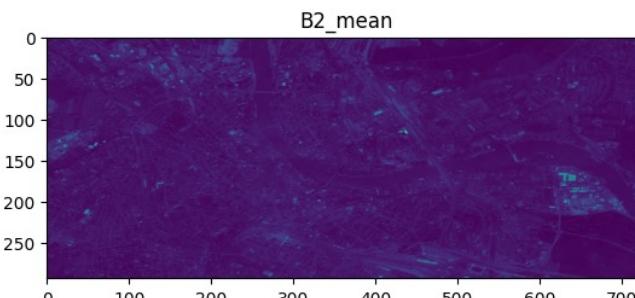
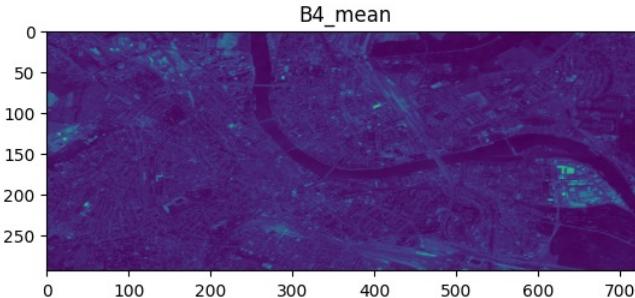
Tabular/Timeseries | Using Google Earth Engine to extract rainfall data as a time series

	country	long	lat	maincause	displaced	date	0_30d_t0_rainfall_m	1_30d_t1_rainfall_m	2_30d_t2_rainfall_m	3_30d_t3_rainfall_m
4967	India	79.226600	17.344400	Monsoonal Rain	2100	2020-10-16 00:00:00	0.619084	1.037469	0.242630	0.302843
4964	Vietnam	107.833000	15.854800	Heavy Rain	900000	2020-10-06 00:00:00	NaN	NaN	NaN	NaN
4963	Mozambique	38.127600	-14.581300	Heavy Rain	4000	2020-10-02 00:00:00	NaN	NaN	NaN	NaN
4966	Togo	0.777334	9.710160	Heavy Rain	16000	2020-09-15 00:00:00	0.099871	0.075639	0.060188	0.513882
4965	Nigeria	6.258020	7.774120	Heavy Rain	25000	2020-09-15 00:00:00	4.037949	1.330742	4.710331	3.221587
4962	USA	-86.579600	32.650800	Tropical Storm Sally	0	2020-09-15 00:00:00	0.017025	0.019261	0.060149	0.044176
4961	India	95.037700	27.848600	Monsoonal Rain	0	2020-09-13 00:00:00	NaN	NaN	NaN	NaN
4950	Afghanistan	67.720900	34.916000	Torrential Rain	0	2020-08-25 00:00:00	0.843814	0.881901	0.289106	0.776603
4951	Pakistan	68.215700	26.957300	Monsoonal Rain	1300	2020-08-24 00:00:00	0.730026	2.174792	0.494846	0.542660
4960	USA	-70.941700	19.564200	Tropical Storm Laura	600000	2020-08-22 00:00:00	0.391254	0.358152	0.596996	0.543690
4954	Haiti	-70.904600	18.907900	Tropical Storm Laura	0	2020-08-21 00:00:00	0.493560	0.363616	0.596426	0.601292
4959	India	77.469800	24.709500	Monsoonal Rain	60	2020-08-21 00:00:00	0.195524	0.190712	0.053476	0.149555
4953	Kenya	37.202100	2.750120	Dam Release and Heavy Rain	5000	2020-08-20 00:00:00	NaN	NaN	NaN	NaN
4955	Uganda	30.333400	0.385265	Torrential Rain	0	2020-08-19 00:00:00	0.006394	0.006034	0.010361	0.001781
4952	Chad	14.646100	10.199000	Heavy Rain	38000	2020-08-10 00:00:00	0.609008	0.030559	0.008041	0.000344

Rainfall (30d periods from observation date)

DIS

Computer Vision | From coordinates to images for Computer Vision applications



- ~275x750 images
- 1 pixel = 15 meters
- Lat and Long images
- RGB reconstructions

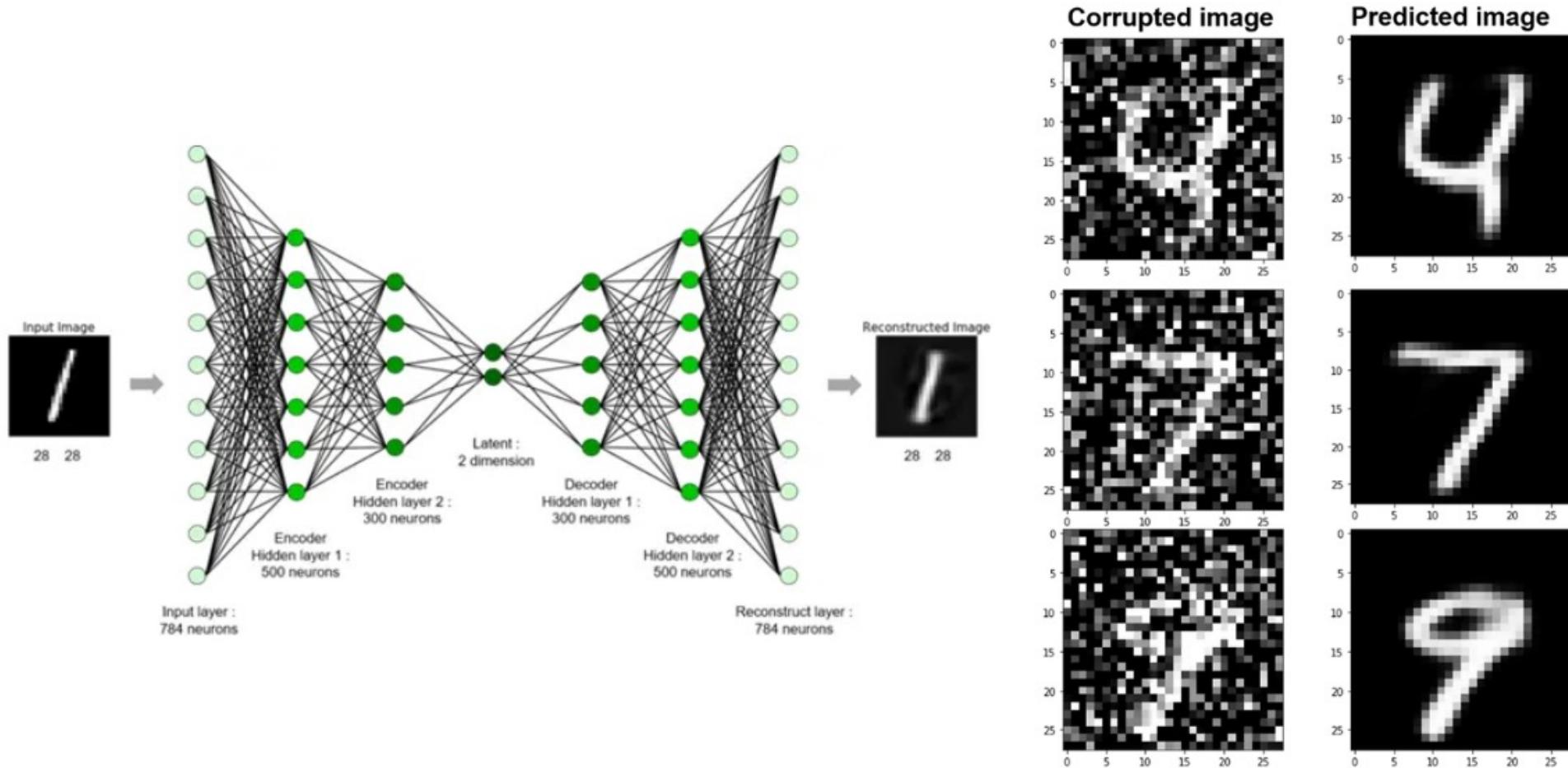
2.2 | Recommendations with embeddings

Helping sales teams find their ideal customers

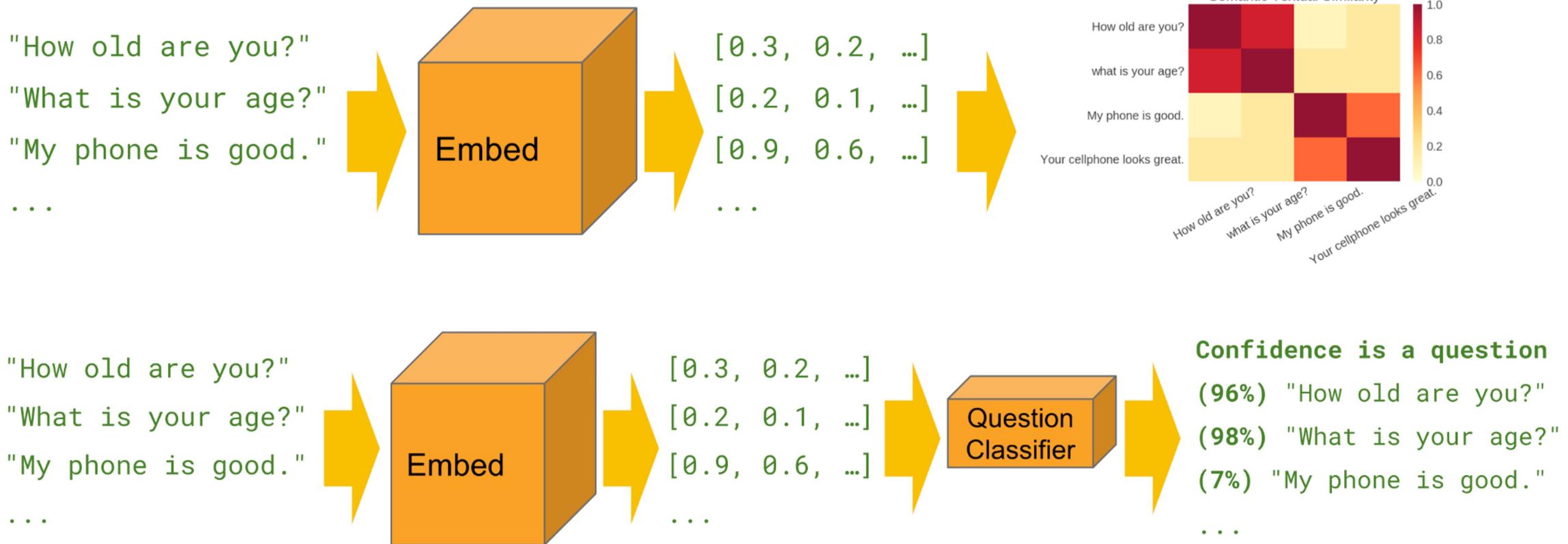
- Lead qualification is manual
- Lots of time spent qualifying
- How can we support this process?

Company Name	Description	Potential Customer?
Novo Nordisk	The Novo nordisk foun..	✓
Facebook	A social media..	✗
Budweiser	We are a bever..	✓
Nike	World leader in..	✓
Google	At Google, we're..	✗
...

First, a quick introduction to embeddings

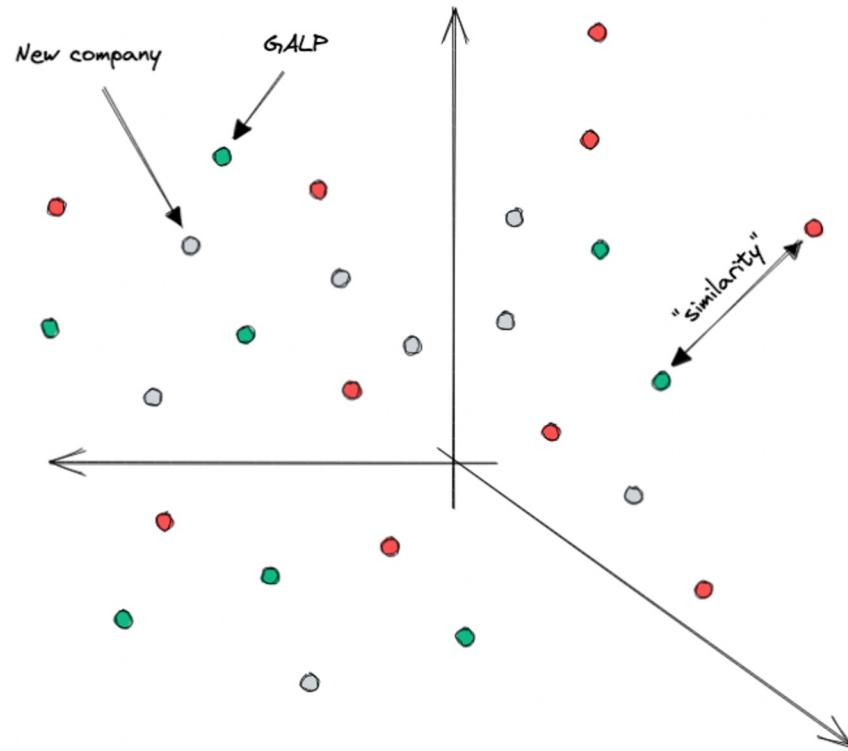


There are a lot of ways to use embeddings in real-world ML problems



1

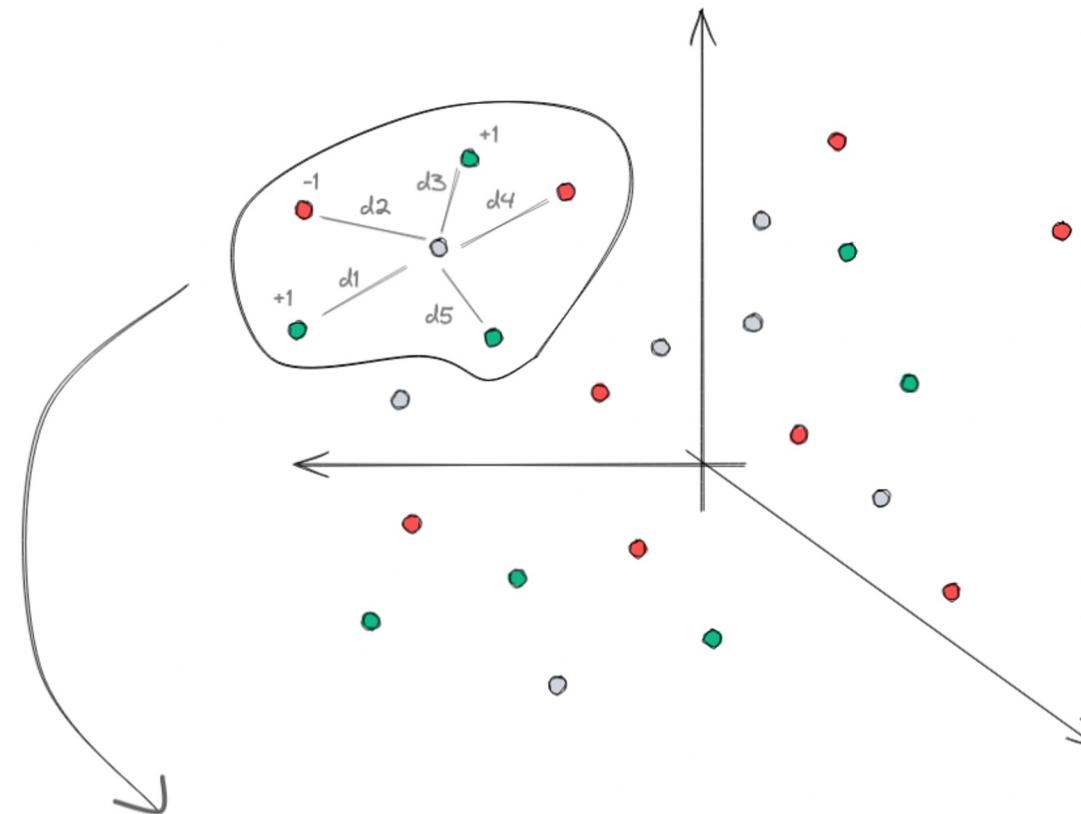
Companies in Space



- Liked
- Disliked
- Unknown

2

Scoring a new company



```
score = np.average(neighbour_score, weights=distances)
```

Bonus, find what is wrong on this formula..

DIS

KNN can be more *contextual* than a traditional binary classifier

- Embeddings can be re-used - if well chosen
- Recommendations are interpretable
- We can weigh different factors
- Multilingual with unbalanced data
- Generally easier to deploy

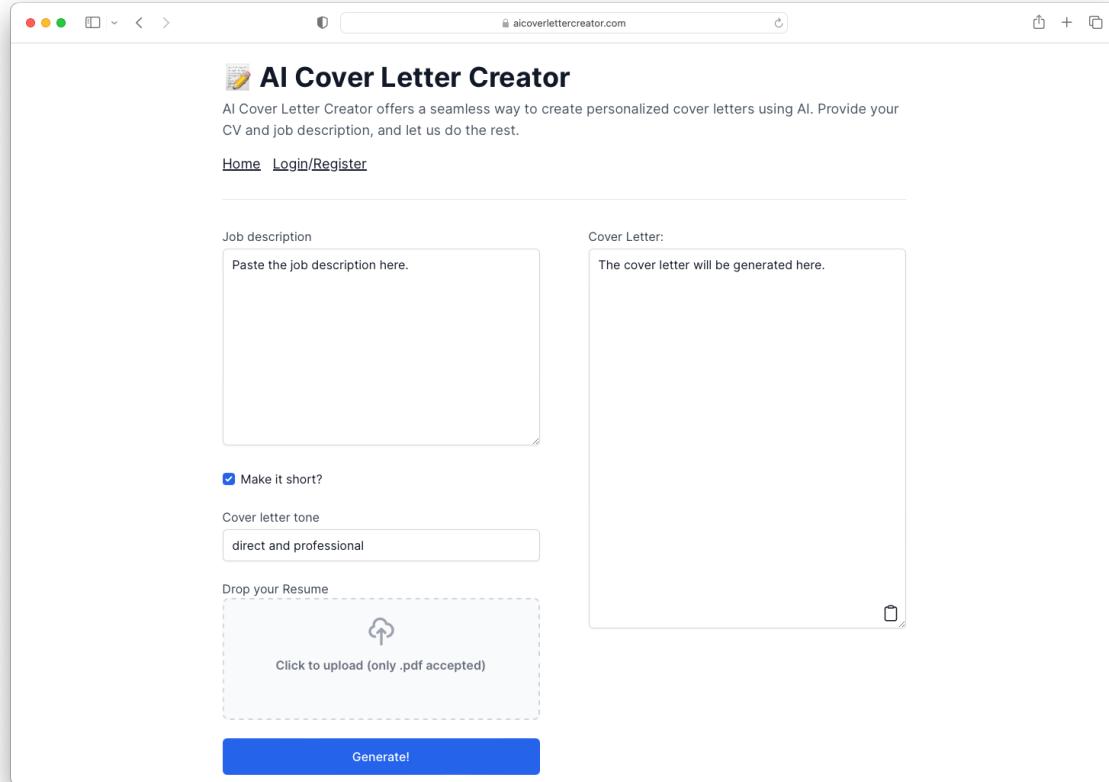
Company Name	Most similar	Rec Score
Budweiser	[... N ...]	0.879
Olx	[... N ...]	0.789
Unbabel	[... N ...]	0.678
Novo Nordisk	[... N ...]	0.001
LA Firefighters	[... N ...]	-0.995
...

Most similar neighbors

Company score

2.3 | Using generative AI

It's becoming simpler to build cool apps fast – but generative is not for everything!



aicoverlettercreator.com

- Django for the web interface
- Some Javascript for interactions
- OpenAI models
- Some cool prompts

Code Blame 22 lines (18 loc) · 839 Bytes

Raw

```
1 Act like an assistant helping the user write personalized cover letter.
2
3 The cover letter should respect the following instructions:
4
5 Tone and style:
6 - Use a {{ tone }} tone.
7 - Do not use cliche sentences like "to express my interest", "I am hereby expressing".
8 - Be concise and direct.
9
10 Format:
11 - Do not use template tags.
12 - The cover letter should be easy to read.
13 {% if short == True %}- Use 2 paragraphs this is extremely important.
14 - Use 150 words or less for the cover letter.
15 {% else %}- Use 3 paragraphs this is extremely important.
16 - Use 250 words or less for the cover letter.
17 {% endif %}
18
19 Content:
20 - Be honest and truthful, do not make things up.
21 - Start with a strong, original, and punchy opening statement to catch the attention of recruiters.
22 - Use information about the user to tailor the cover letter for the job poster.
```

The system prompt

Code Blame

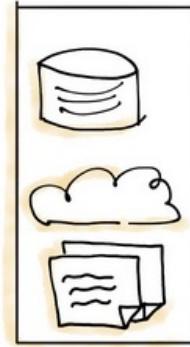
Raw

```
1 Information about the user:
2 ----
3 {{ resume_text }}
4 ----
5
6 Information about the job description:
7 ----
8 {{ job_description }}
9 ----
10
11 Personalized cover letter:
12 ---
```

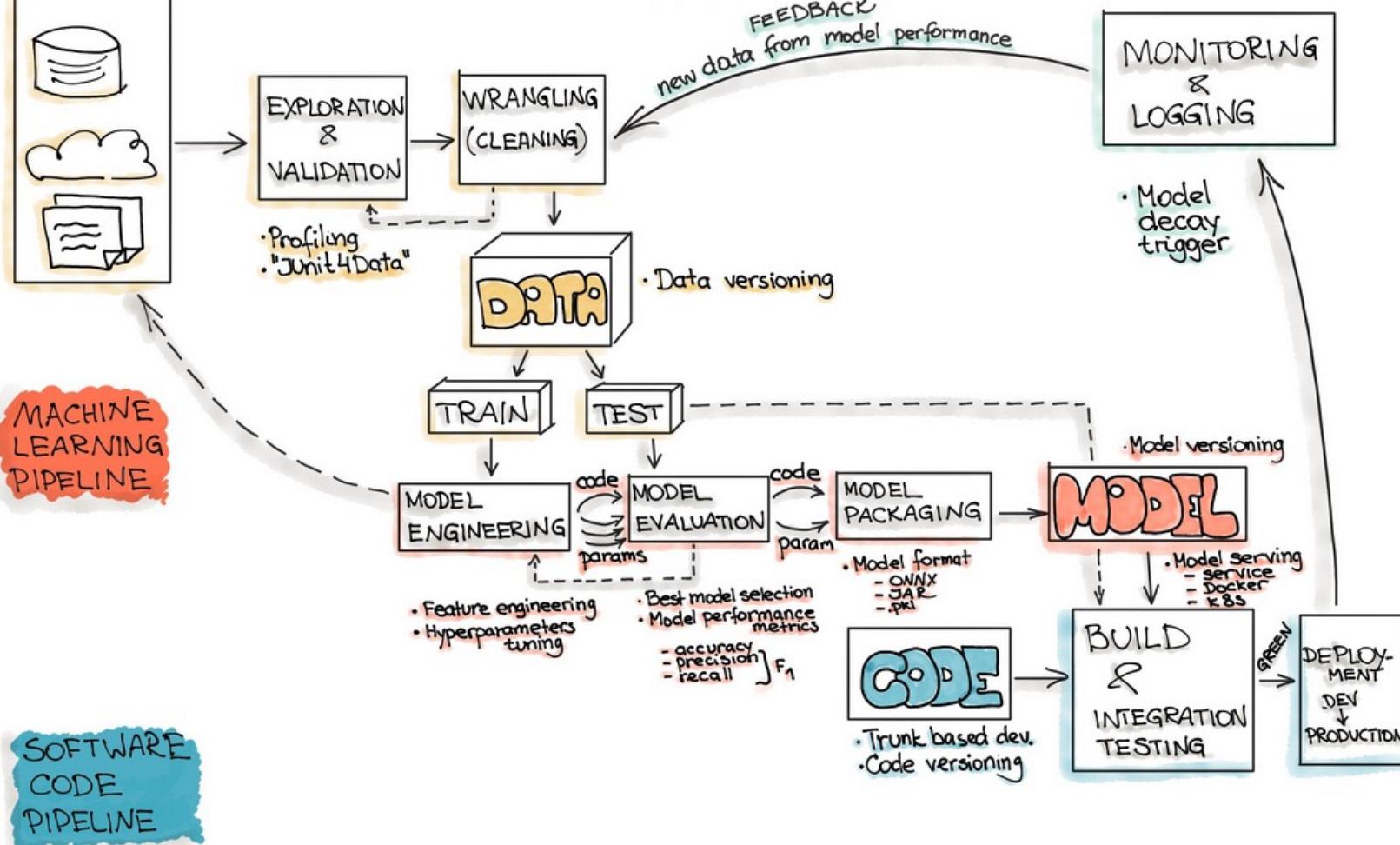
The user prompt

3 | “MLOps”

DATA PIPELINE



MACHINE LEARNING ENGINEERING.



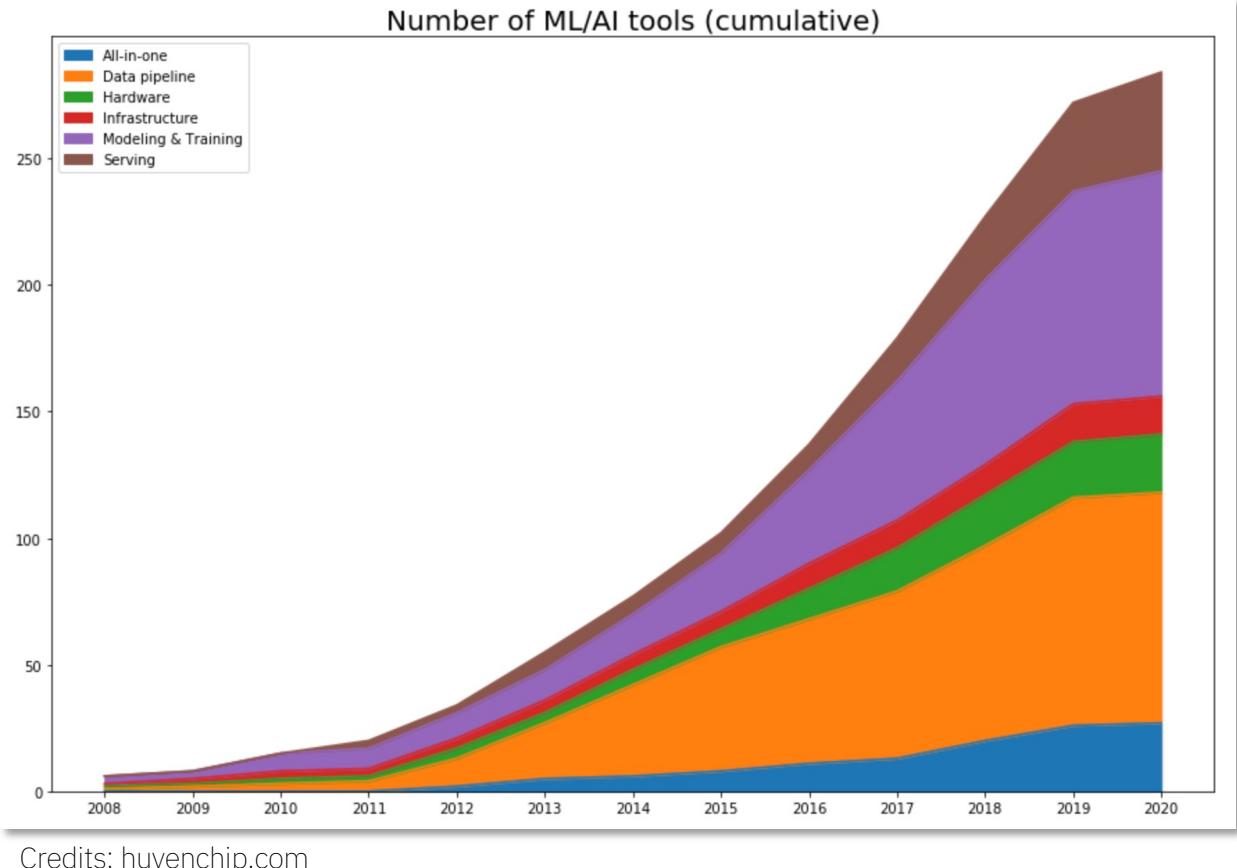
Credits: ml-ops.org

duarteocarmo.com - @duarteocarmo

DIS

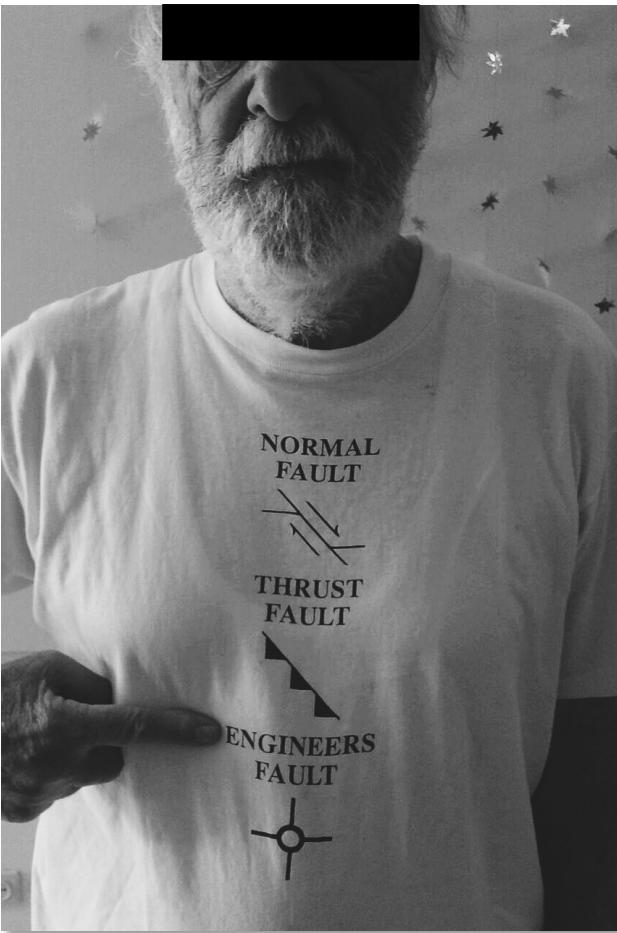
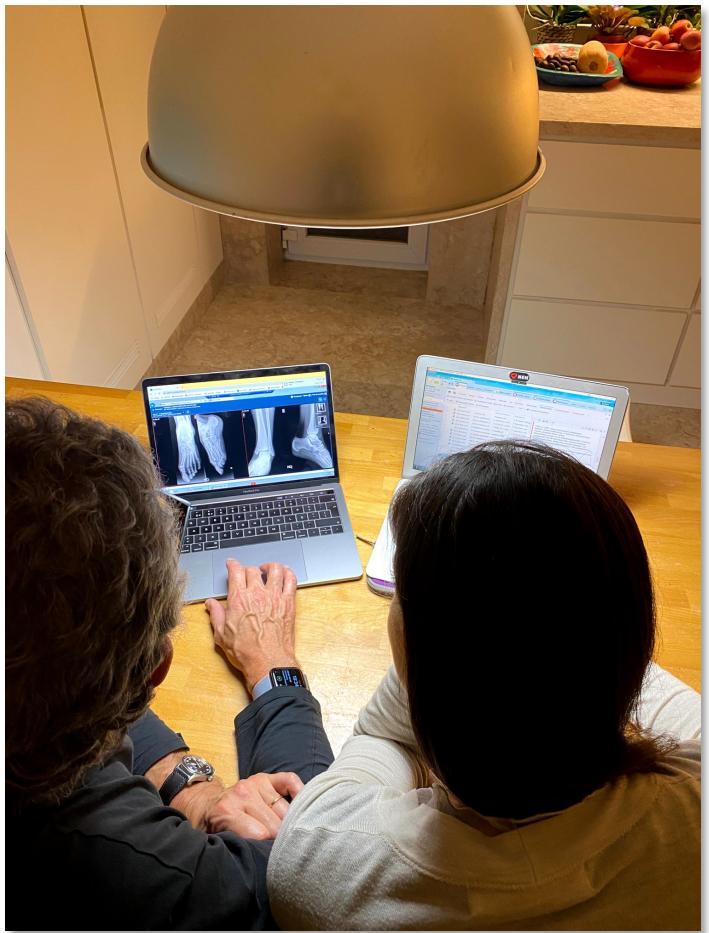
MLOps is not about adopting tools, it's about delivering value

- Gold Rush Age
- FOMO
- Spam emails
- Focus on tools
- 22% have put a model in production
- The real problem: Providing value.



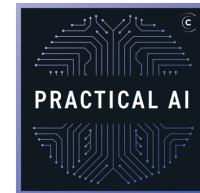
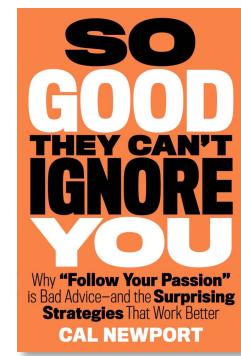
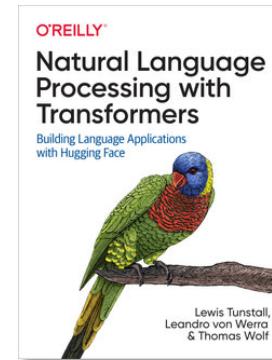
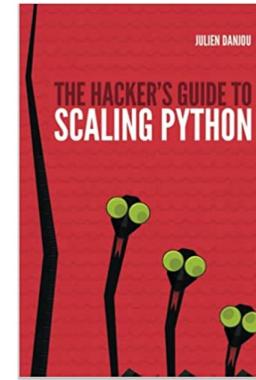
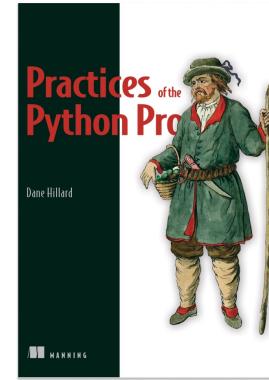
4 | Always be learning

ML is our craft



We should be masters of our craft

- Study
- Stay up-to-date
- Learn regularly
- Build things
- Give back and write



An OCD list of resources

Books

Practices of the Python Pro
Hacker's guide to scaling Python
Designing Data-Intensive Applications
Serious Python

Tutorials

Flask Mega-tutorial
RealPython
Stack Abuse
Kaggle + GitHub

YouTube

CodingTech
Sentdex
Abhishek Thakur
MLOPs Community

Podcasts

Talk Python to Me
Python Bytes
Podcast.__init__
Practical AI

News

PyCoder's Weekly
Medium
Awesome Python Weekly
Reddit RSS

...

Thank you, questions?

(reach out anytime via email: me@duarteocarmo.com)