

Credit Card Fraud - Prevention through Data Science

Team 3 - Brianne Bell, Michael Nguyen, and Dave Friesen

ADS-505-02-FA22

GitHub link: <https://github.com/davefriesen/cc-fraud-prevention>

```
In [1]: # Authors = ['Brianne Bell', 'Michael Nguyen', 'Dave Friesen']
        __contact__ = ['bbell1@csandiego.edu', 'michaelsinguyen@csandiego.edu', 'dfriesen@csandiego.edu']
        __date__ = '2022-10-19'
        __version__ = '1.0.0'
```

Setup

```
In [2]: # Set working directory for 'custom' (profiler, model_process) .py find and data access
        # Importing to CUDASDKTOR
        import os, sys

        src_dir = 'C:/Users/breel.B-E-BELL/OneDrive/Documents/USD_MastersAppliedDataScience/ADS-505 ADSforBusiness/FI
        data_dir = 'C:/Users/breel.B-E-BELL/OneDrive/Documents/USD_MastersAppliedDataScience/ADS-505 ADSforBusiness/F

        src_dir = 'Users/davidfriesen/Desktop/OneDrive/projects/cc-fraud-protection/src/'
        data_dir = 'Users/davidfriesen/Desktop/OneDrive/projects/cc-fraud-protection/data/'

        sys.path.append(src_dir)
        os.chdir(data_dir)
```

```
In [3]: # Import basic libraries
        import numpy as np
        import pandas as pd

        # Import visualization libraries
        import matplotlib.pyplot as plt
        %matplotlib inline

        # Import custom XBA library
        from profiler import profile, profile_cat

        # Import sklearn and performance evaluation libraries
        from sklearn.model_selection import train_test_split

        from sklearn.linear_model import LogisticRegression, Perceptron
        from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
        from sklearn.svm import LinearSVC

        from sklearn.neighbors import KNeighborsClassifier

        from sklearn.tree import DecisionTreeClassifier
        from sklearn.ensemble import RandomForestClassifier
        from sklearn.ensemble import AdaBoostClassifier

        from sklearn.neural_network import MLPClassifier

        # Import custom model library
        from model_process import ModelProcess

        # Import utility libraries
        import copy
        from utils import ctrl_count
```

```
In [4]: # Set basic np, pd, and plt output defaults (keeping this code 'clean')
        %run -i 'src_dir/defaults.py'
```

Set simple notebook style and auto-avoid scrollable output windows

```
In [5]: %html
<style>
    h1, h2, h3 {font-size: 0.9em; color: #0070C0;}
    h1, h2, h3 {line-height: 1.2em !important; margin-top: .5em !important;}
    h4, p {font-size: 1.2em !important; line-height: 1.2em !important;}
</style>
```

```
In [6]: %javascript
    %Ipython.OutputArea.prototype._should_scroll = function(lines) {
        return false;
    }
```

Data Load and Validation

Fraud e-commerce: 'Fraud_Data' and 'IpAddress_to_Country'

```
In [7]: # Ref: https://www.kaggle.com/datasets/vb1nh002/fraud-e-commerce?resource=download

        # Create and confirm fraud dataframe
        fraud_fname = data_dir + 'c-Fraud_Data.csv'
        fraud_ctrl = ctrl_count(fraud_fname)
        fraud_df = pd.read_csv(fraud_fname, on_bad_lines = 'skip', low_memory = False)

        print('\nLoad: file=0d, import=0d, delta=0d')
        (fraud_ctrl, len(fraud_df), fraud_ctrl - len(fraud_df))
        print('\nFraud rows, columns: ', fraud_df.shape)
        print('\nFraud dataframe:')
        print('\nFraud data:\n', fraud_df.head(5))

        # Confirm no duplicate rows, if have save training/test sets - n/a here
        print('\nIs naive', np.array_equal(fraud_data, fraud_df.columns))

        Fraud: file=151113, import=151112, delta=1

        Fraud rows, columns: (151112, 11)

        Fraud dataframe:
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 151112 entries, 0 to 151111
        Data columns (total 11 columns):
         #   Column                Non-Null Count  Dtype
        ---  ---
         0   user_id                151112 non-null    int64
         1   signup_time            151112 non-null    object
         2   purchase_time          151112 non-null    object
         3   purchase_value         151112 non-null    int64
         4   device_id              151112 non-null    object
         5   source                 151112 non-null    object
         6   browser                151112 non-null    object
         7   sex                    151112 non-null    object
         8   age                    151112 non-null    int64
         9   ip_address             151112 non-null    float64
        10   class                  151112 non-null    int64
        dtypes: float64(1), int64(4), object(6)
        memory usage: 12.7+ MB

        Fraud data:
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 0 entries, 0 to 0
        Data columns (total 11 columns):
         #   Column                Non-Null Count  Dtype
        ---  ---
         0   user_id                0 non-null         object
         1   signup_time            0 non-null         object
         2   purchase_time          0 non-null         object
         3   purchase_value         0 non-null         object
         4   device_id              0 non-null         object
         5   source                 0 non-null         object
         6   browser                0 non-null         object
         7   sex                    0 non-null         object
         8   age                    0 non-null         object
         9   ip_address             0 non-null         object
         10  class                  0 non-null         object
        dtypes: object(11)
        memory usage: 0+ MB
```

```
In [8]: # Ref: https://www.kaggle.com/datasets/vb1nh002/fraud-e-commerce?resource=download

        # Create and confirm IP dataframe
        ip_fname = data_dir + 'c-ipAddress_to_Country.csv'
        ip_ctrl = ctrl_count(ip_fname)
        ip_df = pd.read_csv(ip_fname, on_bad_lines = 'skip', low_memory = False)

        print('\nIP: file=0d, import=0d, delta=0d')
        (ip_ctrl, len(ip_df), ip_ctrl - len(ip_df))
        print('\nIP rows, columns: ', ip_df.shape)
        print('\nIP dataframe:')
        print('\nIP data:\n', ip_df.head(5))

        IP: file=138847, import=138846, delta=1

        IP rows, columns: (138846, 3)

        IP dataframe:
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 138846 entries, 0 to 138845
        Data columns (total 3 columns):
         #   Column                Non-Null Count  Dtype
        ---  ---
         0   lower_bound_ip_address 138846 non-null    float64
         1   upper_bound_ip_address 138846 non-null    int64
         2   country                138846 non-null    object
        dtypes: float64(1), int64(1), object(1)
        memory usage: 3.2+ MB

        IP data:
         lower_bound_ip_address  upper_bound_ip_address  country
        0      16777216.000000      16777741.000000      Australia
        1      16777472.000000      16777727.000000              China
        2      16777728.000000      16778239.000000              China
        3      16778240.000000      16779263.000000      Australia
        4      16792640.000000      16781311.000000              China
```

Data Profiling

```
In [9]: # Profile base dataframe(s)
        profile(fraud_df)
        profile(ip_df)
```

	Dtype	count	unique	na	na%	mean	std	min	max	skew(>=3)	<v0.01	WIF(>=10)
user_id	int64	151112	151112			220717.0	115369.3	2.0	400000.0			
signup_time	object	151112	151112									
purchase_time	object	151112	150679									
purchase_value	int64	151112	122			36.9	18.3	9.0	1540			
device_id	object	151112	137956									
source	object	151112	3									
browser	object	151112	5									
sex	object	151112	2									
age	int64	151112	58			33.1	8.6	18.0	76.0			
ip_address	float64	151112	143512			21521453931.0	1248497030.1	52093.5	42948504987			
class	int64	151112	2			0.1	0.3		1.0			

	Dtype	count	unique	na	na%	mean	std	min	max	skew(>=3)	<v0.0
lower_bound_ip_address	float64	138846	138846			2724531562.5	897521519.7	16777216.0	3768096128.0		
upper_bound_ip_address	int64	138846	138846			2724557082.2	897497915.5	16777471.0	3758096383.0		
country	object	138846	235								

Univariate Analysis and Data Preparation

Attribute Names (confirm/update)

```
In [10]: # Rename column to simplify and/or avoid usage issues later - tbd

        # example: 'homeowner dummy': 'homeowner', 'gender dummy': 'gender'
        funder_df.rename(columns=cols, inplace=True)
```

Target and Preliminary Feature Identification (preliminary)

```
In [11]: # Preliminarily identify attributes and types for best further processing
        id_ordinal_cols = ['user_id', 'device_id', 'ip_address']

        cat_nominal_cols = ['source', 'browser', 'sex']
        cat_ordinal_cols = []
        cat_binary_cols = []

        date_interval_cols = ['signup_time', 'purchase_time']

        num_interval_cols = ['age']
        num_ratio_cols = ['purchase_value']

        # Set classification target
        target_cls_col = ['class']

        # Set regression target (n/a for this exercise)
        target_reg_col = []

        # Set preliminary feature reduction (removed below)
        reduce_x_cols = id_nominal_cols + date_interval_cols + target_cls_col + target_reg_col
```

Data Types (confirm/update)

```
In [12]: # Making columns in question float type
        ip_df['lower_bound_ip_address'] = ip_df['lower_bound_ip_address'].astype('float') # technically already is float
        ip_df['upper_bound_ip_address'] = ip_df['upper_bound_ip_address'].astype('float') # integer, needs to be float
        df['ip_address'] = fraud_df['ip_address'].astype('float') # technically already is float
        # checking
        print(ip_df.info())

        # other examples:
        df['cat_nominal_cols'] = df['cat_nominal_cols'].astype('Int64')
        df['date_interval_cols'] = df['date_interval_cols'].astype('datetime64[ns]')

        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 138846 entries, 0 to 138845
        Data columns (total 3 columns):
         #   Column                Non-Null Count  Dtype
        ---  ---
         0   lower_bound_ip_address 138846 non-null    float64
         1   upper_bound_ip_address 138846 non-null    float64
         2   country                138846 non-null    object
        dtypes: float64(2), object(1)
        memory usage: 3.2+ MB
        None
```

Data Enrichment

Combining dataframes where IP in fraud_df is within bounds of the IP bounds of ip_df to get Country name

Filling in Country values

If the IP address doesn't fit into the bounds of the ip dataframe then it is assumed they are from VPN users.

```
In [13]: # Function takes in ip_address (fraud_df) and compares it to upper/lower bounds in ip_df to return country name
        def ip_to_country(ip):
            try:
                return ip_df.country[
                    (ip_df.lower_bound_ip_address <= ip) &
                    (ip_df.upper_bound_ip_address >= ip)].iloc[0]
            except IndexError:
                return 'VPNUser'

        # calling the function
        fraud_df['ip_country'] = fraud_df['ip_address'].apply(ip_to_country)

        # saving to csv file to speed up when opening notebook next time
        fraud_df.to_csv('C:/Users/breel.B-E-BELL/OneDrive/Documents/USD_MastersAppliedDataScience/ADS-505 ADSforBusiness/
        # checking
        fraud_df.head(5)
```

```
Out[13]:
```

user_id	signup_time	purchase_time	purchase_value	device_id	source	browser	sex	age	ip_address	class	ipCountry	
0	22058	2015-02-24 22:55:49	2015-04-18 02:47:11	34	QVPSFPUJOCKZAR	SEO	Chrome	M	39	732758368.80	0	Ja
1	333320	2015-06-07 20:39:50	2015-06-08 01:38:54	16	EOGFQPIZPYXFX	Ads	Chrome	F	53	350311387.87	0	Un
2	1359	2015-01-01 18:52:44	2015-01-01 18:52:45	15	YSSKYQSJHPPLJ	Opera	Opera	M	53	2621473820.11	1	Un
3	150084	2015-04-28 21:13:25	2015-05-04 13:54:50	44	ATOTXKYKUUDQUN	SEO	Safari	M	41	3840542443.91	0	VPN
4	221365	2015-07-21 07:09:52	2015-09-09 18:40:53	39	NAUITBZFJKHWW	Ads	Safari	M	45	415583117.45	0	Un

```
In [14]: # Making how many rows unknown
        print('There were %i unknown Country values (likely using a VPN service)' % fraud_df['ipCountry'].value_counts()
        print('This is %.1f percent of the data' % (fraud_df['ipCountry'].value_counts() / len(fraud_df) * 100))

        There were 21966 unknown Country values (likely using a VPN service)
        This is 14.5 percent of the data
```

Datetime Manipulations

```
In [15]: # calculating time difference between signing up and purchasing
        fraud_df['time_diff_day'] = (pd.to_datetime(fraud_df['purchase_time']) -
                                     pd.to_datetime(fraud_df['signup_time'])).astype('timedelta64[D]')

        fraud_df.head(5)
```

```
Out[15]:
```

user_id	signup_time	purchase_time	purchase_value	device_id	source	browser	sex	age	ip_address	class	ipCountry	
0	22058	2015-02-24 22:55:49	2015-04-18 02:47:11	34	QVPSFPUJOCKZAR	SEO	Chrome	M	39	732758368.80	0	Ja
1	333320	2015-06-07 20:39:50	2015-06-08 01:38:54	16	EOGFQPIZPYXFX	Ads	Chrome	F	53	350311387.87	0	Un
2	1359	2015-01-01 18:52:44	2015-01-01 18:52:45	15	YSSKYQSJHPPLJ	SEO	Opera	M	53	2621473820.11	1	Un
3	150084	2015-04-28 21:13:25	2015-05-04 13:54:50	44	ATOTXKYKUUDQUN	SEO	Safari	M	41	3840542443.91	0	VPN
4	221365	2015-07-21 07:09:52	2015-09-09 18:40:53	39	NAUITBZFJKHWW	Ads	Safari	M	45	415583117.45	0	Un

```
In [16]: #fraud_df = fraud_df.drop(['Unnamed: 0'], axis=1)
        # renaming
        df = fraud_df.copy()
        df.rename(columns = {'signup_time':'signup', 'purchase_time':'purchase', 'ipCountry':'ip_country'})
        df = df.drop(['user_id'], axis=1)
```

```
In [17]: # converting from date to datetime (date and time separate columns) instead of string
        df['signup_date'] = pd.to_datetime(df['signup'])
        df['purchase_date'] = pd.to_datetime(df['purchase'])
        df['signup_time'] = pd.to_datetime(df['signup']).dt.strftime('%H:%M:%S')
        df['purchase_time'] = pd.to_datetime(df['purchase']).dt.strftime('%H:%M:%S')

        #checking
        df.head(3)
```

```
In [18]: # making date a month instead of full date
        df['purchase_month'] = pd.to_datetime(df['purchase']).dt.month_name()
        df['purchase_date'] = pd.to_datetime(df['purchase']).dt.month_name()

        # checking, can comment out
        df.head(2)
```

```
In [19]: # looking days of week
        df['signup_day'] = pd.to_datetime(df['signup_date']).dt.day_name()
        df['purchase_day'] = pd.to_datetime(df['purchase_date']).dt.day_name()

        # zooming in on just weekend (saturday, sunday) or weekday (monday-friday)
        # np.where(condition, when true it's this, when not it's this)
        df['signup_weektype'] = np.where(df['signup_day'] == 'Saturday' | #first condition followed by | (or)
                                         df['purchase_day'] == 'Sunday', # second condition
                                         'Weekend', 'Weekday') # if true response, if false response

        df['purchase_weektype'] = np.where(df['purchase_day'] == 'Saturday' | #first condition followed by | (or)
                                         df['purchase_day'] == 'Sunday', # second condition
                                         'Weekend', 'Weekday') # if true response, if false response

        # checking
        df.head(5)
```

```
In [20]: # splitting time into categories instead of times for leaving data
        # AM: 0-11, PM: 12-23
        sign_weektype = (pd.to_datetime(df['signup_time']).dt.hour.between(0, 11), # am
                        (pd.to_datetime(df['signup_time']).dt.hour.between(12, 23)), # pm
                        )
        time_values = ['AM', 'PM']

        # bringing the splitting into the df3 via np.select (timewise, time_values):
        df['signup_timegroup'] = np.select(sign_weektype, time_values)

        # purchase
        purch_weektype = (pd.to_datetime(df['purchase_time']).dt.hour.between(0, 11), # am
                        (pd.to_datetime(df['purchase_time']).dt.hour.between(12, 23)), # pm
                        )

        # bringing the splitting into the df3 via np.select (timewise, time_values):
        df['purchase_timegroup'] = np.select(purch_weektype, time_values)

        #checking
        df.head(3)
```

```
Out[21]:
```

signup	purchase	purchase_value	device_id	source	browser	sex	age	ip_address	class	...	purchase_time	signu
0	2015-02-24 22:55:49	2015-04-18 02:47:11	34	QVPSFPUJOCKZAR	SEO	Chrome	M	39	732758368.80	0	...	02:47:11
1	2015-06-07 20:39:50	2015-06-08 01:38:54	16	EOGFQPIZPYXFX	Ads	Chrome	F	53	350311387.87	0	...	01:38:54
2	2015-01-01 18:52:44	2015-01-01 18:52:45	15	YSSKYQSJHPPLJ	Opera	Opera	M	53	2621473820.11	1	...	18:52:45

3 rows x 25 columns

```
In [22]: # dropping signup, purchase, signup_date, signup_time, purchase_date, purchase_time, ip_address columns, device_id
        df = df.drop(['signup', 'purchase', 'signup_date', 'signup_time', 'purchase_date', 'purchase_time', 'ip_address', 'device_id'],
                      axis=1)

        # we can drop the day of the week one if we'd like so then we just have weekday vs weekend

        # checking
        df.head(3)
```

```
Out[22]:
```

purchase_value	source	browser	sex	age	class	ip_country	signup_month	purchase_month	signup_day	purchase_day	signup
0	34	SEO	Chrome	M	39	0	Japan	February	April	Tuesday	Saturday
1	16	Ads	Chrome	F	53	0	United States	June	June	Sunday	Monday
2	15	SEO	Opera	M	53	1	United States	January	January	Thursday	Thursday

Categorical Data Updates

```
In [23]: # 'Invald Country Name: 'Bonisire Sint Eustatius; Saba' from running country conversion code so dropping that
        df.drop(df[df['ip_country'] == 'Bonisire Sint Eustatius; Saba'], index, inplace=True)
        df.drop(df[df['ip_country'] == 'Reunion'], index, inplace=True)
        df.shape

        (151108, 16)
```

```
Out[23]:
```

```
In [24]: # making a copy of dataframe in case it goes weird
        df2 = df.copy()
```

```
In [25]: # changing formatting for troublesome names
        df2.loc[df2['ip_country'] == 'Korea Republic of', 'ip_country'] = 'Korea, Republic of'
        df2.loc[df2['ip_country'] == 'Taiwan; Republic of China (ROC)', 'ip_country'] = 'Taiwan'
        df2.loc[df2['ip_country'] == 'Iran (Islamic Republic of)', 'ip_country'] = 'Iran'
        df2.loc[df2['ip_country'] == 'Moldova Republic of', 'ip_country'] = 'Moldova, Republic of'
        df2.loc[df2['ip_country'] == 'Croatia (LOCAL Name: Hrvatska)', 'ip_country'] = 'Croatia'
        df2.loc[df2['ip_country'] == 'Slovakia (SLOVAK Republic)', 'ip_country'] = 'Slovakia'
        df2.loc[df2['ip_country'] == 'Cote D'Ivoire', 'ip_country'] = 'Ivory Coast'
        df2.loc[df2['ip_country'] == 'Virgin Islands (U.S.)', 'ip_country'] = 'United States'
        df2.loc[df2['ip_country'] == 'Tanzania United Republic of', 'ip_country'] = 'Tanzania, United Republic of'
        df2.loc[df2['ip_country'] == 'Congo The Democratic Republic of The', 'ip_country'] = 'Congo'
```

```
In [26]: import pycountry_convert as pc

        def country_to_continent(country_name):
            if country_name == 'VPNUser':
                return 'VPNUser'
            if country_name == 'European Union':
                return 'Europe'
            if country_name == 'Libyan Arab Jamahiriya':
                return 'Africa'
            if country_name == 'Bosnia and Herzegovina':
                return 'Europe'
            if country_name == 'Palestinian Territory Occupied':
                return 'Africa'
            if country_name == 'Curaçao':
                return 'South America'
            if country_name == 'Bonisire Sint Eustatius; Saba':
                return 'North America'
            if country_name == 'Reunion':
                return 'Africa'

            country_alpha2 = pc.country_name_to_country_alpha2(country_name)
            country_code = pc.country_alpha2_to_continent_code(country_alpha2)
            country_continent_code = pc.convert_continent_code_to_continent_name(country_continent_code)

            return country_continent_name

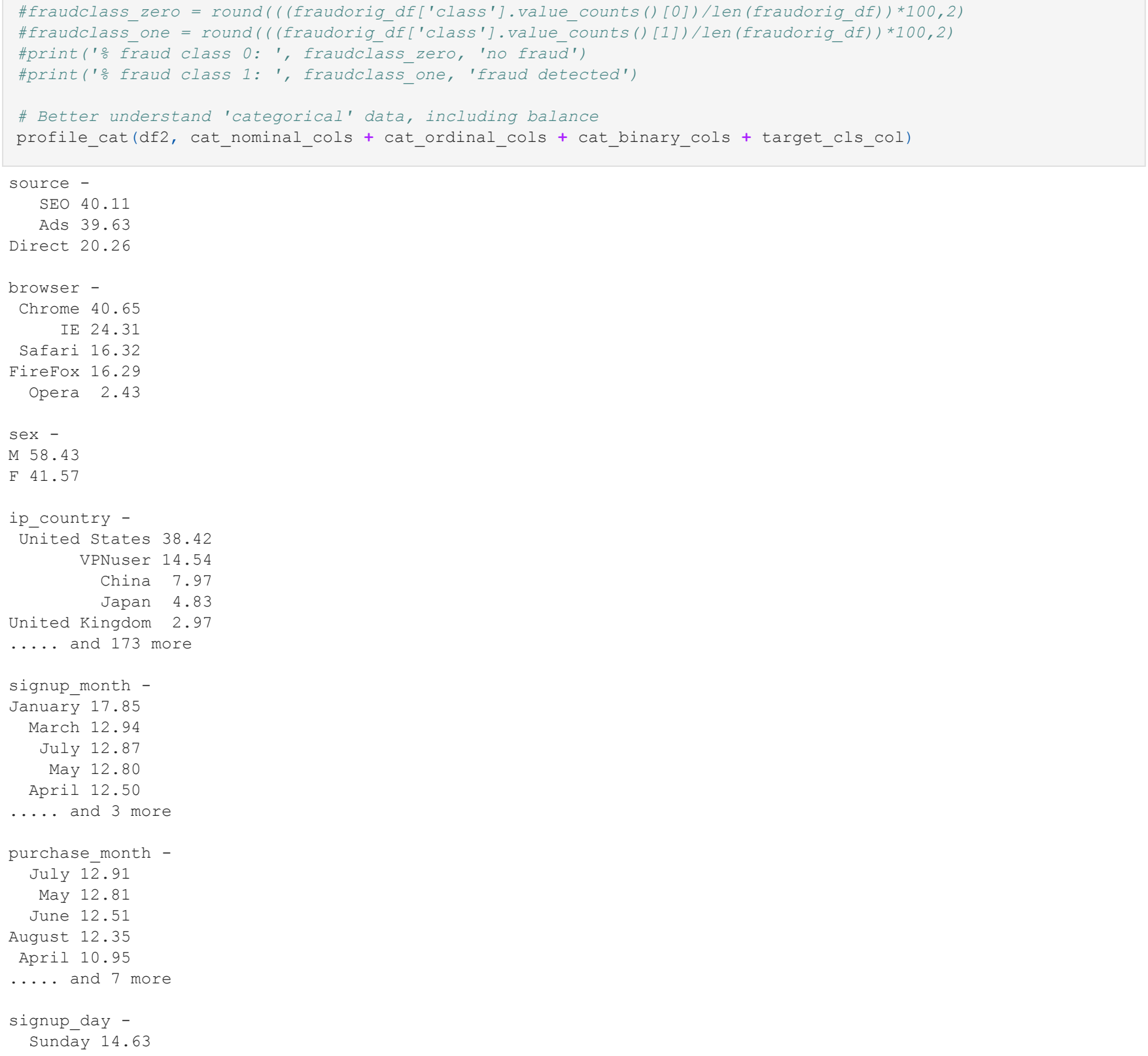
        # calling the function
        df2['continent'] = df2['ip_country'].apply(lambda x: country_to_continent(x))
```

```
In [27]: # checking for unique continents (should be no surprises)
        df2['continent'].unique()
```

```
Out[27]:
```

```
array(['Asia', 'North America', 'VPNUser', 'South America', 'Europe',
       'Africa', 'Oceania'], dtype=object)
```

```
In [28]: plt.hist(df2['continent'])
        plt.show()
```



```
In [29]: # NOTE: DEFERRING HERE SINCE CATEGORICAL ENCODING HANDLED IN MODELPROCESS PIPELINE

        # marking categorical columns (browser, source, sex, signup_month, purchase_month,
        # signup_day, purchase_date, signup_weektype, purchase_time, 'purchase_date',
        # 'purchase_time', 'time_diff_day', 'ip_address', 'device_id'),
        # axis=1)

        # we can drop the day of the week one if we'd like so then we just have weekday vs weekend

        # checking
        df2.head(3)
```

```
Out[29]:
```

purchase_value	source	browser	sex	age	class	ip_country	signup_month	purchase_month	signup_day	purchase_day	signup
0	34	SEO	Chrome	M	39	0	Japan	February	April	Tuesday	Saturday
1	16	Ads	Chrome	F	53	0	United States	June	June	Sunday	Monday
2	15	SEO	Opera	M	53	1	United States	January	January	Thursday	Thursday
3	44	SEO	Safari	M	41	0	VPNUser	April	May	Tuesday	Monday
4	39	Ads	Safari	M	45	0	United States	July	September	Tuesday	Wednesday
...
15107	43	SEO	Chrome	M	28	1	United States	January	March	Tuesday	Sunday
15108	35	SEO	Safari	M	32	0	Netherlands	May	May	Friday	Tuesday
15109	40	SEO	IE	F	26	0	Japan	March	May	Friday	Wednesday
15110	46	SEO	Chrome	M	37	0	United States	July	September	Thursday	Monday
15111	20	Direct	IE	M	38	0	VPNUser	June	July	Wednesday	Tuesday

151108 rows x 17 columns

Feature Updates

```
In [33]: cat_nominal_cols = cat_nominal_cols + \
        ['ip_country', 'signup_month', 'purchase_month', 'signup_day',
        'purchase_day', 'signup_weektype', 'purchase_weektype',
        'signup_timegroup', 'purchase_timegroup', 'time_diff', 'continent']
```

Categorical Data Profiling

```
In [34]: #fraudclass zero = round(((fraudorig_df['class']).value_counts() / len(fraudorig_df)) * 100,2)
        #fraudclass one = round(((fraudorig_df['class']).value_counts() / len(fraudorig_df)) * 100,2)
        #fraudclass two = round(((fraudorig_df['class']).value_counts() / len(fraudorig_df)) * 100,2)
        #fraudclass three = round(((fraudorig_df['class']).value_counts() / len(fraudorig_df)) * 100,2)

        # better understand 'categorical' data, including balance
        profile_cat(df2, cat_nominal_cols + cat_ordinal_cols + target_cls_col)

        source
        20.11
        Ads 39.63
        Direct 20.26

        browser -
        Chrome 40.65
        IE 24.31
        Safari 16.32
        Firefox 16.29
        Opera 2.43

        sex:
        M 58.43
        F 41.57

        ip_country =
        United States 38.42
        VPNUser 14.54
        Asia 20.55
        Japan 4.83
        United Kingdom 2.97
        ... and 15 more

        signup_month -
        Lower 69.70
        March 12.94
        July 12.87
        May 12.80
        June 12.51
        August 12.35
        April 12.28
        ... and 7 more

        purchase_month -
        July 12.91
        May 12.80
        June 12.51
        August 12.35
        April 12.28
        ... and 7 more

        signup_weektype =
        PM 50.06
        AM 49.94

        purchase_timegroup =
        PM 50.21
        AM 49.79

        time_diff -
        Lower 69.70
        weeks 16.74
        day_of 5.80
        days 1.91
        month 2.31
        ... and 1 more

        continent =
        North America 41.36
        Asia 20.55
        Europe 17.22
        VPNUser 14.54
        South America 3.50
        ... and 2 more

        class =
        0 90.64
        1 9.36
```

Data Partitioning

```
In [35]: # Create predictor and target dataframes
        fraud_x = df2[['browser', 'source', 'sex', 'continent', 'signup_month', 'purchase_month',
        'signup_day', 'purchase_day', 'signup_weektype', 'purchase_weektype',
        'signup_timegroup', 'purchase_timegroup', 'time_diff']]
        fraud_y = df2[target_cls_col]
```

```
In [36]: # Split data and confirm proportions
        train_ratio = 0
```



```
In [38]: # Simple function to remove column-level outliers (note this needs expanding
# to row-level and/or conversion, for best results)
def remove_outliers(df_col):
    q1 = df_col.quantile(0.25)
    q3 = df_col.quantile(0.75)
    iqr = q3 - q1
    ubound = q1 + 1.5 * iqr
    lbound = q3 - 1.5 * iqr
    df_out = df_col[(df_col >= lbound) & (df_col <= ubound)]
    removed = len(df_col) - len(df_out)
    if removed > 0:
        print(df_col.name, 'outliers removed: ', removed)
    return df_out

# Per note above, holding this out here but example:
# for c in num_ratio_cols:
#     fraud_df[c] = remove_outliers(fraud_df[c])
```

Correlation Analysis

```
In [39]: # Show -relatively strong correlations to consider, using approach similar to
# what is used in profile() function to highlight >= threshold
corr = 0.5
corr_df = pd.DataFrame(
    X_train[num_interval_cols + num_ratio_cols].corr())
# apply lambda x: abs(x) >= CORR_TH
corr_df = round(corr_df, 1).astype(str)
corr_df.replace(['0', '0.0', 'nan', 'False'], '', inplace=True)
corr_df
```

```
Out[39]:
```

	age	purchase_value
age	1.0	
purchase_value		1.0

```
In [40]: # To consider for categorical value correlation?

def cramers_corrected_stat(confusion_matrix):
    """ calculate Cramers V statistic for categorical-categorical association.
    uses correction from Bergsma and Wicher,
    Journal of the Korean Statistical Society 42 (2013): 323-328
    """
    chi2 = ss.chi2_contingency(confusion_matrix)[0]
    n = confusion_matrix.sum().sum()
    phi2 = chi2/n
    r,k = confusion_matrix.shape
    phi2corr = max(0, phi2 - ((k-1)*(r-1))/(n-1))
    rcorr = r - ((r-1)**2)/(n-1)
    kcorr = k - ((k-1)**2)/(n-1)
    return np.sqrt(phi2corr / min((kcorr-1), (rcorr-1)))

#cols = bin_cols
#corrM = np.zeros((len(cols),len(cols)))
# There's probably a nice pandas way to do this
# for col1, col2 in itertools.combinations(cols, 2):
#     idx1, idx2 = cols.index(col1), cols.index(col2)
#     corrM[idx1, idx2] = cramers_corrected_stat(pd.crosstab(building_df[col1], building_df[col2]))
#     corrM[idx2, idx1] = corrM[idx1, idx2]

#corr = pd.DataFrame(corrM, index=cols, columns=cols)
#fig, ax = plt.subplots(figsize=(7, 6))
#ax = sns.heatmap(corr, annot=True, ax=ax); ax.set_title("Cramer V Correlation between Variables");
```

Data Mining (Unsupervised)

Including a placeholder in case we see an opportunity to apply e.g., clustering from coursework...

Modeling

Model Setup (selection)

```
In [53]: # Set feature cols for appropriate pipeline preprocessing
cat_cols = cat_nominal_cols + cat_ordinal_cols + cat_binary_cols # one-hot encoding, imputing (if necc)
num_cols = num_interval_cols + num_ratio_cols # scaling, imputing (if necc)

# Set model list
mp_queue = [
    (LogisticRegression(), {'random_state': 42}),
    (Perceptron(), {'class_weight': 'balanced'}),
    (LinearDiscriminantAnalysis(), None),
    (LinearSVC(), {'max_iter': 500}),

    (KNeighborsClassifier(), {'n_neighbors': 3}),
    (KNeighborsClassifier(), {'n_neighbors': 5}),
    (KNeighborsClassifier(), {'n_neighbors': 7}),

    (DecisionTreeClassifier(), {'max_depth': 4, 'random_state': 42}),
    (DecisionTreeClassifier(), {'max_depth': 5, 'random_state': 42}),

    (RandomForestClassifier(), {'max_depth': 4, 'random_state': 42}),
    (RandomForestClassifier(), {'max_depth': 5, 'random_state': 42}),
    (AdaBoostClassifier(), {'n_estimators': 10, 'random_state': 42}),
    (MLPClassifier(), {'random_state': 42}),
]

# Iterate models (note use of 'copy' is to preserve mutable elements
# of model_queue tuple for possible later use)
mp_df = pd.DataFrame(mp_queue, columns=['algorithm', 'params'])
mp_df['mp'] = mp_df.apply(
    lambda mp: ModelProcess(copy.deepcopy(mp['algorithm']), None,
                             copy.copy(mp['params']),
                             X_train, y_train,
                             X_val, y_val,
                             X_test, y_test,
                             cat_cols, num_cols).train_validate_test(), axis=1)

# Compile, sort, and display results
mp_df[['train_acc', 'train_f1', 'train_time',
       'val_acc', 'val_f1', 'val_time',
       'test_acc', 'test_f1', 'test_time']] = \
    mp_df['mp'].apply(lambda sumf: sumf)
mp_df.sort_values(by=['train_acc', 'val_acc', 'test_acc'],
                  ascending=[False, False, False], inplace=True)
mp_df.loc[:, mp_df.columns != 'mp']
```

```
Out[54]:
```

	algorithm	params	train_acc	train_f1	train_time	val_acc	val_f1	val_time	test_acc	test_f1	test_time
9	MLPClassifier()	{'random_state': 42}	0.96	0.85	150.90	0.95	0.82	0.34	0.95	0.84	0.17
5	DecisionTreeClassifier()	{'max_depth': 5, 'random_state': 42}	0.96	0.84	1.58	0.95	0.83	0.18	0.96	0.84	0.09
4	DecisionTreeClassifier()	{'max_depth': 4, 'random_state': 42}	0.96	0.84	1.44	0.95	0.83	0.18	0.96	0.84	0.09
6	RandomForestClassifier()	{'max_depth': 4, 'random_state': 42}	0.96	0.84	5.30	0.95	0.83	0.47	0.96	0.84	0.22
7	RandomForestClassifier()	{'max_depth': 5, 'random_state': 42}	0.96	0.84	5.88	0.95	0.83	0.50	0.96	0.84	0.24
0	LogisticRegression()	{'random_state': 42}	0.95	0.83	2.22	0.95	0.82	0.18	0.96	0.84	0.09
3	LinearSVC()	{'max_iter': 500}	0.95	0.82	7.11	0.95	0.81	0.09	0.95	0.83	0.04
8	AdaBoostClassifier()	{'n_estimators': 10, 'random_state': 42}	0.95	0.82	3.24	0.95	0.81	0.29	0.95	0.83	0.15
2	LinearDiscriminantAnalysis()	None	0.95	0.82	4.51	0.95	0.81	0.17	0.95	0.83	0.09
1	Perceptron()	{'class_weight': 'balanced'}	0.95	0.81	1.23	0.94	0.81	0.09	0.95	0.81	0.04

```
In [ ]: #mp_df.loc[3]['mp'].confusion_matrix('val')
#mp_df.loc[3]['mp'].summary('val')
```

```
In [ ]:
```