

Deploying Machine Learning Models

Incorporating Software Development Practises and Reproducibility in Data Science

Carl Stanley



Data Science is one of the most popular professions in the last few years, often being called a "hyped" job. It's especially popular in my flat (hi Jake).

- Data science is a pretty new field in IT
- Comprises of solving business problems with code however you want to tackle it
- A few different areas that can be specialised in:
 - NLP
 - Deep Learning
 - Etc.



Data Science business use cases range from fraud detection and cyber-security, to incident prevention and recommending products



Vodafone & Telefonica(O2) "Cognitive lligence"



I UBER 🗢 \$ 100% 1455 Market Street ~ laim These Discounts

More Restaurants



Machine Learning models to suggest restaurants and menu items

UberEats uses

real time

Uber use Machine Learning to **predict** where rider demand and driver-partner availability will be at various places and times in the future.

Google Flights uses Machine Learning to warn users about delays, and even predict them before they're V 1 8 18% puncedne Meet at 1456 Market St Bank of Am Connecting you to nearby drivers

2:29 am

7:22 am

0

I GA

Learning model to predict ETA routing errors and then use the predicted error to make a correction, to improve accuracy and user



Experimentation is Key in Data Science – the more curious you are, the better your results.

- Jupyter notebooks give you a really easy way to interact with data
- Quick plots, analysis, and markdown integrated
- Can export it as HTML to chat through

🖯 Jupytei	UMAP-post_reduction_clustering Last Checkpoint: a day ago (unsaved changes)		Logout
	View Insert Cell Kernel Help	Not Trusted	No Kernel 🔵
8 + ×	2 🚯 ✔ HRun ■ C >> Code >> 📼 O git nbdiff		
	athena - AthenaQuerier()		
	data_in = athena.execute_query(query, return_df-True) print("Data loaded in with shape ()",format(data in.shape))		
	2021-03-09 13:27:41,324 [INFO] starting Athena query 2021-03-09 13:20:51,305 [INFO] loading ae@efeae.995a-4369-941_bef956f986C8.csv 2021-03-09 13:37:05,306 [INFO] Athena query complete: returning dataframe		
	Data loaded in with shape (6413541, 152)		
	<pre>from sklearn.model_selection import train_test_split</pre>		
	data = data_in.copy()		
	<pre>data[flag_cols] = data[flag_cols].fillna(0) # nmint/"counting tagaget")</pre>		
	# data["TARGET"] = np.where(data["ee_paym"]==1, # data["TARGET"] = np.where(data["ee_paym"]==1,		
	# , np.where(data["ee_simp_contract"]==1, 3, # , np.where(data["ee_simp_contract"]==1, 4, 1)))		
	<pre># , 0) # , 0) # , 0) data = data dron(column=["an naum" "an handsat" "an sime contract" "an sime contract" "an sime contract".</pre>		
	nmint("Dealing with dear data")		
	<pre>print(Searang meet out Searang) for col in dam_cols: data[col] = data[col].fillna(data[col].mode()[0])</pre>		
	<pre>print("Dealing with Experian data") for col in flu cols: data[col] - data[col].fillna(data[col].mode()[0]) data[col] - data[col].astype(int)</pre>		
	<pre>data["TARGET"] = data["segment"] print("Removing bad columns") for col in cols_to_remove: try: data = data.drop(columns-col) except KeyError: pass</pre>		
	<pre>data = stratified_sample(data, ["TARGET"], 100000, seed-42, keep_index-True) # data["index"] = data.index</pre>		
	<pre># data = data.sort_values(by="TARGET", ascending=False) # data = data[0:len(data[data["TARGET"]!-0])*2]</pre>		
	X = data.drop(columns=["TARGET"]) y = data[["TARGET"]]		
	X_xfm = preprocess(X)		
	Dealing with daar data Dealing with Experian data Removing bad columns		
	Preprocessing Data		
	Fetching column lists for transformations Shortening postcodes Fitting discrete columns Fitting discrete columns Filling strings with Nok dealt with		



When you're done with experimentation, you need to make your code reproducible by packaging your work

- Once ready to use a model at scale, you need to create a python package for it
- This is not optional (most of the time)!
- You can automate a lot with CI/CD but the end product should be a callable package to run code against your model



BT

Docker is industry standard for ensuring developers utilise the same exact environment, no matter where they are

- Docker enables developers to produce a reproducible environment through a concept called "containers".
- Traditionally, you set up your env, run your code, then run tests; Docker does all of this at once
- Docker provides a standard framework to ship / deploy / scale your code

What does this mean?

Building 1 docker image for a use case (i.e. <u>ROOT Conda</u>) ensures identical set-up each time, reducing issues with tertiary versioning etc., then overlay your model code and push!



Docker Hub is the source of all of truth for your images – it's GitHub for containers



An Aside – Docker isn't just useful for data science! It can be used to speed up development on any analysis



Jupyter Notebook > .py > Docker Example

~

Notebook (.ipynb)

How Code is Run

- In Cells
- Allows for rapid experimentation
- Works quickly and in modular fashion

Who can Run the Code

• Person with the notebook with GUI

How Tests are Done

• They aren't

Results

- Quick prototyping
- Interactive shell
- Visual outputs for quick analysis

Package (.py)

How Code is Run

- From terminal
- Functional / OOP
- Runs as a package (if you're good at writing code)

Who can Run the Code

- Anyone with the same env as you
- Projects usually contain a requirements file and a README (not always...)

How Tests are Done

• Manually with pytest (if you've written them)

Results

- End to end (e2e) process, run all at once
- Open ended regarding deployment



Dockerfile



How Code is Run

- From terminal, inside container
- Code often cloned straight from Git

Who can Run the Code

• Anyone with docker

How Tests are Done

- Tests are already done when the container is made
- Can add more tests with pytest if you want

Results

• Shippable product that anyone can use with 1 line



CI can enable trust in your work with less active effort, and CD enables agile development and lets developers learn quickly



Data Science isn't just a modelling career; you can't claim you've made a meaningful model until you're actually helping a business achieve it's objectives



Version control + containers + CI/CD together form a great foundation for agile, reproducible data science

Example – Using ML to reduce noisy data in an experiment

Step 1 – Understand the business problem

Step 2 – Gather data

Step 3 – Experiment with feature engineering & Preprocessing

Step 4 – Modelling

Step 5 (where most people stop) – Package the project for quick retraining

Step 6 – Containerise and deploy model to serve at scale

Step 7 – Monitor performance

BT

Step 1: Understanding the business problem is (obviously) key, but most of the time the stakeholders won't know what they want when they ask for it

- Stakeholders have a *very* different perspective on business problems
- Solution design is a process don't rush it
- You're selling your capabilities



Step 2: Gather data – This is where the most value can be gained from a model, so get creative; you're trying to best describe your base for the problem

- This is where the bulk of the work takes place
- Creating your dataset can take a long time and feel like you're not making progress
- Skipping step 3 as it's very case-by-case sensitive



Step 4: Modelling – Unfortunately the simplest solution is *always* the best one, especially in a business context

- This is normally the easiest part of the job especially if you've done a good job in step 2/3
- Keep in mind what you're trying to achieve from a business perspective
- Keep asking yourself 2 questions:
 - How is this improving on the current business process
 - How is this going to be consumed



Step 5: Packaging – If you've built a big fancy model that takes 6 hours to provide inference, you're going to have a bad time

- This is where you enable your model to be called in 1 line of code
- It opens up a world of possibility for how it can be consumed
- Packaging your model and wrapping it in a Flask server enables live predictions at scale

```
"""Frontend for flask app."""
  from flask import Flask, request, jsonify
  from flask sqlalchemy import SQLAlchemy
  import pandas as pd
  from simoxsell.predict import make prediction, run batch scoring
  from simoxsell.train import run training
  from simoxsell import version as version
  from aws tools.s3 tools import send scoring output
  from aws tools.cloudwatch logging import logger as logger
  server = Flask( name )
  @server.route("/simoxsell/helloworld")
v def helloworld():
      """Return Hello World for API Testing."""
      logger.debug("Successfully called route.")
      return "Hello World! You're on the API for simoxsell!"
  @server.route("/simoxsell/make prediction", methods=["GET"])
v def make_test_pred_from_flat_file():
      """Make a test prediction to demonstrate model."""
      _logger.debug("Successfully called route.")
      logger.info("Making test prediction from Athena...")
      data = pd.read csv("/app/data/cs test data.csv")
      preds = make prediction(data)
      logger.info(preds)
```

Step 6: Containerise & Deploy – Most deployment options opt for Docker Containers anyway, so you're a step ahead of the game

- Unfortunately, most businesses disable 90% of Cloud Infrastructure features
- There is a grey area between data science & data engineering here – the ML Engineer is here to fill the boots
- Microservices are cool and hip right now



Stage 7: Monitoring – This is a very poorly done step in general; if this is done right it really helps the lazy developer move on with their life

Why is Monitoring not done as much?

- It's boring
- Businesses aren't mature enough to challenge data science effectively
- It takes a lot of effort to do right

What should we be monitoring?

- Output metrics: The prediction results as & when they're made to check they're reasonable
- Online performance: A comparison of the prediction to the real result, so that we can understand real performance
- Machine-oriented metrics: API response times, CPU usage, etc.



P. M.	API Decisional - Direfere X +								
4	C D O tiet Ser	ente 6565/30000/08	rgf.m.Derm an davriden	otherpits Weathers	ter-briddis-trine			ΨŪΔ	0
5	II ML API Dashboard -					🗮 (A) 🛃			
				Math	Territore				
+	Title					distant sector		1.000	Takes (
	2025-01-16 10:15:15	LASSO 0.0.4444 0		1 3 4	0.007.0005740 0.007 4 5 4 2 8 9				1.00
		Average House Price Prediction Rels (veccord)							
۰	200 K 130 K				- 00				
	374	+							
	n 1919 - Annua Parlitico Arrian (2)	····	2006-01-14 Vol.17-15 Nge Prediction Amount (3) 1	103 K 110	1.0 - Junio Pedition An	uer march			
	Biandard Error of the Maser (SEM)				Bandard Score (2 Score)				
	1004				*				
	-								
		1417 MIL	e tere		12 1010				1925
					2 doors				
	Min Prediction				Max Production				
						1			
-				SALSER					
100									
0									



Thanks for listening!

Questions?



Fig 1. Marcella's EPP class, 2016, colourised