


```
In [39]: # Show relatively strong correlations to consider, using approach similar to
# what is used in profile() function to highlight >= threshold
CORR_TH = 0.5
corr_df = pd.DataFrame()
X_train[num_interval_cols + num_ratio_cols].corr() \
    .apply(lambda x: abs(x*(x==CORR_TH)))
corr_df_round(corr_df, 1).astype(str)
corr_df.replace('1.0', 'g,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 Bergma 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
    phi1corr = 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)

This is a placeholder for future opportunity, e.g., pre-supervised model clustering to support feature selection.

Modeling

Model Setup (selection)

```
In [41]: # 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_intervl_cols + num_ratio_cols # scaling, imputing (if necc)

# Set model list
mp_queue = [
    (LogisticRegression(), ('random_state': 42)),
    (Perceptron(), ('class_weight': 'balanced')),
    (LinearDiscriminantAnalysis(), ('solver': 'lbfgs', 'tol': 1e-5)),
    (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', 'param'])
mp_df['mp'] = mp_df.apply(
    lambda mp: ModelProcess(copy.deepcopy(mp['algorithm']), None,
        copy.copy(mp['param']),
        X_train, y_train,
        X_val, y_val,
        X_test, y_test,
        None,
        cat_cols, num_cols).train_validate_test(1, 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 mp: sum(list(map(
            lambda dataset: mp.score(dataset) + [mp.time[dataset], ('train', 'val', 'test')], [])).tolist()
        ), axis=1)
mp_df.sort_values(by=['train_acc', 'val_acc', 'test_acc'],
    ascending=[False, False, False], inplace=True)
mp_df.loc[:, mp_df.columns != 'mp'].to_csv('results_table.csv')
```

LogisticRegression: train...	done in 3.65s.
LogisticRegression: val...	done in 0.18s.
LogisticRegression: test...	done in 0.09s.
Perceptron: train...	done in 0.96s.
Perceptron: val...	done in 0.09s.
Perceptron: test...	done in 0.05s.
LinearDiscriminantAnalysis: train...	done in 3.54s.
LinearDiscriminantAnalysis: test...	done in 0.08s.
LinearSVC: train...	done in 9.13s.
LinearSVC: val...	done in 0.08s.
LinearSVC: test...	done in 0.04s.
KNeighborsClassifier: train...	done in 78.73s.
KNeighborsClassifier: test...	done in 22.39s.
KNeighborsClassifier: val...	done in 11.16s.
KNeighborsClassifier: train...	done in 78.46s.
KNeighborsClassifier: test...	done in 22.28s.
KNeighborsClassifier: val...	done in 11.04s.
KNeighborsClassifier: train...	done in 77.02s.
KNeighborsClassifier: test...	done in 21.82s.
KNeighborsClassifier: val...	done in 10.99s.
DecisionTreeClassifier: train...	done in 3.35s.
DecisionTreeClassifier: val...	done in 0.15s.
DecisionTreeClassifier: test...	done in 0.08s.
DecisionTreeClassifier: train...	done in 1.47s.
DecisionTreeClassifier: val...	done in 0.15s.
DecisionTreeClassifier: test...	done in 0.08s.
RandomForestClassifier: train...	done in 5.18s.
RandomForestClassifier: val...	done in 0.45s.
RandomForestClassifier: test...	done in 0.22s.
RandomForestClassifier: train...	done in 5.75s.
RandomForestClassifier: val...	done in 0.47s.
RandomForestClassifier: test...	done in 0.23s.
AdaBoostClassifier: train...	done in 3.13s.
AdaBoostClassifier: val...	done in 0.26s.
AdaBoostClassifier: test...	done in 0.14s.
MLPClassifier: train...	done in 446.35s.
MLPClassifier: val...	done in 0.19s.
MLPClassifier: test...	done in 0.11s.

	algorithm	train_acc	train_f1	train_time	val_acc	val_f1	val_time	test_acc	test_f1	test_time
12	MLPClassifier()	0.96	0.77	446.35	0.94	0.84	0.19	0.95	0.86	0.11
4	KNeighborsClassifier()	0.96	0.71	78.73	0.95	0.65	22.39	0.95	0.67	11.16
8	DecisionTreeClassifier()	0.96	0.70	1.47	0.95	0.68	0.15	0.96	0.71	0.08
7	DecisionTreeClassifier()	0.96	0.70	1.35	0.95	0.68	0.15	0.96	0.71	0.08
10	RandomForestClassifier()	0.96	0.70	5.75	0.95	0.68	0.47	0.96	0.71	0.23
9	RandomForestClassifier()	0.95	0.69	5.18	0.95	0.67	0.45	0.96	0.70	0.22
0	LogisticRegression()	0.95	0.69	3.65	0.95	0.67	0.18	0.95	0.70	0.09
5	KNeighborsClassifier()	0.95	0.69	78.46	0.95	0.66	22.28	0.95	0.69	11.04
6	KNeighborsClassifier()	0.95	0.68	77.02	0.95	0.66	21.82	0.95	0.68	10.99
3	LinearSVC()	0.95	0.67	713	0.95	0.66	0.08	0.95	0.68	0.04
11	AdaBoostClassifier()	0.95	0.67	313	0.95	0.66	0.26	0.95	0.68	0.14
2	LinearDiscriminantAnalysis()	0.95	0.67	3.54	0.95	0.66	0.15	0.95	0.68	0.08
1	Perceptron()	0.47	0.22	0.96	0.46	0.22	0.09	0.48	0.22	0.05

Feature Importance Review for Key (Candidate) Algorithms

```
In [43]: # Show top features for Logistic Regression (this could be 'baked' into the
# ModelProcess() class or similar in future)
features = mp_df.loc[8]['mp'].pipe(['preprocessor'].transformers_[0][1]['onehotencoder']).get_feature_names_out()
fi = pd.concat([pd.DataFrame(features, columns=['feature']),
    pd.DataFrame(num_cols, columns=['feature'])],
    ignore_index=True)

fi['importance'] = model.coef_[0]
fi['abs_importance'] = abs(fi['importance'])
fi.sort_values(by='abs_importance', 'feature',
    ascending=[False, True], inplace=True)
fi[fi['abs_importance'] >= 0.50][['feature', 'importance']]

Out [43]:
      feature  importance
227  time_diff_day_of      2.77
197  purchase_month_January      2.00
36   ip_country_Bulgaria      -1.39
111  ip_country_Malta      1.34
154  ip_country_Slovenia      -1.16
98   ip_country_Latvia      1.14
118  ip_country_Morocco      -1.08
228  time_diff_days      -1.04
12   ip_country_Algeria      0.97
232  time_diff_nan      -0.91
30   ip_country_Bolivia      0.82
16   ip_country_Armenia      0.76
231  time_diff_weeks      -0.69
80   ip_country_Iceland      0.67
92   ip_country_Kazakhstan      0.61
78   ip_country_Hong Kong      0.59
25   ip_country_Belgium      0.59
142  ip_country_Qatar      -0.57
131  ip_country_Oman      -0.57
19   ip_country_Azerbaijan      0.54
95   ip_country_Kuwait      0.53
96   ip_country_Kyrgyzstan      -0.51
143  ip_country_Romania      -0.51
170  ip_country_Turkmenistan      0.51
159  ip_country_Sudan      0.50
```

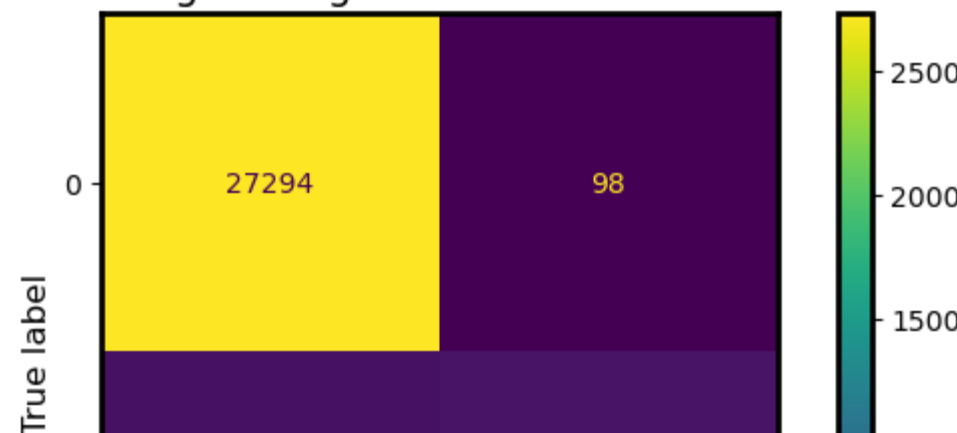
```
In [44]: # Show top features for Decision Tree with max depth of 5
model = mp_df.loc[8]['mp'].model
features = mp_df.loc[8]['mp'].pipe(['preprocessor'].transformers_[0][1]['onehotencoder']).get_feature_names_out()
fi = pd.concat([pd.DataFrame(features, columns=['feature']),
    pd.DataFrame(num_cols, columns=['feature'])],
    ignore_index=True)

fi['importance'] = model.feature_importances_
fi['abs_importance'] = abs(fi['importance'])
fi.sort_values(by='abs_importance', 'feature',
    ascending=[False, True], inplace=True)
fi[fi['abs_importance'] >= 0.05][['feature', 'importance']]

Out [44]:
      feature  importance
227  time_diff_day_of      0.88
197  purchase_month_January      0.12
```

Selected Model Performance Review

```
In [57]: # Show confusion matrix and summary for Logistic Regression
mp_df.loc[0]['mp'].confusion_matrix('val')
mp_df.loc[0]['mp'].summary('val')
```



LogisticRegression - val dataset				
	precision	recall	f1-score	support
0	0.95	1.00	0.97	27392
1	0.94	0.52	0.67	2830
accuracy			0.95	30222
macro avg	0.95	0.76	0.82	30222
weighted avg	0.95	0.95	0.95	30222

Selected Model Gains Review (in business terms)

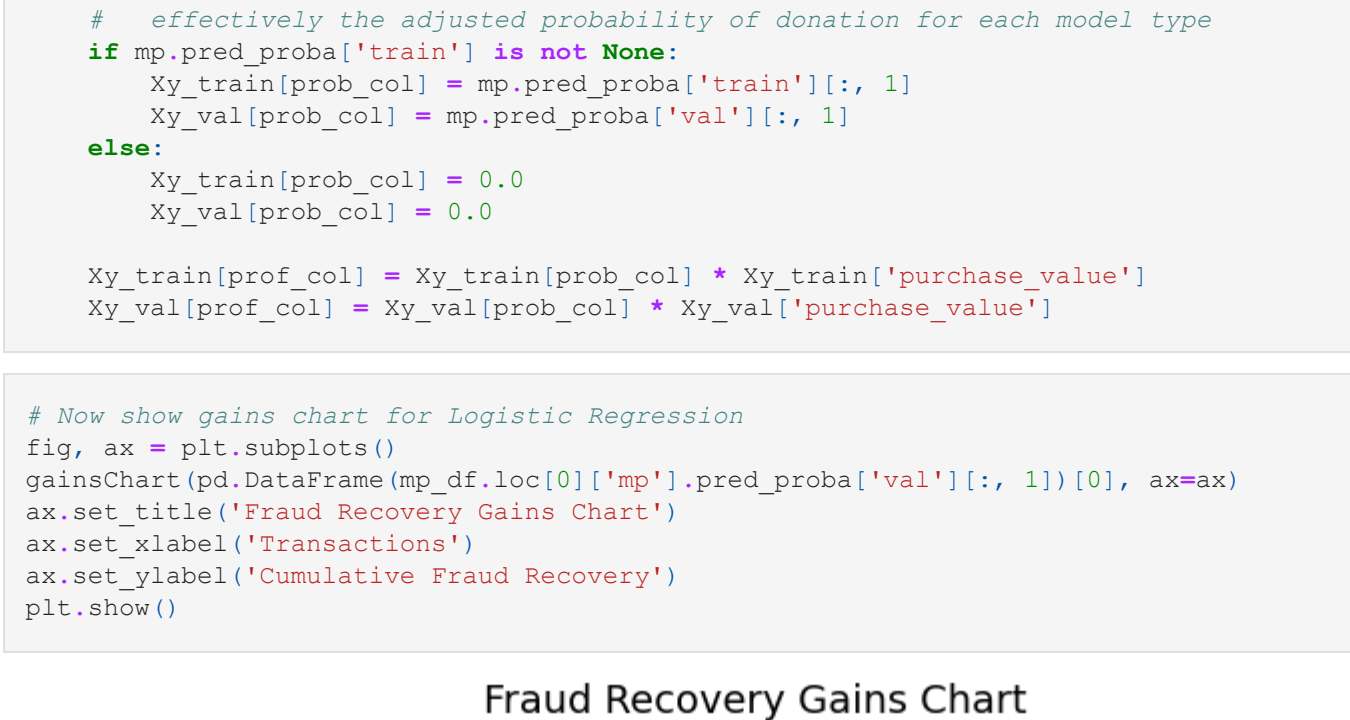
```
In [58]: # Recombine predictors and target
Xy_train = pd.concat([X_train, y_train], axis=1)
Xy_val = pd.concat([X_val, y_val], axis=1)

# Iterate model ('method') results to calc fraud savings
for index, row in mp_df.iterrows():
    mp = row['mp']
    prob_col = str(index)+'_prob'
    prof_col = str(index)+'_prof'

    # Add prediction probabilities for each model ('method') to each row; these are
    # effectively the adjusted probability of donation for each model type
    if mp_pred_proba['train'] is not None:
        Xy_train[prob_col] = mp_pred_proba['train'][:, 1]
        Xy_val[prob_col] = mp_pred_proba['val'][:, 1]
    else:
        Xy_train[prob_col] = 0.0
        Xy_val[prob_col] = 0.0

    Xy_train[prof_col] = Xy_train[prob_col] * Xy_train['purchase_value']
    Xy_val[prof_col] = Xy_val[prob_col] * Xy_val['purchase_value']

In [59]: # Now show gains chart for Logistic Regression
fig, ax = plt.subplots()
gainsChart(pd.DataFrame(mp_df.loc[0]['mp'].pred_proba['val'][:, 1])[0], ax=ax)
ax.set_title('Fraud Recovery Gains Chart')
ax.set_xlabel('Transactions')
ax.set_ylabel('Cumulative Fraud Recovery')
plt.show()
```



```
In [60]: # Show fraud recovery dollars for Logistic Regression
print("\nFraud recovery estimate: $({0:,.0f})".format(sum(Xy_val['0_prof'])))
print("\nFraud recovery estimate as percentage of total: ({0:.1f})\n".format(sum(Xy_val['0_prof']) / sum(Xy_val['purchase_value']) * 100.0))

Fraud recovery estimate: $103,435

Fraud recovery estimate as percentage of total: 9.3%
```

```
In [ ]:
```