

Predictive Health Care

Project Description

PwC is working closely with a network of doctors from Germany. The doctors are interested in finding the possible adverse effects of drugs used to treat neurological pain.

Besides treatment with usual drugs such as gabapentin or lyrica, doctors have found tramal to be a good remedy against neurological pain. If all of these have a similar effect, it would be interesting for the doctors to compare the adverse effects of tramal to gabapentin or lyrica. Depending on the patient, they could then prescribe the medication with fewer adverse effects.

The US has a central database, hosted by the FDA. The FDA Adverse Event Reporting System (FAERS) is a database that contains information on adverse event and medication error reports submitted to FDA. The database is designed to support the FDA's post-marketing safety surveillance programme for drug and therapeutic biologic products.

Data Description

Code for a patient outcome (See table below)

CODE MEANING TEXT

DE Death

LT Life-Threatening

HO Hospitalization - Initial or Prolonged

DS Disability

CA Congenital Anomaly

RI Required Intervention to Prevent Permanent Impairment/Damage

OT Other Serious (Important Medical Event)

NOTE: The outcome from the latest version of a case is provided. If there is more than one outcome, the codes will be line listed.

Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import random

import statsmodels.api as sm
import statsmodels.formula.api as smf
from statsmodels.formula.api import logit

import datetime
from datetime import datetime, timedelta

import scipy.stats

import pandas_profiling
from pandas_profiling import ProfileReport

#matplotlib inline
#sets the default autosave frequency in seconds
%autosave 60
sns.set_style('dark')
sns.set(font_scale=1.2)

plt.rc('axes', titlesize=9)
plt.rc('axes', labelsize=14)
plt.rc('xtick', labelsize=12)
plt.rc('ytick', labelsize=12)

import warnings
warnings.filterwarnings('ignore')

# Web scraping
import requests
from bs4 import BeautifulSoup

# Use Folium library to plot values on a map.
import folium

# Use Feature-Engine library
from feature_engine.encoding import OrdinalEncoder
from feature_engine.outlier_removers import Winsorizer
from feature_engine.import_encoders as fe
from feature_engine.discretisation import EqualFrequencyDiscretiser, EqualFrequencyDiscretiser, DecisionTreeDiscretiser
from feature_engine.encoding import OrdinalEncoder

pd.set_option('display.max_columns',None)
pd.set_option('display.max_rows',None)
pd.set_option('display.width',1000)
pd.set_option('display.float_format', '{:1.2f}'.format)

random.seed(0)
np.random.seed(0)
np.set_printoptions(suppress=True)

Autosaving every 60 seconds
```

2019 Quarter 1

```
In [2]: df_19Q1 Drug = pd.read_table("DRUG19Q1.txt",sep='$')

In [3]: df_19Q1 Drug

Out [3]:
```

	primaryid	caseid	drug_seq	role_cod	drugname	prod_ai	val_vbm	route	dose_vbm	cum_dose_chr	
0	1000661812	10006618	1	PS	LIPITOR	ATORVASTATIN CALCIUM	1	NaN	UNK	nan	
1	1000661812	10006618	2	SS	BENTYL	DICYCLOMINE HYDROCHLORIDE	1	NaN	UNK	nan	
2	1000661812	10006618	3	SS	DICYCLOMINE	DICYCLOMINE HYDROCHLORIDE	1	NaN	UNK	nan	
3	1000808590	10008085	1	PS	SANDOSTATIN LAR DEPOT	OCTREOTIDE ACETATE	1	Intramuscular	20 MG, BIW	nan	
4	1000808590	10008085	2	SS	SANDOSTATIN LAR DEPOT	OCTREOTIDE ACETATE	1	Intramuscular	30 MG, QMO	nan	
...	
1648982	99978793	9997879	3	C	DOLIPRANE	ACETAMINOPHEN	1	NaN	UNK	nan	
1648983	99978793	9997879	4	C	INDOCOLLYRE	INDOMETHACIN	1	NaN	UNK	nan	
1648984	99978793	9997879	5	C	LEDERFOUL	LEUCOVORN CALCIUM	1	NaN	UNK	nan	
...	
1648985	99998153	9999815	1	PS	SUTENT	SUNITINIB MALATE	1	Oral	37.5 MG, CYCLIC (DAILY, 4 WEEKS ON, 2 WEEKS OFF)	nan	
...	
1648986	99998153	9999815	2	C	L-THYROXIN (LEVOTHYROXINE SODIUM)		NaN	2	NaN	UNK	nan

1648987 rows x 20 columns

```
In [4]: df_19Q1 Drug.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1648987 entries, 0 to 1648986
Data columns (total 20 columns):
#   Column              Non-Null Count  Dtype  ---
0   primaryid            1648987 non-null   int64
1   caseid               1648987 non-null   int64
2   drug_seq             1648987 non-null   int64
3   role_cod             1648987 non-null   object
4   drugname             1648977 non-null   object
5   prod_ai              1617259 non-null   object
6   val_vbm              1648987 non-null   int64
7   route                1585139 non-null   object
8   dose_vbm             985266 non-null   object
9   cum_dose_chr          47009 non-null    float64
10  cum_dose_unit         47023 non-null    object
11  decal                 856330 non-null   object
12  rechal               258523 non-null   object
13  lot_num              289346 non-null   object
14  exp_dt               3496 non-null     object
15  nda_num              531373 non-null   float64
16  dose_amt             656764 non-null   float64
17  dose_unit            656769 non-null   object
18  dose_freq            708301 non-null   object
19  dose_freq            417653 non-null   object
dtypes: float64(3), int64(4), object(13)
memory usage: 251.6+ MB

In [5]: df_19Q1 Drug.describe()

Out [5]:
```

	primaryid	caseid	drug_seq	val_vbm	cum_dose_chr	nda_num	dose_amt
count	1648987.00	1648987.00	1648987.00	47009.00	531373.00	656764.00	
mean	190114742.75	15453701.21	6.92	1.02	1859674.13	95447814.91	103145
std	19861637.69	1119314.35	11.64	0.14	1859674.13	13611129397.52	170499.67
min	37624583.00	3762458.00	1.00	1.00	0.00	0.00	0.00
25%	15649482.00	15530068.00	1.00	1.00	75.00	21829.00	5.00
50%	158806961.00	15867207.00	4.00	1.00	600.00	90071.00	30.00
75%	160128414.00	15997378.50	8.00	1.00	4980.00	202439.00	150.00
max	1593052612.00	16168081.00	397.00	2.00	39000000.00	3400955762713.00	7800000.00

```
In [6]: df_19Q1 Drug.columns

Index(['primaryid', 'caseid', 'drug_seq', 'role_cod', 'drugname', 'prod_ai', 'val_vbm', 'route', 'dose_vbm', 'cum_dose_chr', 'dose_unit', 'dose_freq', 'dose_amt'], dtype='object')

In [7]: df_19Q1 Drug.drop(['caseid', 'drug_seq', 'prod_ai', 'val_vbm', 'route', 'dose_vbm', 'cum_dose_chr', 'cum_dose_unit', 'dose_freq', 'dose_amt'], dtype='object')

In [8]: df_19Q1 Drug

Out [8]:
```

	primaryid	role_cod	drugname
0	1000661812	PS	LIPITOR
1	1000661812	SS	BENTYL
2	1000661812	SS	DICYCLOMINE
3	1000808590	PS	SANDOSTATIN LAR DEPOT
4	1000808590	SS	SANDOSTATIN LAR DEPOT
...
1648982	99978793	C	DOLIPRANE
1648983	99978793	C	INDOCOLLYRE
1648984	99978793	C	LEDERFOUL
1648985	99998153	PS	SUTENT
1648986	99998153	C	L-THYROXIN (LEVOTHYROXINE SODIUM)

1648987 rows x 3 columns

Load the outcome file

```
In [9]: df_19Q1_OUT = pd.read_table("OUTC19Q1.txt", sep='$')

In [10]: df_19Q1_OUT

Out [10]:
```

	primaryid	caseid	outc_cod
0	1000808590	10008085	HO
1	1000808590	10008085	OT
2	100107484	10010748	HO
3	100107484	10010748	HO
4	100171328	10017132	OT
...
310657	99885794	9988579	DS
310658	99885794	9988579	OT
310659	99924278	9992427	HO
310660	99978793	9997879	HO
310661	99998153	9999815	HO

310662 rows x 3 columns

```
In [11]: plt.figure(figsize=(20,10))
sns.countplot(x=df_19Q1_OUT.outc_cod)
plt.suptitle('Patient Outcome 2019 Q1', x=0.5, y=1.02, ha='center', fontsize=20)
plt.tight_layout()
plt.show()
```

Patient Outcome 2019 Q1

```
In [12]: dfQ1 = pd.merge(left=df_19Q1_Drug, right=df_19Q1_OUT, how='inner', on='primaryid')

In [13]: dfQ1

Out [13]:
```

	primaryid	role_cod	drugname	caseid	outc_cod
0	1000808590	PS	SANDOSTATIN LAR DEPOT	10008085	OT
1	1000808590	SS	SANDOSTATIN LAR DEPOT	10008085	HO
2	1000808590	SS	SANDOSTATIN LAR DEPOT	10008085	HO
3	1000808590	SS	SANDOSTATIN LAR DEPOT	10008085	OT
4	1000808590	SS	SANDOSTATIN LAR DEPOT	10008085	HO
...
1684340	99978793	C	DOLIPRANE	9997879	HO
1684341	99978793	C	INDOCOLLYRE	9997879	HO
1684342	99978793	C	LEDERFOUL	9997879	HO
1684343	99998153	PS	SUTENT	9999815	HO
1684344	99998153	C	L-THYROXIN (LEVOTHYROXINE SODIUM)	9999815	HO

1684345 rows x 5 columns

```
In [14]: dfQ1["drugname"].value_counts()

Out [14]:
```

drugname	count
HUMIRA	18293
NEKIMIN	17213
PRENIZONE	16547
XOLAIR	15772
METFORMIN	14815
...	...
ALIVIA	1
TRISOLIN (ASCORBIC ACID)BIOTIN;CALCIUM PANTOTHENATE;CYANOCOBALAMIN;ER	1
CIMICIFUGA	1
morphine sul	1
Multi vitamin infusion (Ascorbic acid, Pyridoxine hydrochloride, Tocoph	1

Name: drugname, Length: 50808, dtype: int64

```
In [15]: (dfQ1["drugname"] == "LYRICA").value_counts()

Out [15]:
```

drugname	count
LYRICA	1677652
TRAMAL	6093

Name: drugname, dtype: int64

```
In [16]: (dfQ1["drugname"] == "TRAMAL").value_counts()

Out [16]:
```

drugname	count
TRAMAL	1684095
LYRICA	230

Name: drugname, dtype: int64

```
In [17]: (dfQ1["drugname"] == "GABAPENTIN").value_counts()

Out [17]:
```

drugname	count
GABAPENTIN	1675073
LYRICA	8772

Name: drugname, dtype: int64

Select the three medications

```
In [18]: dfQ2 = dfQ1[(dfQ1["drugname"] == "LYRICA") | (dfQ1["drugname"] == "TRAMAL") | (dfQ1["drugname"] == "GABAPENTIN")]

In [19]: dfQ2

Out [19]:
```

	primaryid	role_cod	drugname	caseid	outc_cod
194	100442633	SS	GABAPENTIN	10044263	OT
195	100442633	SS	GABAPENTIN	10044263	DS
896	100716802	SS	LYRICA	10071680	OT
938	100904069	PS	LYRICA	10090406	HO
939	100904069	PS	LYRICA	10090406	OT
...
1683997	99227823	C	GABAPENTIN	9922783	OT
1684012	99252085	SS	LYRICA	9925208	HO
1684120	99314732	SS	LYRICA	9931473	HO
1684322	99885794	C	TRAMAL	9988579	DS
1684323	99885794	C	TRAMAL	9988579	DS
1684323	99885794	C	TRAMAL	9988579	DS

15715 rows x 5 columns

```
In [20]: dfQ2.groupby(by=["outc_cod", "drugname"])["primaryid"].count()

Out [20]:
```

outc_cod	drugname	count
CA	GABAPENTIN	15
CA	LYRICA	2
DE	GABAPENTIN	1358
DE	TRAMAL	352
DS	GABAPENTIN	213
DS	LYRICA	162
DS	TRAMAL	4
HO	GABAPENTIN	2814
HO	LYRICA	2200
HO	TRAMAL	86
LT	GABAPENTIN	223
LT	LYRICA	201
LT	TRAMAL	22
OT	GABAPENTIN	4147
OT	TRAMAL	101
OT	LYRICA	2

Name: primaryid, dtype: int64

```
In [21]: dfQ2.to_csv("dfQ2.csv", index=False)

In [22]: tabledfQ2 = pd.pivot_table(data=dfQ2, index=["outc_cod"], columns=["drugname"], aggfunc="count")["caseid"]

In [23]: tabledfQ2

Out [23]:
```

drugname	GABAPENTIN	LYRICA	TRAMAL
CA	15.00	2.00	nan
DE	1358.00	352.00	37.00
DS	213.00	162.00	4.00
HO	2814.00	2200.00	86.00
LT	223.00	201.00	22.00
OT	4147.00	3776.00	101.00
RI	2.00	nan	nan

Data Visualization

```
In [24]: plt.figure(figsize=(20,20))

g = sns.catplot(x="outc_cod", col="drugname",
               kind="count", data=dfQ2,
               height=5, aspect=1)

g.set_ylabels("Patient Outcome")
g.set_xlabels("Counts")
g.set_xticklabels(rotation=90)

plt.suptitle("Comparison of 3 Drugs 2019 Q1", x=0.5, y=1.2, ha="center", fontsize=20)
plt.show()
```

Figure size 1440x1440 with 0 Axes

Comparison of 3 Drugs 2019 Q1

2019 Quarter 2

```
In [25]: df_19Q2 Drug = pd.read_table("DRUG19Q2.txt",sep='$')

In [26]: df_19Q2 Drug

Out [26]:
```

	primaryid	caseid	drug_seq	role_cod	drugname	prod_ai	val_vbm	route	dose_vbm	cum_dose_chr
0	1000661813	10006618	1	PS	LIPITOR	ATORVASTATIN CALCIUM	1	NaN	UNK	na
1	1000661813	10006618	2	SS	BENTYL	DICYCLOMINE HYDROCHLORIDE	1	NaN	UNK	na
2	1000661813	10006618	3	SS	DICYCLOMINE	DICYCLOMINE HYDROCHLORIDE	1	NaN	UNK	na
3	1000808410	10008084	1	PS	INFLIXIMAB, RECOMBINANT	INFLIXIMAB	1	Intravenous (not otherwise specified)	NaN	na
4	1000808410	10008084	2	C	METHOTREXATE	METHOTREXATE	1	Unknown	NaN	na
...
1863005	999981612	9999816	15	C	CALCIUM -VIT D	CALCIUM/ITAMIN D	1	Oral	1 DO, 3X/DAY	na
1863006	999981612	9999816	16	C	RESTASIS	CYCLOSPORINE	1	NaN	1 GTT, 2X/DAY (1 GTT IN EACH AFFECTED EYE EVER...	na
1863007	999981612	9999816	17	C	TRAZODONE	TRAZODONE HYDROCHLORIDE	1	Oral	100 MG, 3X/DAY	na
1863008	999981612	9999816	18	C	TIZANIDINE	TIZANIDINE	1	Oral	8 MG, DAILY	na
1863009	999981612	9999816	19	C	BOTOX	ONABOTULINUMTOXINA	1	NaN	UNK, EVERY 3 MONTHS	na

1863010 rows x 20 columns

```
In [27]: df_19Q2 Drug.info(null_counts=True)

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1863010 entries, 0 to 1863009
Data columns (total 20 columns):
#   Column              Non-Null Count  Dtype  ---
0   primaryid            1863010 non-null   int64
1   caseid               1863010 non-null   int64
2   drug_seq             1863010 non-null   int64
3   role_cod             1863010 non-null   object
4   drugname             1862993 non-null   object
5   prod_ai              182186 non-null   object
6   val_vbm              1863010 non-null   int64
7   route                1521462 non-null   object
8   dose_vbm             1068362 non-null   object
9   cum_dose_chr          47352 non-null    float64
10  cum_dose_unit         47247 non-null    object
11  decal                 95618 non-null   object
12  rechal               285765 non-null   object
13  lot_num              319636 non-null   object
14  exp_dt               5092 non-null     object
15  nda_num              824549 non-null   object
16  dose_amt             578493 non-null   float64
17  dose_unit            783979 non-null   object
18  dose_freq            824549 non-null   object
19  dose_freq            482814 non-null   object
dtypes: float64(3), int64(4), object(13)
memory usage: 284.3+ MB

In [28]: df_19Q2 Drug.describe()

Out [28]:
```

	primaryid	caseid	drug_seq	val_vbm	cum_dose_chr	nda_num	dose_amt
count	1863010.00	1863010.00	1863010.00	47352.00	578493.00	783999.00	
mean	197427186.64	15771903.29	7.21	1.02	46024.45	9722295.22	851.59
std	210867983.25	1135698.67	10.49	0.14	1699622.09	14333987978.39	157949.90
min	39490342.00	3949034.00	1.00	1.00	0.00	0.00	0.00
25%	159405832.00	15802600.00	2.00	1.00	78.00	21344.00	3.00
50%	16233237.50	16204212.50	4.00	1.00	600.00	78787.00	20.00
75%	163698028.50	16349342.00	9.00	1.00	5850.00	202192.00	108.00
max	1636576010.00	16577055.00	326.00	2.00	19900000.00	228000731000.00	7200000.00

```
In [29]: df_19Q2 Drug.columns

Out [29]:
```

	primaryid	caseid	role_cod	drugname	prod_ai	val_vbm	route	dose_vbm	cum_dose_chr
0	1000661813	10006618	PS	LIPITOR	ATORVASTATIN CALCIUM	1	NaN	UNK	na
1	1000661813	10							