

General Assembly
Data Science Immersive
Capstone Project

Jocelyn Ho 26 February 2021

Background

Purpose

- To utilise data on side effects of drugs
 - Little is known about how data from the Yellow Card Scheme is used

Goal

 To predict whether certain types of side effects/ drugs are more likely to result in a severe or fatal side effect

Metrics

- Classification model predictions
- Mean cross-validated accuracy score

MHRA



- Executive Agency of the Department of Health and Social Care
- Acts on behalf of the Ministers
- Protect and promote public health and patient safety
- Ensure medicines and medical devices meet appropriate standards of safety, quality and efficacy

Yellow Card Scheme



- Run by the MHRA
- The UK system for collecting and monitoring information on safety concerns
 - e.g. suspected side effects or adverse incidents
 - Includes medicines and medical devices
- Relies on voluntary reporting
- Reported by health professionals and the public (patients, carers, parents)

Data Source 1: Yellow Card Interactive Drug Analysis Profiles

Web Scrape - Selenium

- 1. Click onto each alphabet
- 2. Click onto each subset of alphabet
- 3. Click onto each drug
- 4. Get unique yellow card ID for each drug
- 5. Download zip file for each drug
- 6. Unzip all folders



A B C D E F G	H I J K L M N	O P Q R S T	<u>u v w x y z</u>	
Aa-Ad Ae-Ah Ai-Al Aı	m-Ap Aq-Au Av-Az			
Abacavir				
Abatacept			A1 1000 TO	
	To download the cor	mma-separated-v	alues (CSV) data fo	r ABACAVIR, click here.
Abciximab				
<u>Abemaciclib</u>				
<u>Abiraterone</u>				
<u>Acalabrutinib</u>				
		•••••		



Summary of each report for the drug, includes patient demographics





Adverse reactions within each report for the drug (multiple events per case)





(Not used)
Method of
administration for the
drug in each report



Create dataset - SQL:

- Join case to events for each drug on report number
- 2. Add yellow card id to each row
- 3. Append combined dataframe to a new SQL table
- 4. Repeat for all drugs

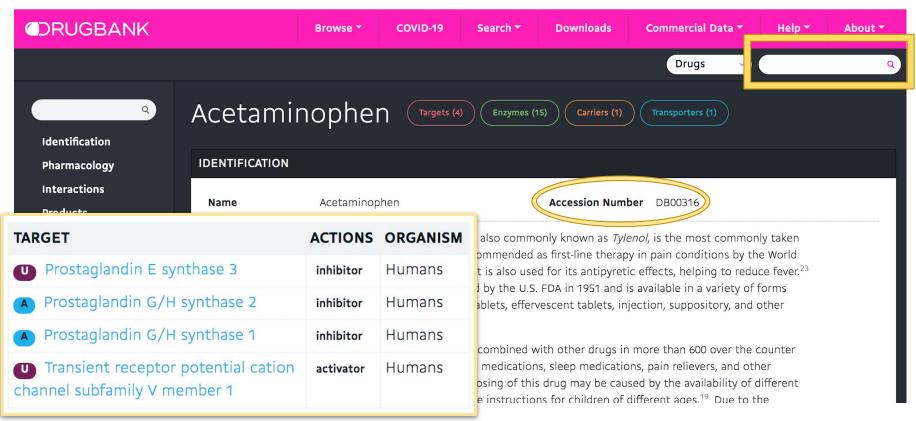
ADR	PT	SOC_ABBREV	FATAL_YN	SEX	AGE_10	RECVD_YEAR	SENDER_TYPE	CONSUMER_YN	HCP_YN	NONSERIOUS_SERIOUS_FATAL_NSF
1	Chills	Genrl	N	Female	40.0	2003	Indirect	N	Υ	s
1	Malaise	Genrl	N	Female	40.0	2003	Indirect	N	Υ	s
1	Vomiting	Gastr	N	Female	40.0	2003	Indirect	N	Υ	s
2	Oral candidiasis	Infec	N	Female	50.0	2003	Direct	N	Υ	N
3	Photosensitivity reaction	Skin	N	Female	60.0	2003	Direct	N	Υ	s

Dataframe 1: Yellow Card

Data Source 2: Drugbank



- 1. Search each drug
- 2. Obtain accession number and extra information on drug targets



Int64Index: 1987 entries, 0 to 2328 Data columns (total 5 columns): Column Non-Null Count Dtype

> 1987 non-null object yc id 1987 non-null int64

object

object

object

Dropped drug_cat column due to insufficient information

target 1568 non-null drug cat 675 non-null db id 1987 non-null

drug

dtypes: int64(1), object(4) memory usage: 93.1+ KB

Dataframe 2: DrugBank

	drug	yc_id	target	db_id
0	abacavir	40046536	['Reverse transcriptase/RNaseH', 'HLA class I	DB01048
1	abatacept	561378321	['T-lymphocyte activation antigen CD80', 'T-ly	DB01281
2	abciximab	231911819	['Integrin beta-3', 'Integrin alpha-IIb', 'Low	DB00054
3	abemaciclib	369408139	['Cyclin-dependent kinase 4', 'Cyclin-dependen	DB12001
4	abiraterone	968368347	['Steroid 17-alpha-hydroxylase/17,20 lyase']	DB05812

Data Source 3: UCI dataset - Drugs.com reviews

Web scrape - Selenium

- 1. Get unique drugNames (brand names)
- 2. Find active ingredients of each product

UCI	and the same				
Machine Learning Repository Center for Machine Learning and Intelligent Systems					

		7				
	drugName	condition	review	rating	date	usefulCount
16374	Mirtazapine	Depression	"I've tried a fe	10	February 28, 2012	22
20647	3 Mesalamine	Crohn's Disease, Mai	"My son has Crohn&#</td><td>8</td><td>May 17, 2009</td><td>17</td></tr><tr><td>15967</td><td>2 Bactrim</td><td>Urinary Tract Infection</td><td>"Quick reduction of</td><td>9</td><td>September 29, 2017</td><td>3</td></tr><tr><td>3929</td><td>contrave</td><td>Weight Loss</td><td>"Contrave combines</td><td>9</td><td>March 5, 2017</td><td>35</td></tr><tr><td>9 76</td><td>3 Cyclafem 1 / 35</td><td>Bi h Control</td><td>"I have been on this</td><td>9</td><td>October 22, 2015</td><td>4</td></tr><tr><td>20808</td><td>Zyon</td><td>Keratosis</td><td>"4 days in on first 2 w</td><td>4</td><td>July 3, 2014</td><td>13</td></tr><tr><td>21589</td><td>Copper</td><td>Birth Control</td><td>"I've had the</td><td>: 6</td><td>June 6, 2016</td><td>1</td></tr><tr><td>16985</td><td>2 Amitriptyline</td><td>Migraine Prevention</td><td>"This has been great</td><td>9</td><td>April 21, 2009</td><td>32</td></tr><tr><td>2329</td><td>Methadone</td><td>Opiate Withdrawal</td><td>"Ive been on Methad</td><td>7</td><td>October 18, 2016</td><td>21</td></tr><tr><td>7142</td><td>Levora</td><td>Birth Control</td><td>"I was on this pill for</td><td>2</td><td>April 16, 2011</td><td>3</td></tr></tbody></table>			

BRAND name

Calpol
Intant
Suspension
Paracetamol

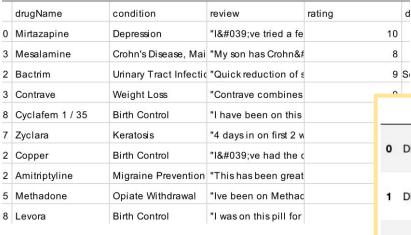
Strawberry flavour

2+ months

DRUG name/ active ingredient

Extract necessary information

- 1. Groupby DrugBank ID
- 2. Obtain unique conditions
- 3. Calculate mean rating for each product





Dataframe 3: Drugs.com reviews

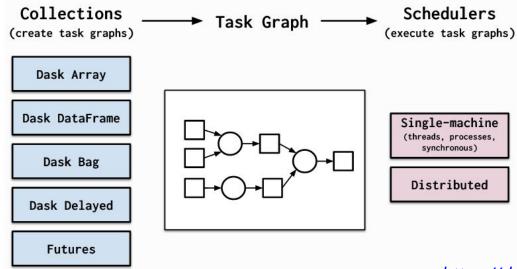
	db_id	drug	condition	review	rating	n_reviews
0	DB00002	['erbitux' 'cetuximab']	('colorectal cancer' 'head and neck cancer' 's	[""i have stage 4 colon cancer with liver mets	6.875000	16
1	DB00003	['dornase alfa' 'pulmozyme']	['cystic fibrosis']	["my outcome/experience with pulmozyme was go	10.000000	2
2	DB00005	['etanercept']	['rheumatoid arthritis' 'psoriatic arthritis'	["took enbrel with methotrexate for four yea	8.246305	203
3	DB00007	['leuprolide' 'lupron depot-ped' 'lupron' 'eli	['uterine fibroids' 'prostate cancer' 'endomet	["i was so skeptical about these shots, scare	6.179487	273
4	DB00008	['pegasys' 'peginterferon alfa- 2a']	['hepatitis c']	["i was treated with this twice on it's	6.791667	24

Merging data - Dask

Dask - Introduction



- Composed of two parts:
 - Dynamic task scheduling
 - 'Big data' collection parallel arrays, dataframes etc.



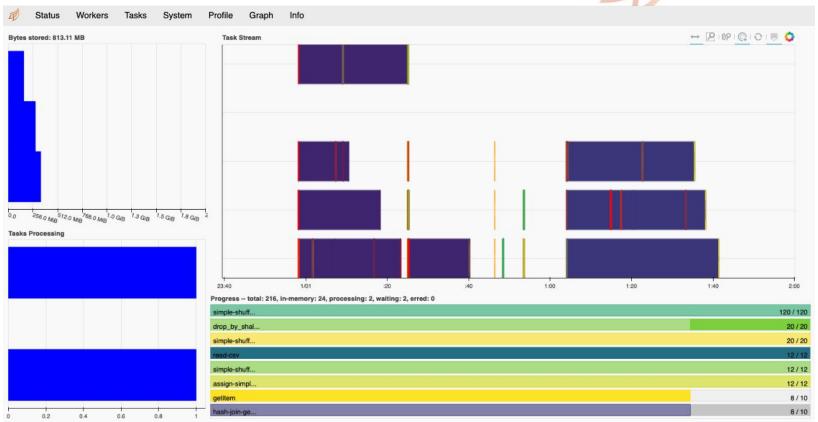
https://docs.dask.org/en/latest/

Dask - Uses and advantages



- Enable efficient parallel computations on single machines
 - Leverages multi-core CPUs and streams data efficiently from disk
 - Allows to swap out the cluster for single-machine schedulers
- Avoids excess memory use:
 - Finds ways to evaluate computations in a low-memory footprint when possible
 - Pulls in chunks of data from disk, does necessary processing, throws away intermediate values asap
- Able to perform computations on moderately large datasets (100GB+) even on relatively low-power laptops.
- Requires no configuration and no setup, installed by default

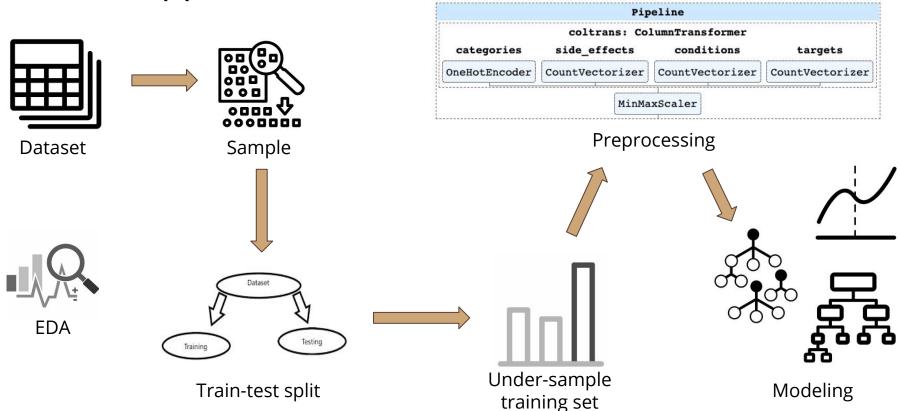




Dask - issues

- Did not close down clients (workers) properly
 - Memory leaked
 - Keep running out of application memory when running models
- Still not enough RAM to process dataset after successfully merging
 - Dropped reviews column
 - Instead used the number of reviews each product received as an indicator
 - Sample 10% of dataset to run models

Overall Approach



Final Dataset

1,755,219

Number of rows

3

Text columns

16

Columns (Prior to processing)

6

Categorical columns

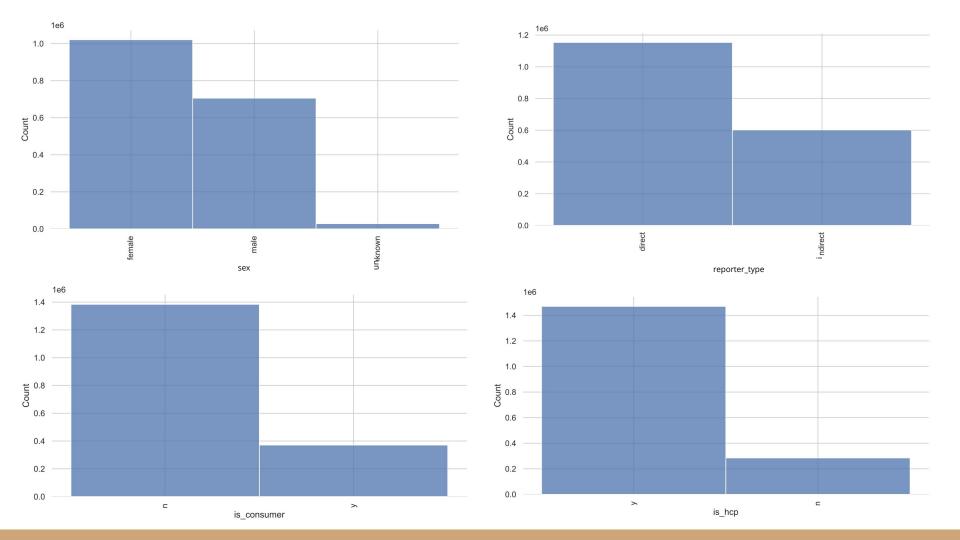
3

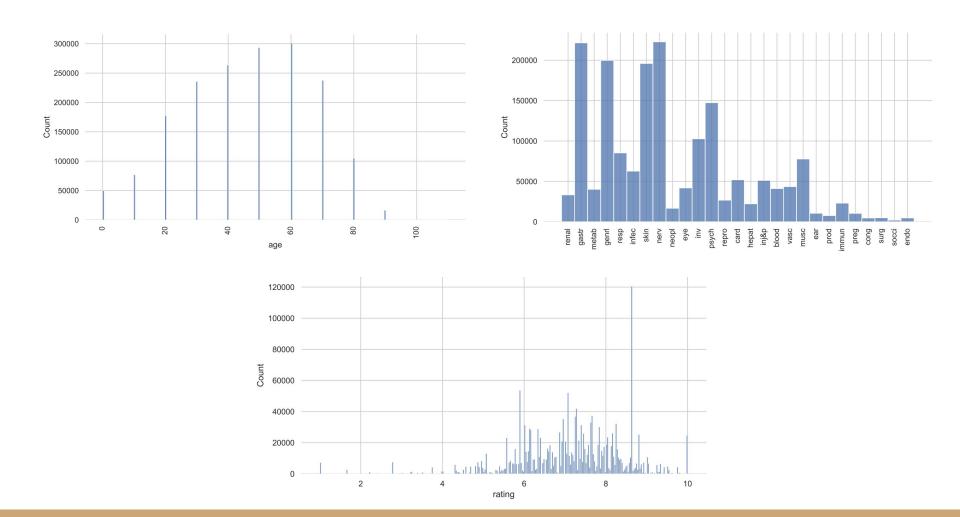
Classes within target variable

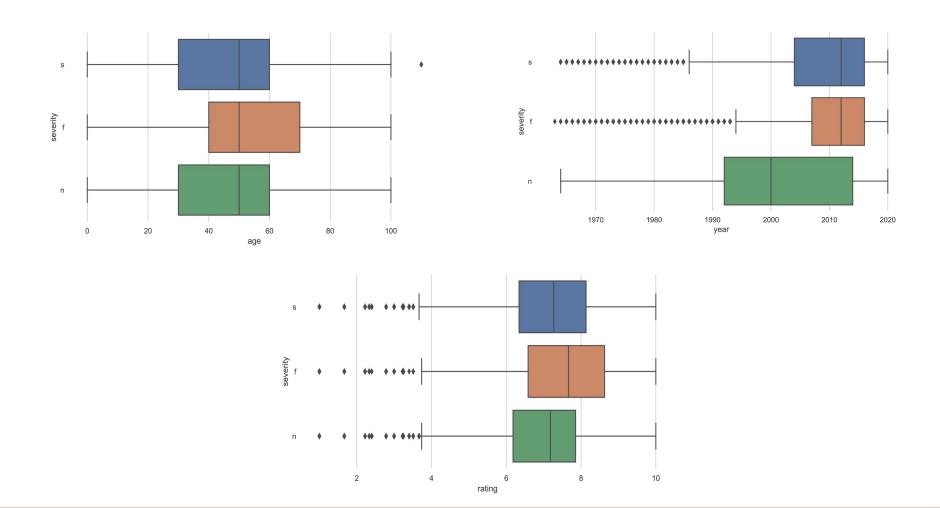
4

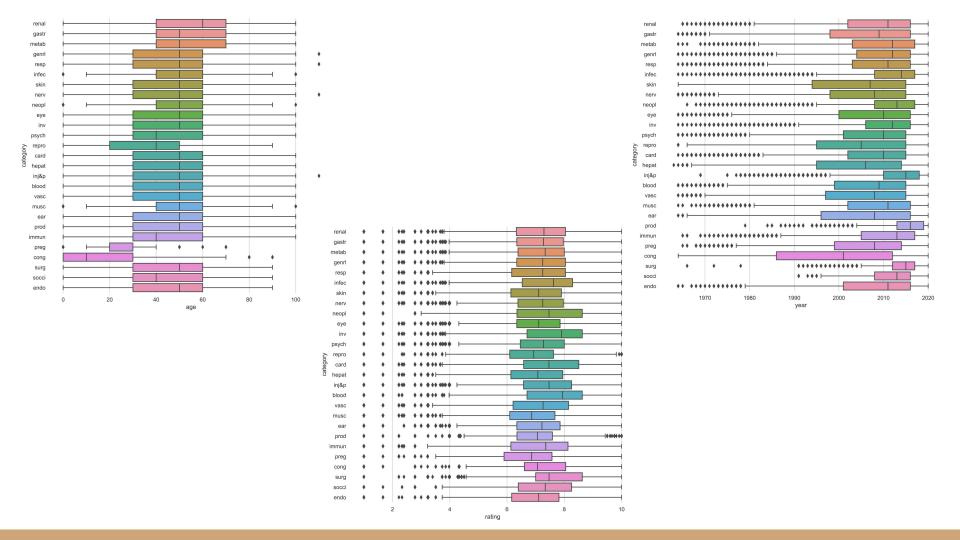
Numeric columns

EDA

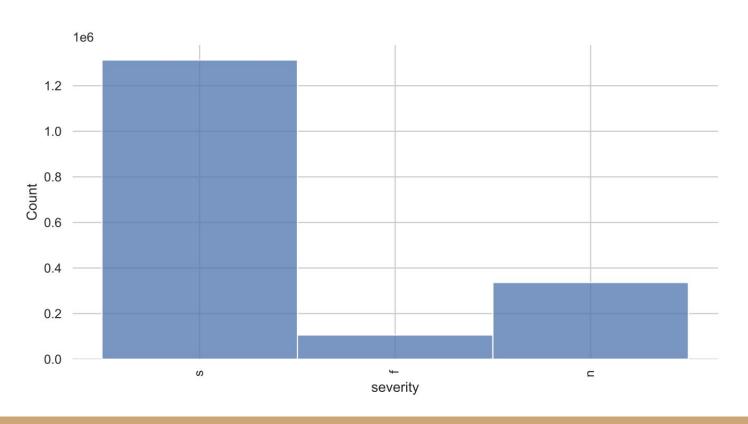








Histogram showing the distribution of target variable classes



NLP

Natural Language Processing

130

130

127

120

117

114

113

111

107

107

103

101

100

97

96

		~	40	40	040
u	ru	2	Ld	12	ets
		0		- 0	

subunit alpha

growth factor

binding protein

calcium channel

channel subunit

voltage dependent

channel subfamily

penicillin binding

sodium dependent

beta adrenergic

protein kinase

dependent type

voltage gated

type opioid

subfamily member

aminobutyric acid

prostaglandin synthase

type calcium

364 muscarinic acetylcholine 238 n b adrenergic alpha 158 142 v 140 n

Drug indications (conditions)

influenza prophylaxis	245	blood pressure
nausea vomiting	119	years ago
birth control	97	started taking
vomiting pregnancy	83	birth control
nausea vomiting pregnancy	83	twice day
blood pressure	76	years old
high blood pressure	76	times day
high blood	76	weight gain
rheumatoid arthritis	73	year old
diabetes type	59	works great
tract infection	58	doctor prescribed
allergic rhinitis	49	feel better
nasal congestion	46	months ago
uterine bleeding	46	dry mouth
abnormal uterine	45	highly recommend
abnormal uterine bleeding	45	stopped taking
urinary tract	44	mood swings
urinary tract infection	39	taking medication
crohn disease	37	worked great
bipolar disorder	35	went away

Drugs.com reviews

63,007784

59.812564

59.187928

50.596942

43.769524

43.570576

41.860535

39.980575

38.192683

35.673650

34.085749

33.624177

32.096512

30.421098

29.078865

28.799457

27.933365

27.421454

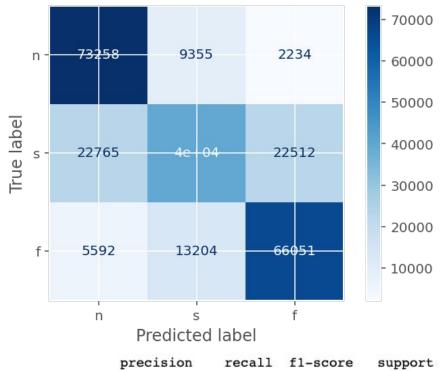
27.337887

25.870941

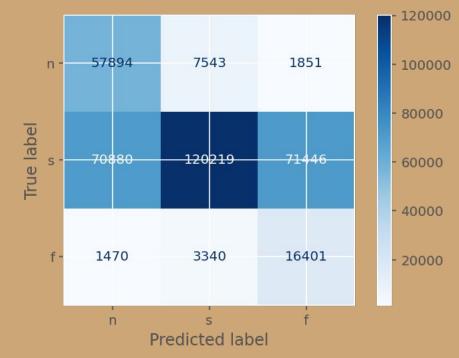
Data Modelling

Model	Best estimators/ parameters	Best Score	Training Score	Test Score		
Random Forest	max_depth = 20, max_features = 90, n_estimators = 500	0.693	0.729	0.507		
Decision Tree Classifier	max_depth = 20, max_features = 70	0.639	0.622	0.524		
AdaBoost	max_depth = 3, max_features = 0.3, n_estimators = 100	0.695	0.699	0.544		
Hist Gradient Boosting Classifier	max_depth = 3	0.682	0.699	0.544		
Naive Bayes (Multinomial)	alpha = 0.0001	0.617	0.616	0.456		
Logistic Regression	'C': 100.0, 'multi_class': 'ovr', 'penalty': 'l2', 'solver': 'saga'	0.699	0.703	0.554		
Baseline accuracy: 74.8% (0.747897)						

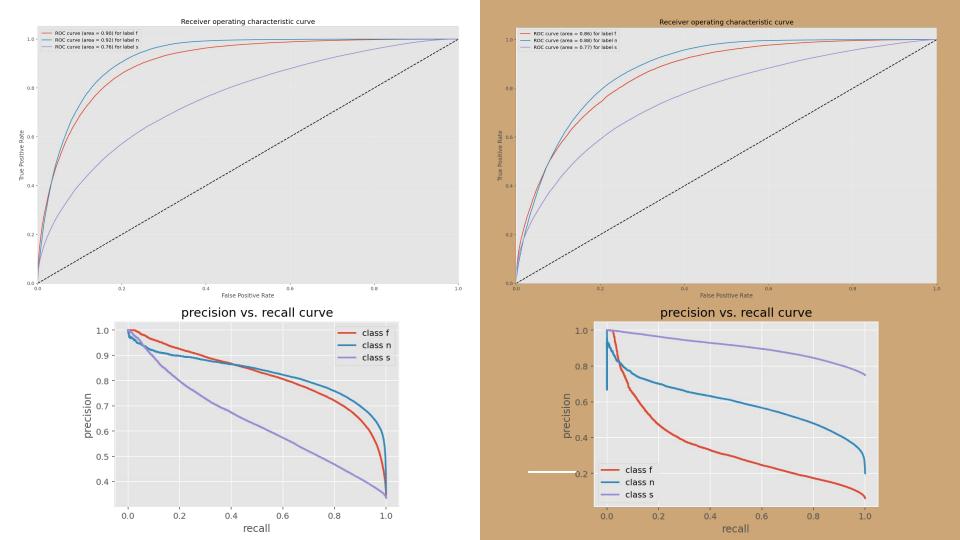
Model: Logistic Regression

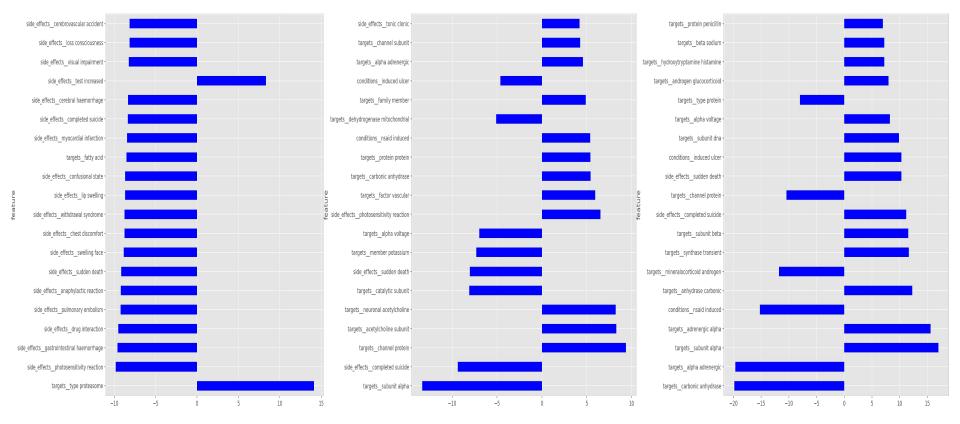


		precision	recall	f1-score	support
	f	0.78	0.73	0.75	90797
	n	0.86	0.72	0.79	101615
	s	0.47	0.64	0.54	62129
accur	асу			0.70	254541
macro	avg	0.70	0.70	0.69	254541
weighted	avg	0.74	0.70	0.71	254541



	precision	recall	f1-score	support
83	0.77	0.18	0.30	89698
81	0.86	0.44	0.59	130244
	0.46	0.92	0.61	131102
accurac	7		0.55	351044
macro av	0.70	0.51	0.50	351044
weighted av	0.69	0.55	0.52	351044





non-serious serious fatal

Limitations

Limitations - Data

Yellow Card reports:

- Relies on voluntary reporting
 - May not truly represent distribution of side effect severity in real life
- Naturally skewed class distribution
 - Low fatal cases:
 - Marketing license approval only if drug safety is established after extensive clinical trials
 - Difficult to pinpoint cause of death comorbidities/ complications/ actual side effects
 - Low non-serious cases:
 - Less likely to report non-serious side effects
- Subjectivity/ ambiguity in determining severity of side effects
 - Can be serious if reporter considers it serious
 - Criteria provided only 1 out of 6 options cover the class 'fatal'

MHRA Yellow Card definitions

"The seriousness criteria... were determined by a working group of the Council for International Organizations of Medical Sciences (CIOMS) and are defined as 6 possible categories... reporters to select one of the following criteria by ticking the appropriate box on the Yellow Card.

- (1) patient died due to reaction
- (2) life threatening
- (3) resulted in hospitalisation or prolonged inpatient hospitalisation
- (4) congenital abnormality
- (5) involved persistent or significant disability or incapacity
- (6) if the reaction was deemed medically significant.

In addition to this, seriousness of reaction terms has also been defined by the MHRA in our medical dictionary. Therefore an ADR report can be serious because the reporter considers the reaction to be serious or because the reaction term itself is considered serious in our medical dictionary."

Limitations

Drugs.com reviews:

- Reviews only up to 2018, no recent reviews
- US brands

DrugBank:

- Information on US-marketed drugs, not the UK
- May not include all the information for each drug
 - e.g. drug classes and drug targets

General:

- Insufficient computational power/ RAM
- Insufficient time
- Fewer data points after under-sampling training set
 - Insufficient data to train model properly

Future work

Data:

- Incorporate drug classes (e.g. ATC codes from WHO) (categories based on what conditions the drug is used for)
- Use ATC codes as unique ID instead of DrugBank number
- Include actual reviews for more NLP/ sentiment analysis
- Dummify drug target column instead of NLP

Modelling:

- Use different under-/over-sampling methods and ratios
- Increase sample size
- Run more extensive models
- Fine tune parameters

Others:

Utilise AWS