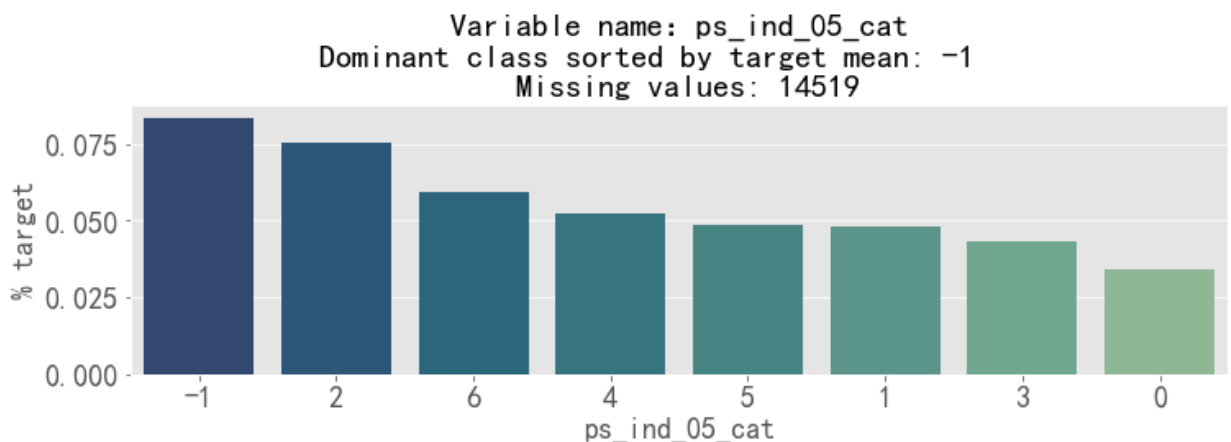
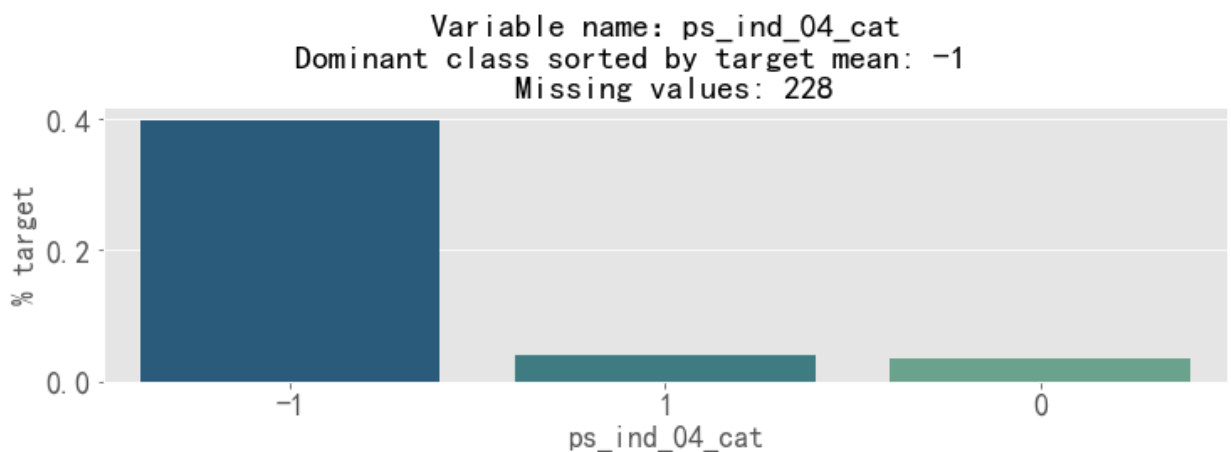
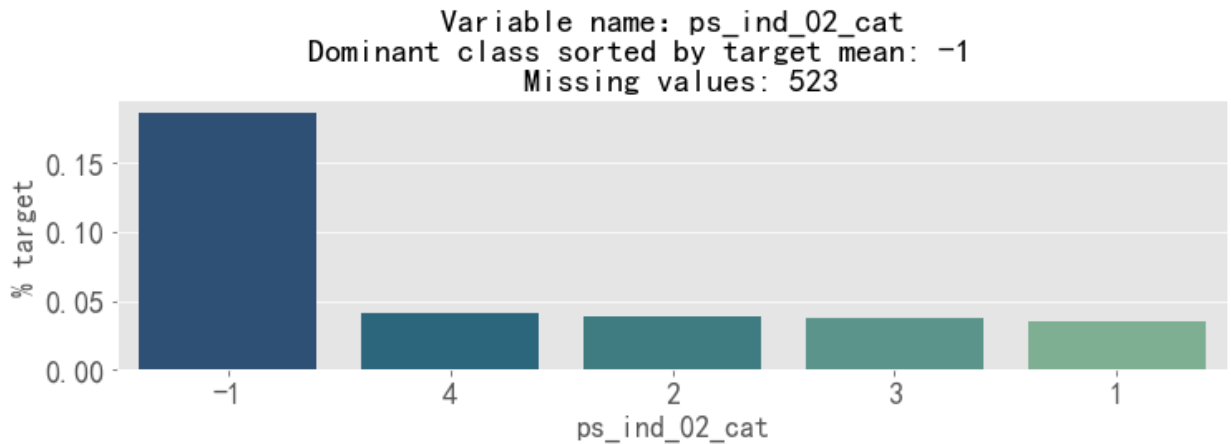


```
In [5]: from data_management import meta
```

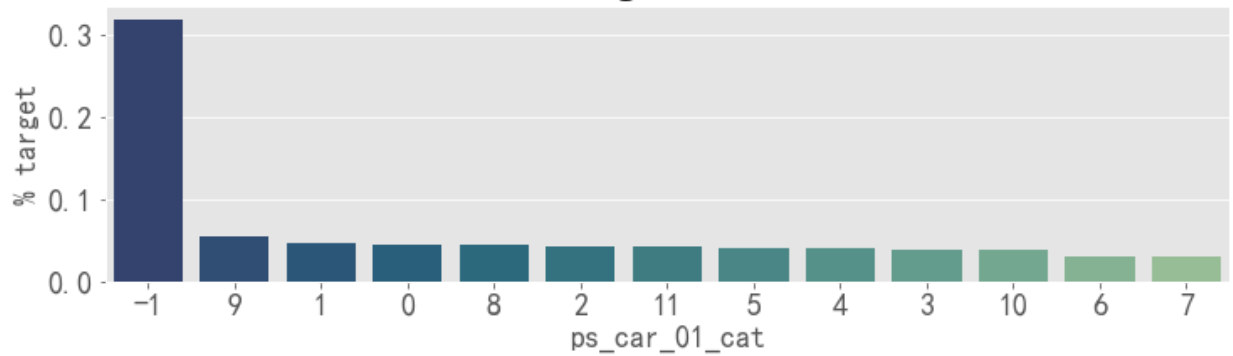
```
In [6]: metadata = meta(train, test)
```

```
In [8]: cat_cols = metadata[(metadata.level == 'nominal') & (metadata.keep)].index
```

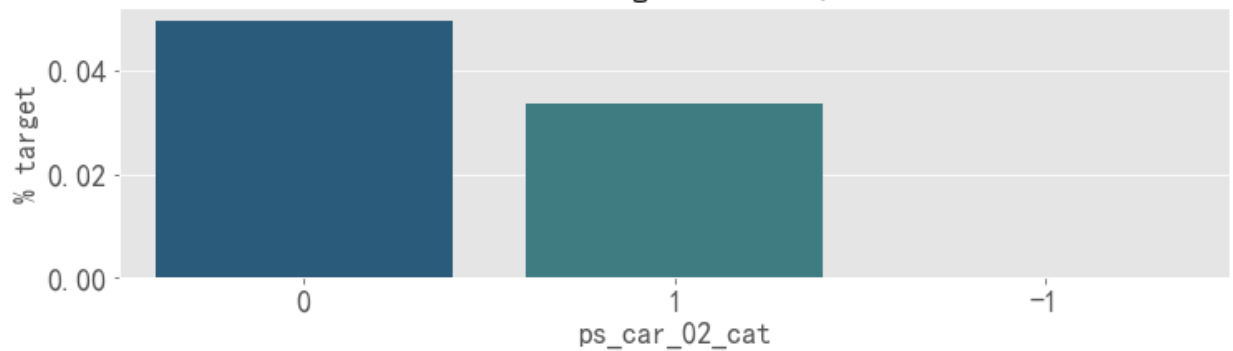
```
In [9]: for col in cat_cols:
fig, ax = plt.subplots(figsize=(12,3))
# Calculate the percentage of target=1 per category value
cat_perc = train[[col, 'target']].groupby([col],as_index=False).mean()
cat_perc.sort_values(by='target', ascending=False, inplace=True)
sns.barplot(ax=ax, x=col, y='target', data=cat_perc, order=cat_perc[col])
plt.ylabel('% target', fontsize=18)
plt.xlabel(col, fontsize=18)
plt.tick_params(axis='both', which='major', labelsize=18)
plt.title(f"Variable name: {col}\n Dominant class sorted by target mean\n\n Missing values: {metadata.loc[col]['missing']}", fontsize = 20)
plt.show();
```



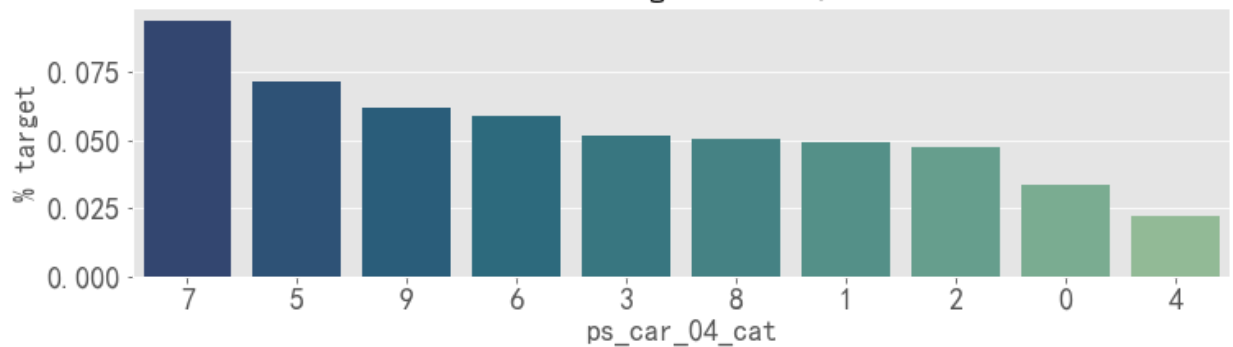
Variable name: ps_car_01_cat
Dominant class sorted by target mean: -1
Missing values: 267



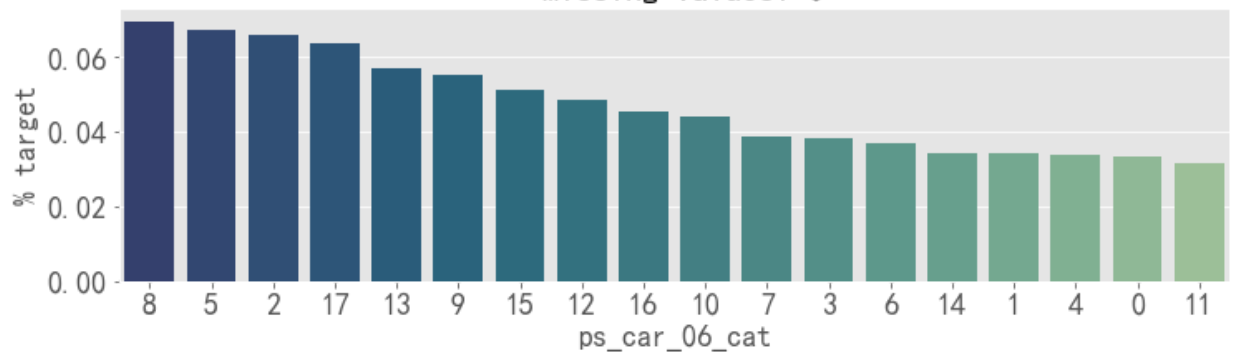
Variable name: ps_car_02_cat
Dominant class sorted by target mean: 0
Missing values: 10

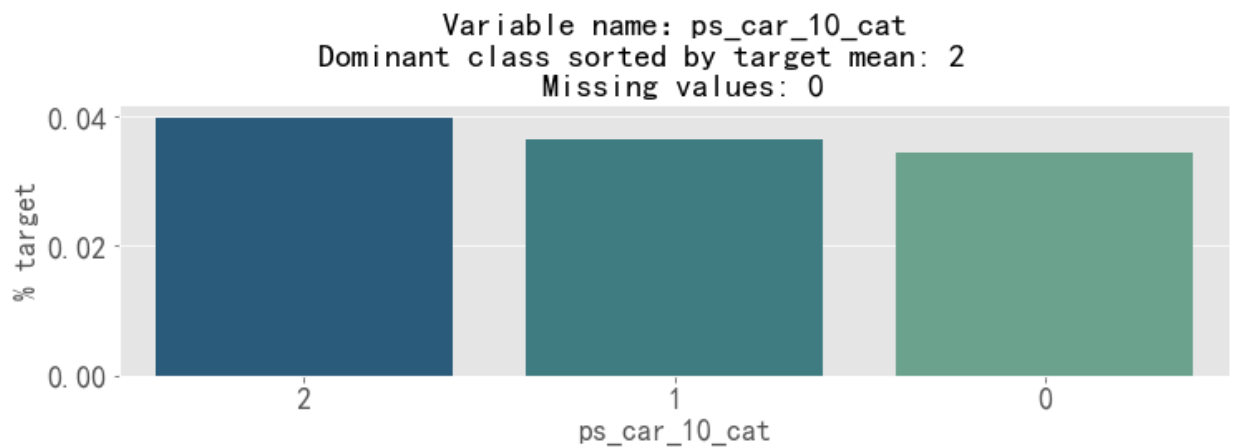
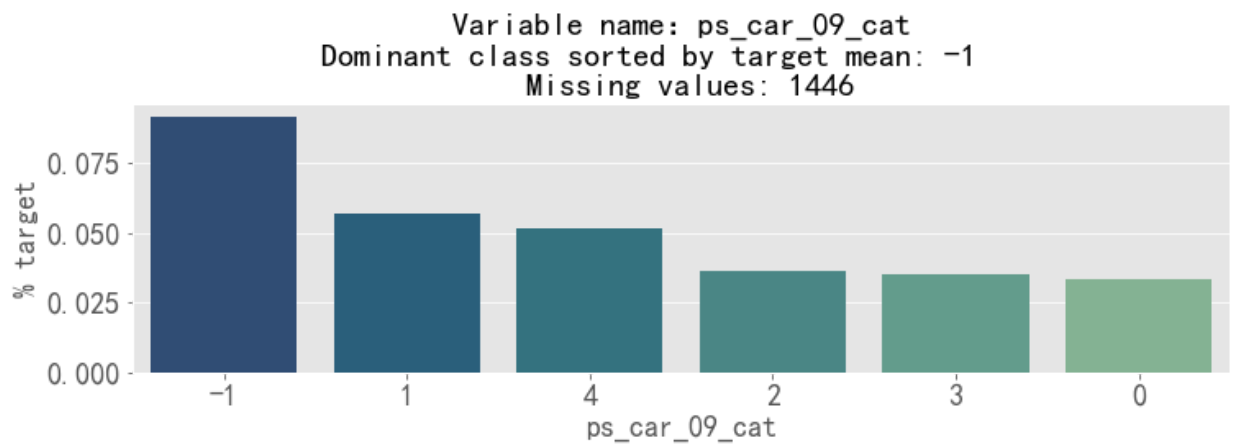
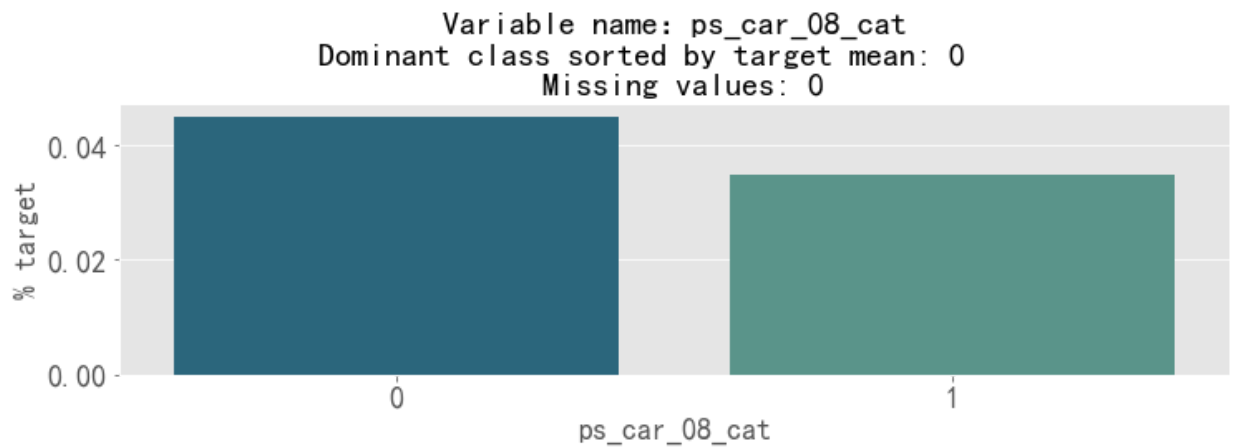
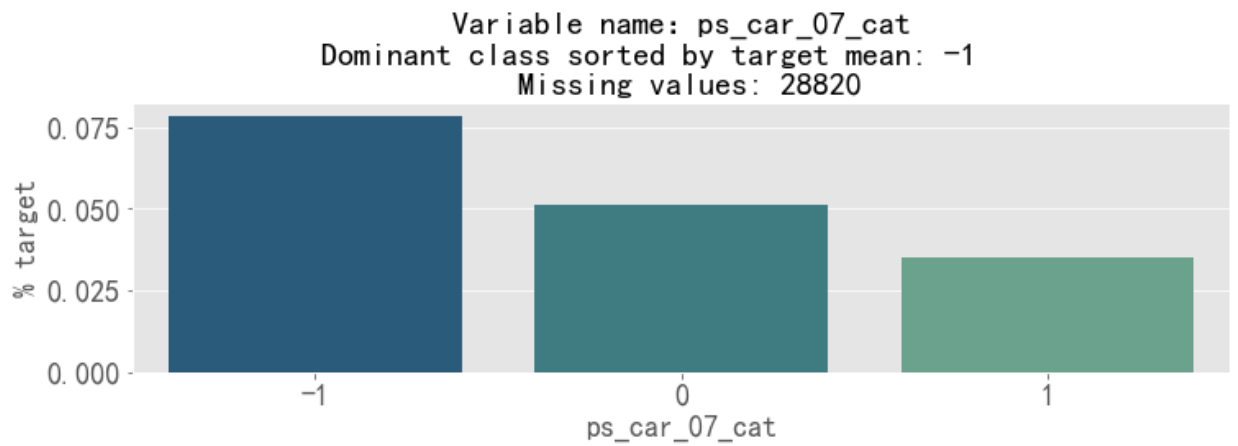


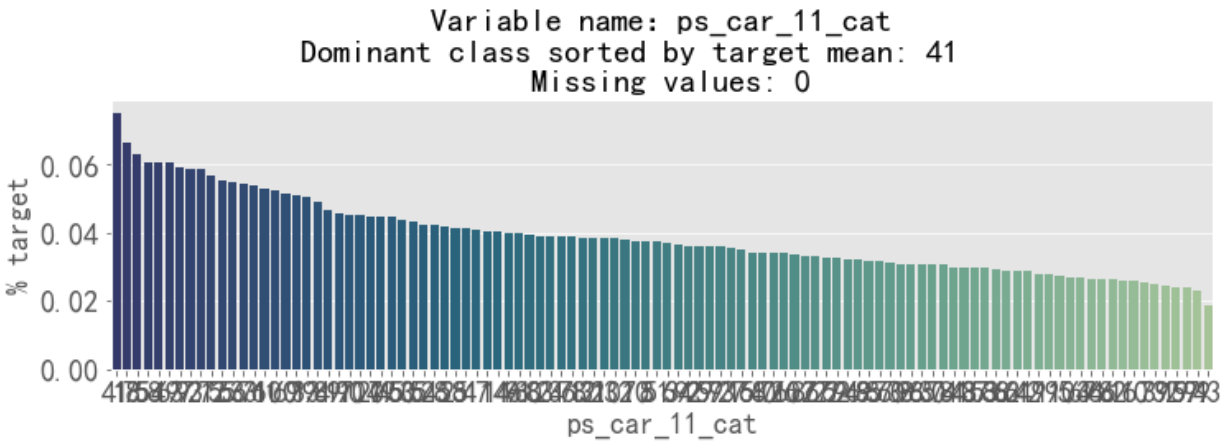
Variable name: ps_car_04_cat
Dominant class sorted by target mean: 7
Missing values: 0



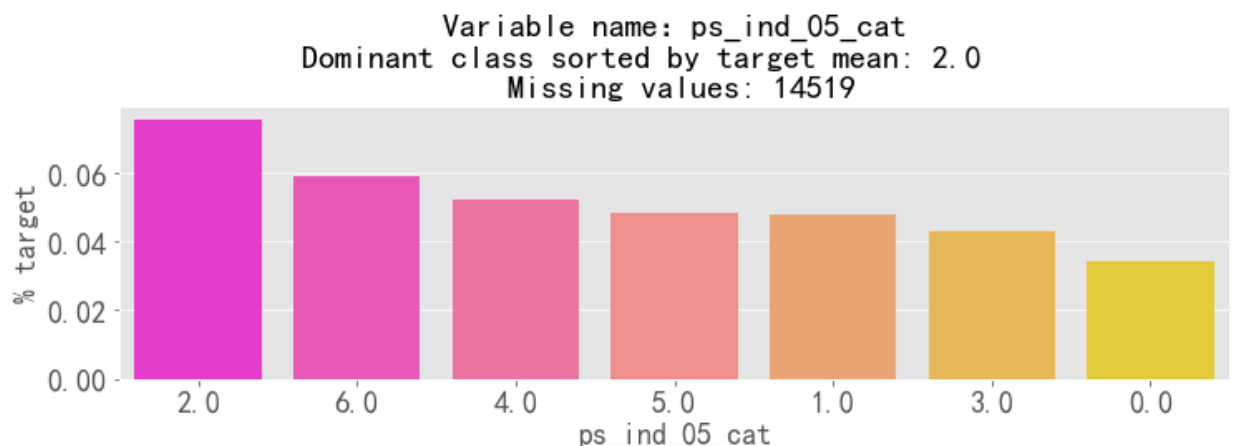
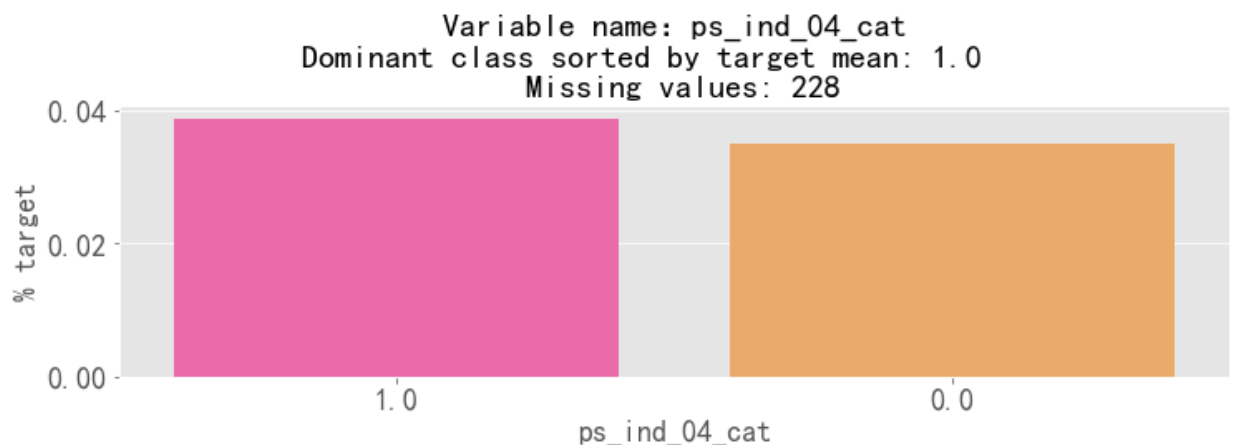
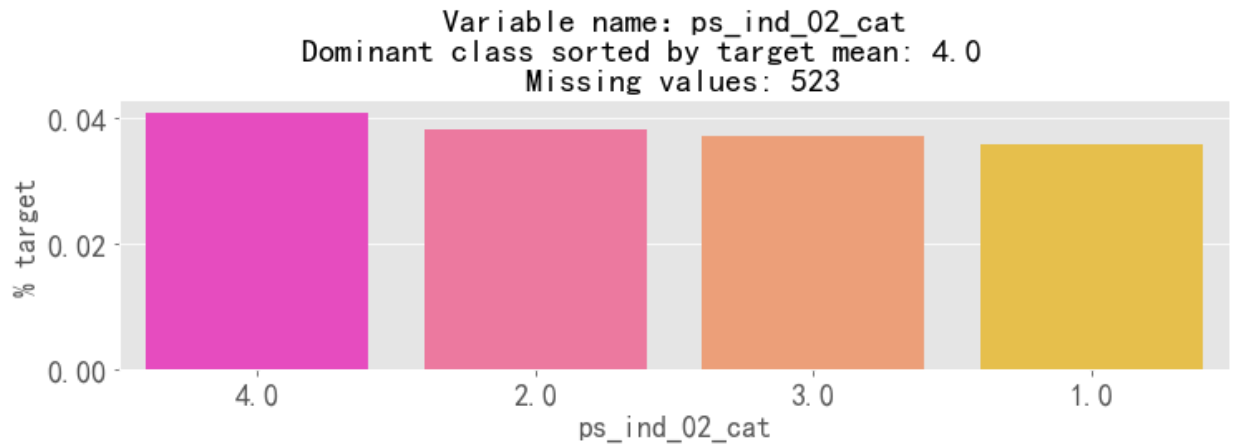
Variable name: ps_car_06_cat
Dominant class sorted by target mean: 8
Missing values: 0



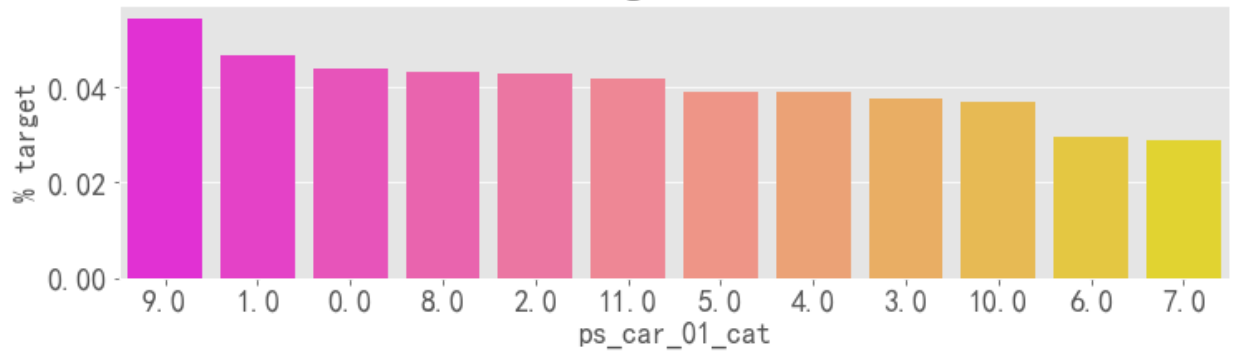




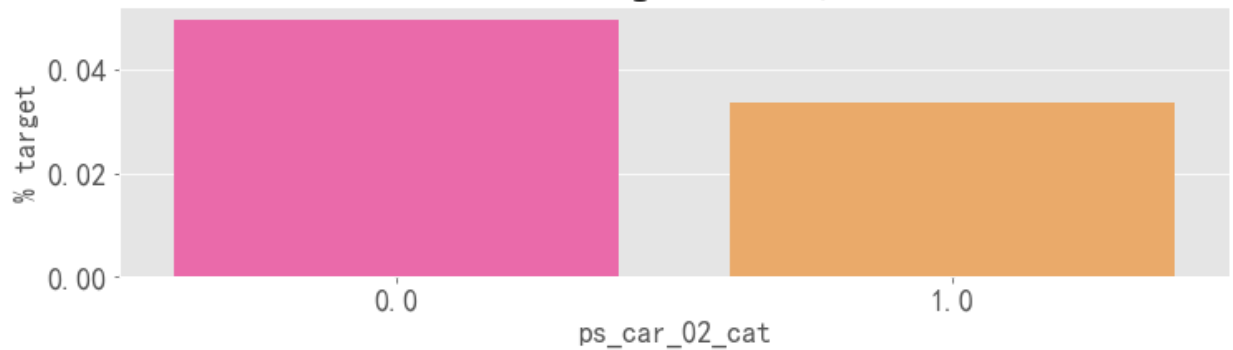
```
In [10]: for col in cat_cols:
fig, ax = plt.subplots(figsize=(12,3))
# Calculate the percentage of target=1 per category value
cat_perc = train_imp[[col, 'target']].groupby([col],as_index=False).mean()
cat_perc.sort_values(by='target', ascending=False, inplace=True)
sns.barplot(ax=ax, x=col, y='target', data=cat_perc, order=cat_perc[col].index)
plt.ylabel('% target', fontsize=18)
plt.xlabel(col, fontsize=18)
plt.tick_params(axis='both', which='major', labelsize=18)
plt.title(f"Variable name: {col}\n Dominant class sorted by target mean\n\n Missing values: {metadata.loc[col]['missing']}", fontsize = 20)
plt.show();
```



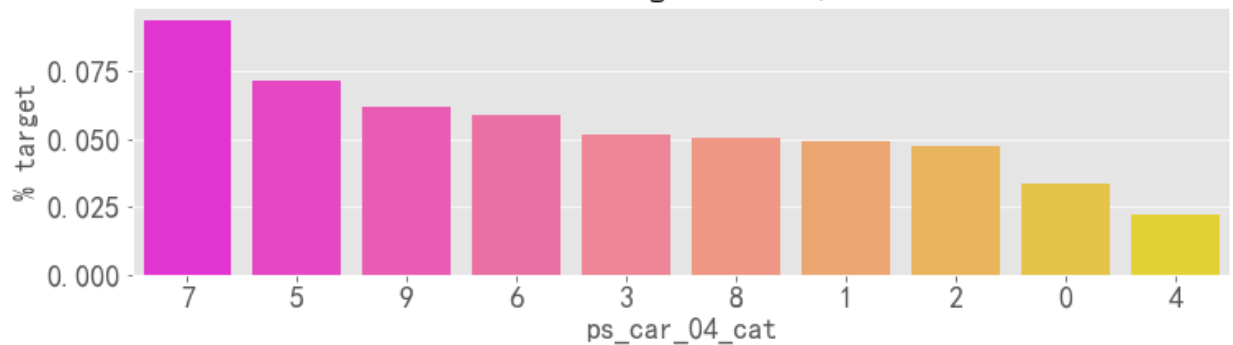
Variable name: ps_car_01_cat
Dominant class sorted by target mean: 9.0
Missing values: 267



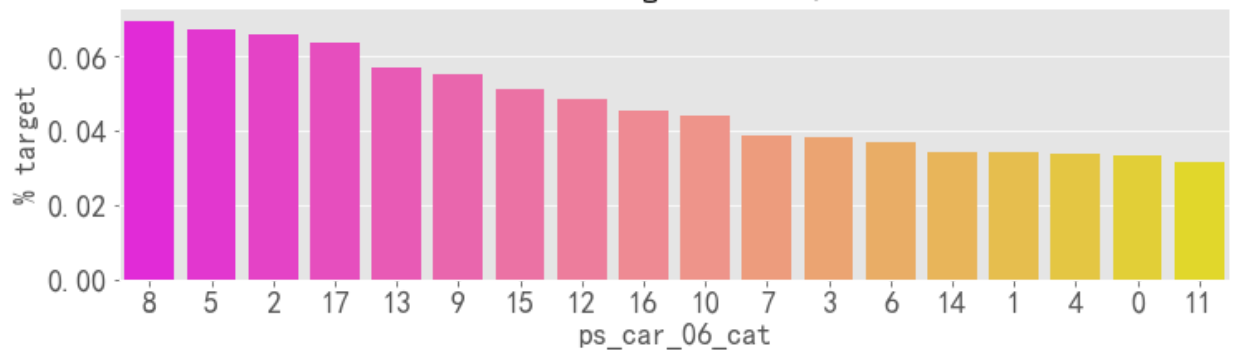
Variable name: ps_car_02_cat
Dominant class sorted by target mean: 0.0
Missing values: 10



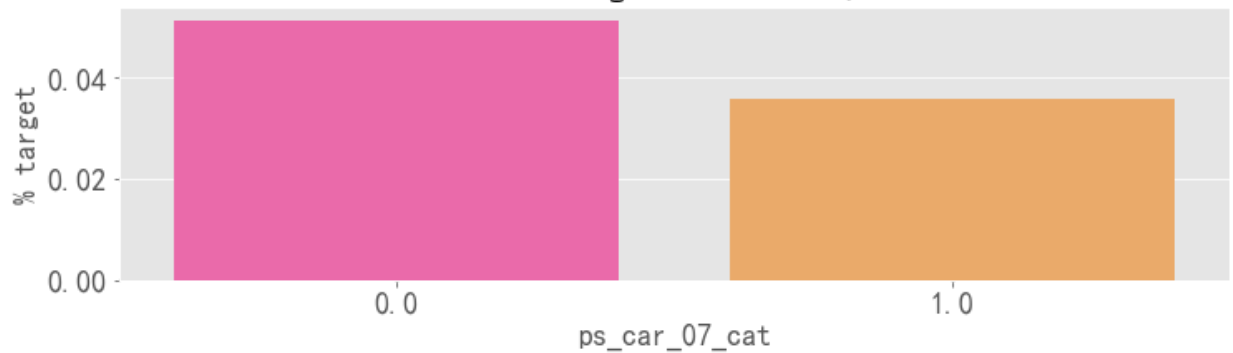
Variable name: ps_car_04_cat
Dominant class sorted by target mean: 7
Missing values: 0



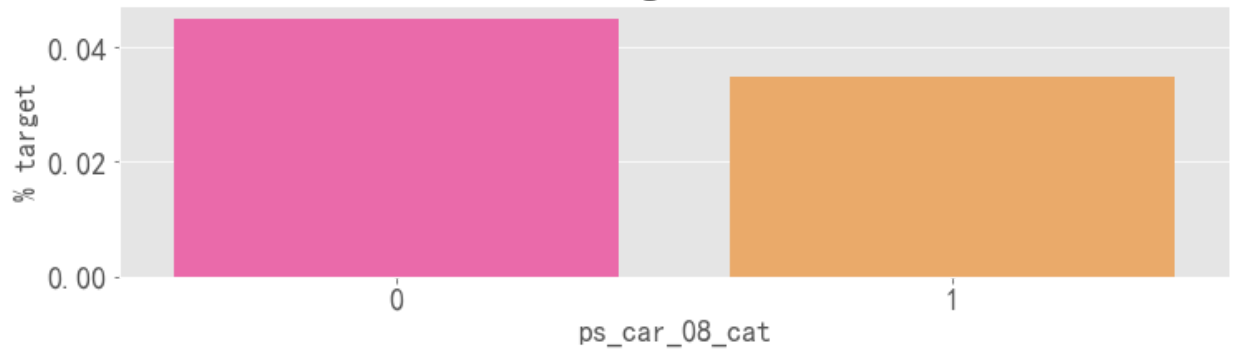
Variable name: ps_car_06_cat
Dominant class sorted by target mean: 8
Missing values: 0



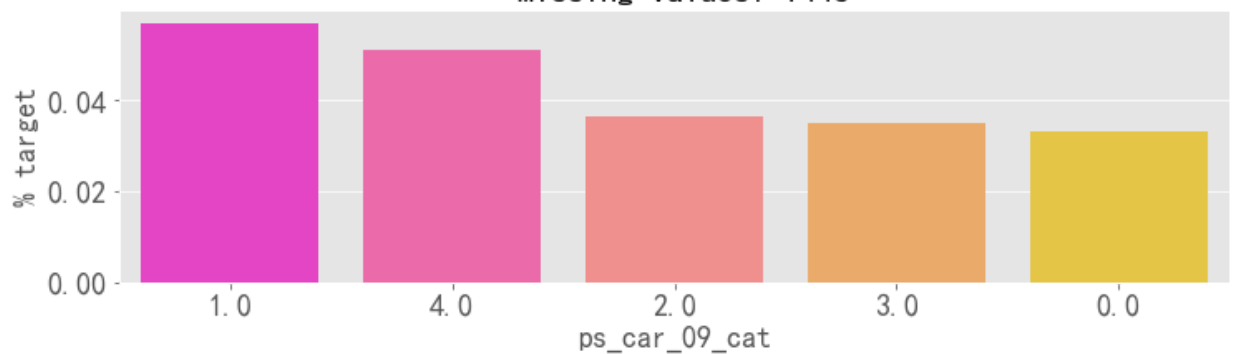
Variable name: ps_car_07_cat
Dominant class sorted by target mean: 0.0
Missing values: 28820



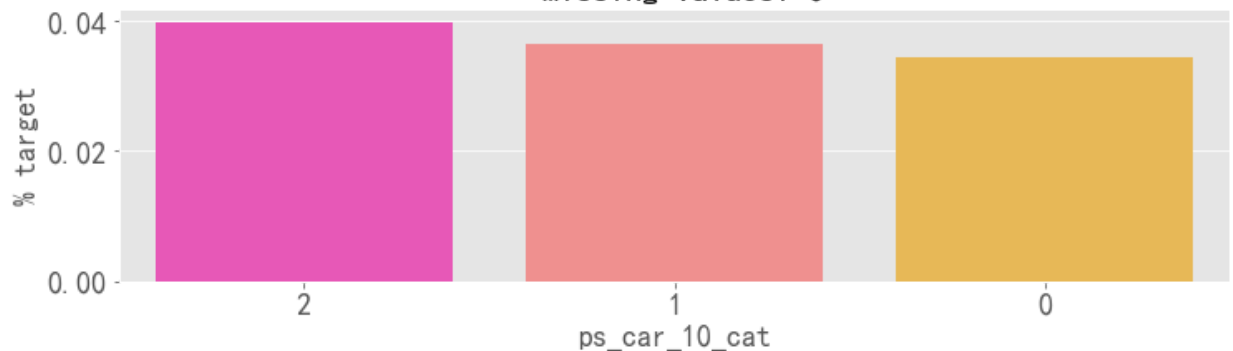
Variable name: ps_car_08_cat
Dominant class sorted by target mean: 0
Missing values: 0

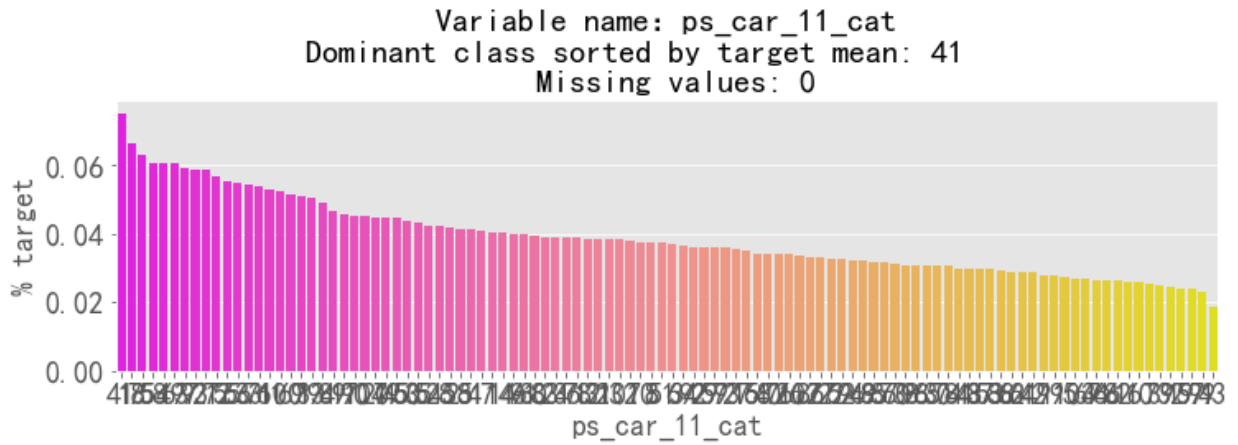


Variable name: ps_car_09_cat
Dominant class sorted by target mean: 1.0
Missing values: 1446



Variable name: ps_car_10_cat
Dominant class sorted by target mean: 2
Missing values: 0

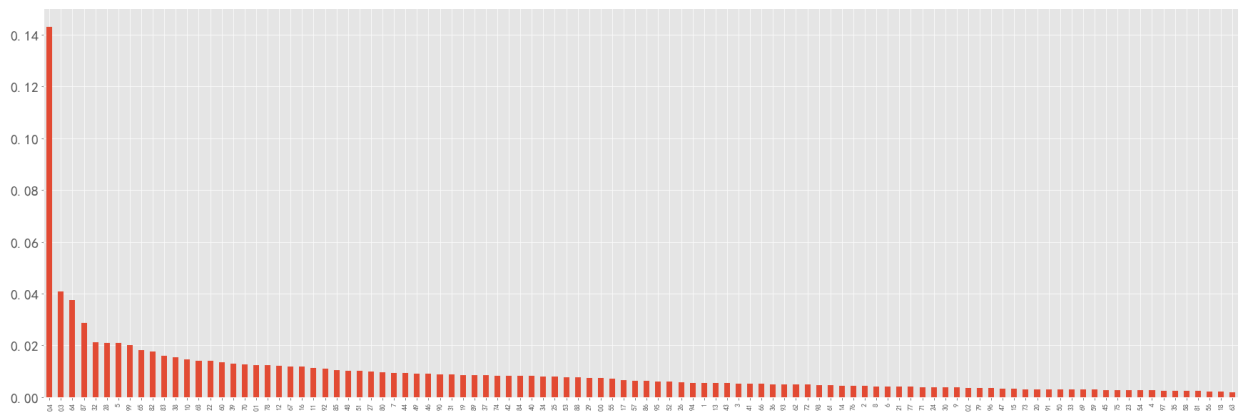




```
In [11]: reverse_trans_cols = [
    'ps_car_07_cat',
    'ps_ind_05_cat',
    'ps_car_09_cat',
    'ps_ind_02_cat',
    'ps_car_01_cat',
    'ps_ind_04_cat',
    ]
```

```
In [12]: train_imp[reverse_trans_cols] = train[reverse_trans_cols]
```

```
In [13]: (train_imp.ps_car_11_cat.value_counts()/train_imp.shape[0]).plot(kind='bar')
plt.tick_params(axis='y', which='major', labelsize=20)
```



Binary

```
In [14]: bin_cols = metadata[(metadata.level == 'binary') & (metadata.keep)].index
```

```
In [15]: zero_list = []
one_list = []
for col in bin_cols:
    zero_list.append((train[col]==0).sum())
    one_list.append((train[col]==1).sum())
```

```
In [16]: trace1 = go.Bar(
    x=bin_cols,
    y=zero_list,
    name='Zero count'
)
trace2 = go.Bar(
    x=bin_cols,
    y=one_list,
    name='One count'
)

data = [trace1, trace2]
layout = go.Layout(
    barmode='stack',
    title='Count of 1 and 0 in binary variables including TARGET'
)

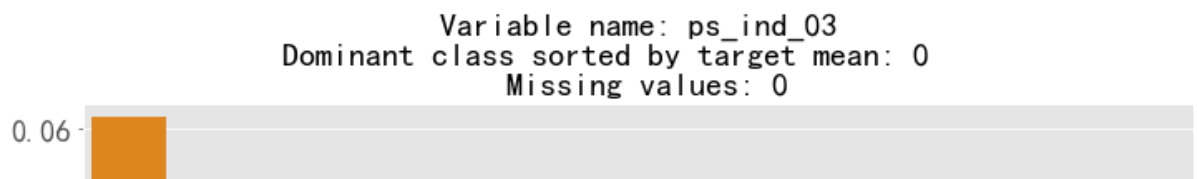
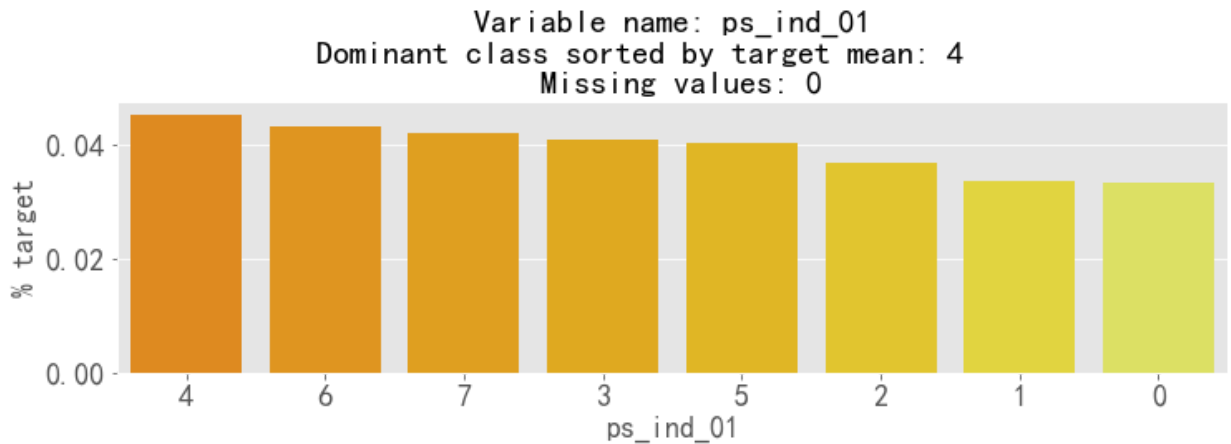
fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='stacked-bar')
```

Here we observe that there are 4 features : ps_ind_10_bin, ps_ind_11_bin, ps_ind_12_bin, ps_ind_13_bin which are completely dominated by zeros. This begs the question of whether these features are useful at all as they do not contain much information about the other class vis-a-vis the target.

Ordinal

```
In [17]: ord_cols = metadata[(metadata.level == 'ordinal') & (metadata.keep)].index
```

```
In [18]: for col in ord_cols:
fig, ax = plt.subplots(figsize=(12,3))
# Calculate the percentage of target=1 per category value
cat_perc = train[[col, 'target']].groupby([col],as_index=False).mean()
cat_perc.sort_values(by='target', ascending=False, inplace=True)
sns.barplot(ax=ax, x=col, y='target', data=cat_perc, order=cat_perc[col])
plt.ylabel('% target', fontsize=18)
plt.xlabel(col, fontsize=18)
plt.tick_params(axis='both', which='major', labels=18)
plt.title(f"Variable name: {col}\n Dominant class sorted by target mean\n\n Missing values: {metadata.loc[col]['missing']}", fontsize = 20)
plt.show();
```



```
In [19]: conti_cols = metadata[(metadata.level.isin(['ratio','interval'])) & (metada
```

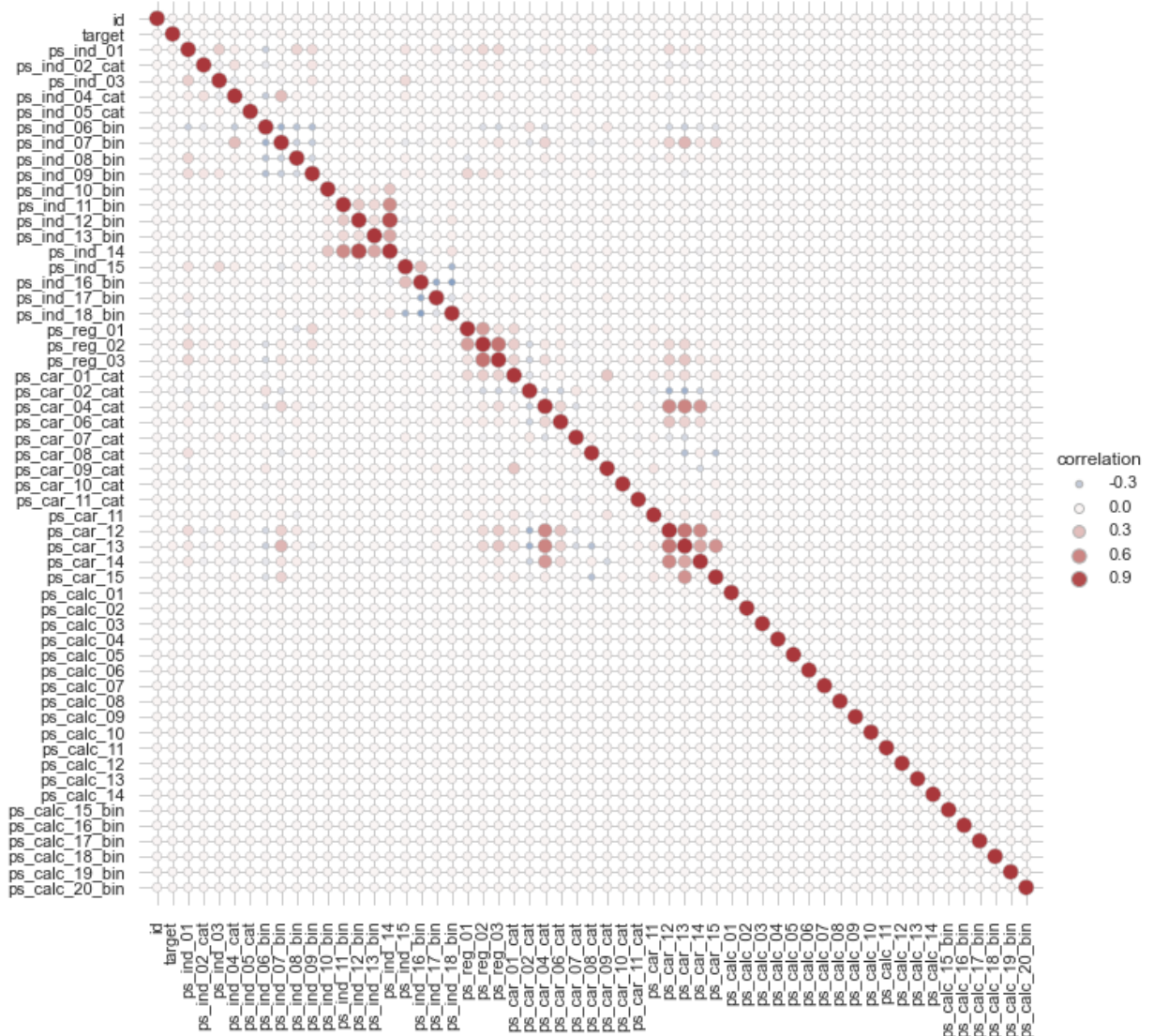
```
In [20]: conti_cols
```

```
Out[20]: Index(['ps_reg_01', 'ps_reg_02', 'ps_reg_03', 'ps_car_12', 'ps_car_13',
               'ps_car_14', 'ps_car_15', 'ps_calc_01', 'ps_calc_02', 'ps_calc_0
               3'],
              dtype='object', name='colname')
```

```
In [21]: sns.set_theme(style="whitegrid")
corr_mat = train_imp.corr().stack().reset_index(name="correlation")

# Draw each cell as a scatter point with varying size and color
g = sns.relplot(
    data=corr_mat,
    x="level_0", y="level_1", hue="correlation", size="correlation",
    palette="vlag", hue_norm=(-1, 1), edgecolor=".7",
    height=10, sizes=(20, 100), size_norm=(-.2, .8),
)

# Tweak the figure to finalize
g.set(xlabel="", ylabel="", aspect="equal")
g.despine(left=True, bottom=True)
g.ax.margins(.02)
for label in g.ax.get_xticklabels():
    label.set_rotation(90)
for artist in g.legend.legendHandles:
    artist.set_edgecolor(".7")
```




```

In [22]: # Generate a mask for the upper triangle

corr_mat = train_imp[conti_cols].corr()
mask = np.triu(np.ones_like(corr_mat, dtype=bool))

# Set up the matplotlib figure
f, ax = plt.subplots(figsize=(11, 9))

# Generate a custom diverging colormap
cmap = sns.diverging_palette(230, -20, as_cmap=True)

# Draw the heatmap with the mask and correct aspect ratio
sns.heatmap(corr_mat, mask=mask, cmap=cmap, vmax=.3, center=0,
            square=True, linewidths=.5, cbar_kws={"shrink": .75}, annot=True

```

Out[22]: <AxesSubplot:>



There are a strong correlations between the variables:

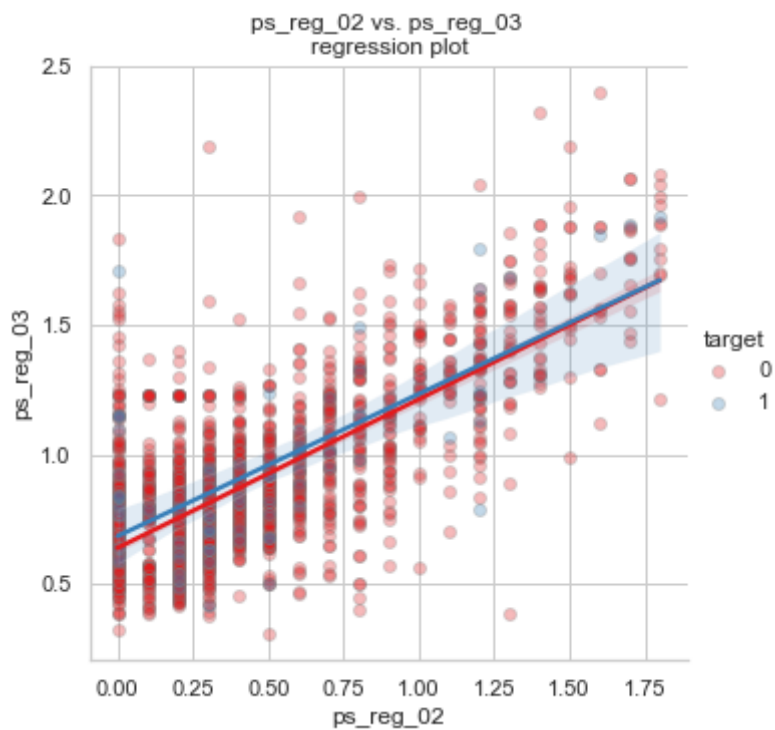
- ps_reg_02 and ps_reg_03 (0.69)
- ps_car_12 and ps_car_13 (0.67)
- ps_car_12 and ps_car_14 (0.58)
- ps_car_13 and ps_car_15 (0.53)

- ps_reg_01 and ps_reg_02 (0.47)

```
In [23]: trn_sample = train_imp.sample(2000)
```

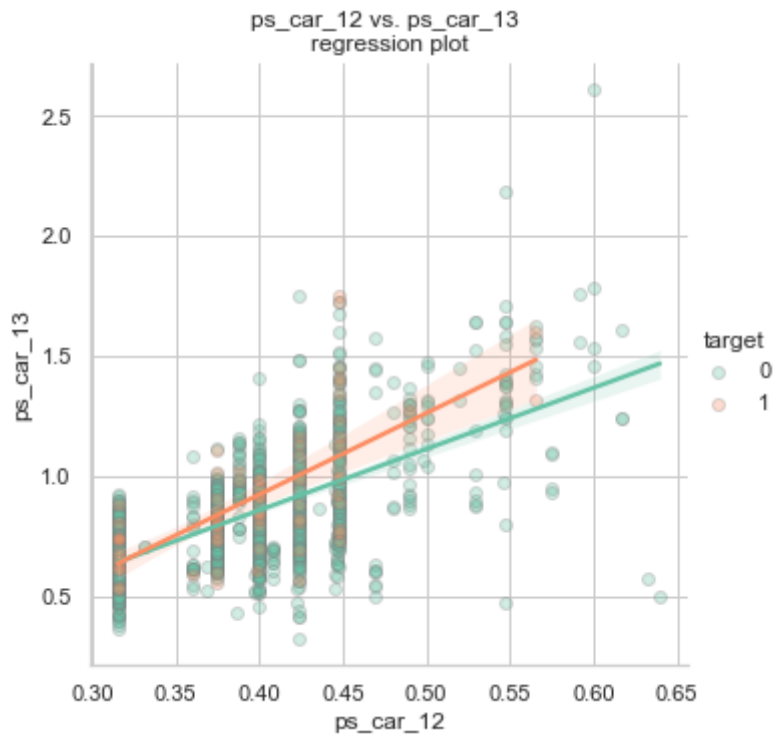
```
In [24]: sns.lmplot(x='ps_reg_02', y='ps_reg_03', data=trn_sample, hue='target',  
                  palette='Set1', height = 5, scatter_kws={'alpha':0.3,'edgecolor'
```

```
Out[24]: <seaborn.axisgrid.FacetGrid at 0x19f00194dc8>
```



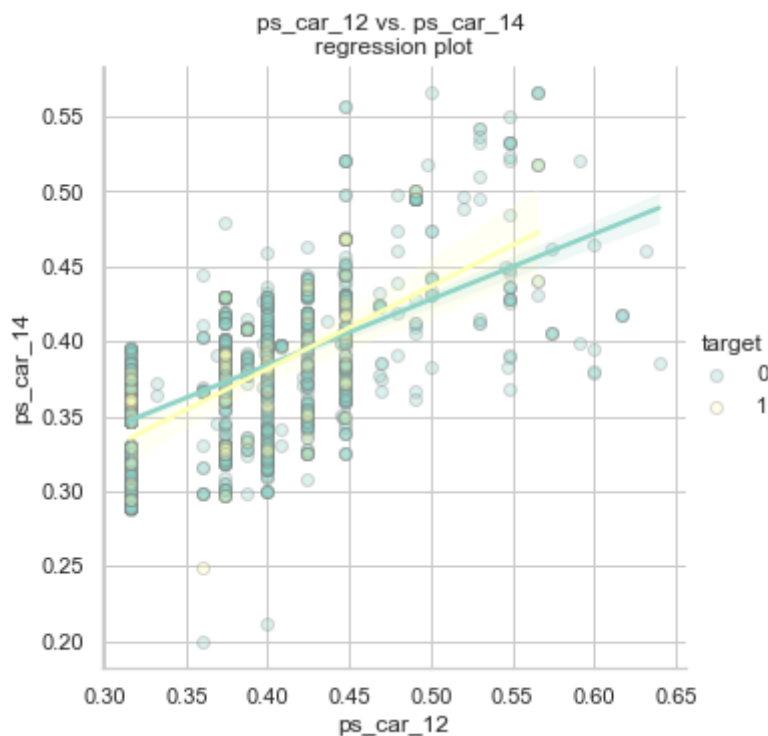
```
In [25]: sns.lmplot(x='ps_car_12', y='ps_car_13', data=trn_sample, hue='target',
                    palette='Set2', height = 5, scatter_kws={'alpha':0.3,'edgecolor'
```

```
Out[25]: <seaborn.axisgrid.FacetGrid at 0x19f0fda9108>
```



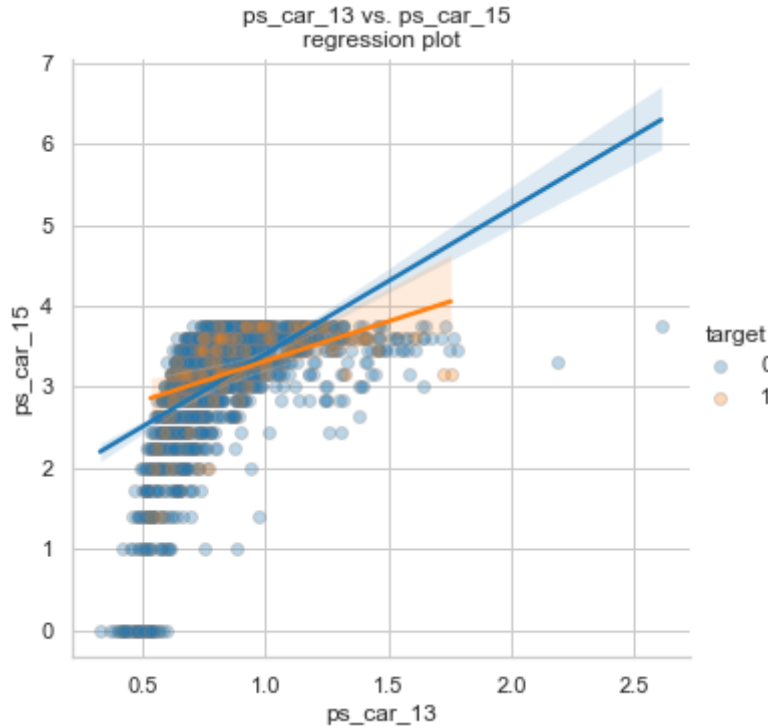
```
In [26]: sns.lmplot(x='ps_car_12', y='ps_car_14', data=trn_sample, hue='target',
                    palette='Set3', height = 5, scatter_kws={'alpha':0.3,'edgecolor'
```

```
Out[26]: <seaborn.axisgrid.FacetGrid at 0x19f003da548>
```



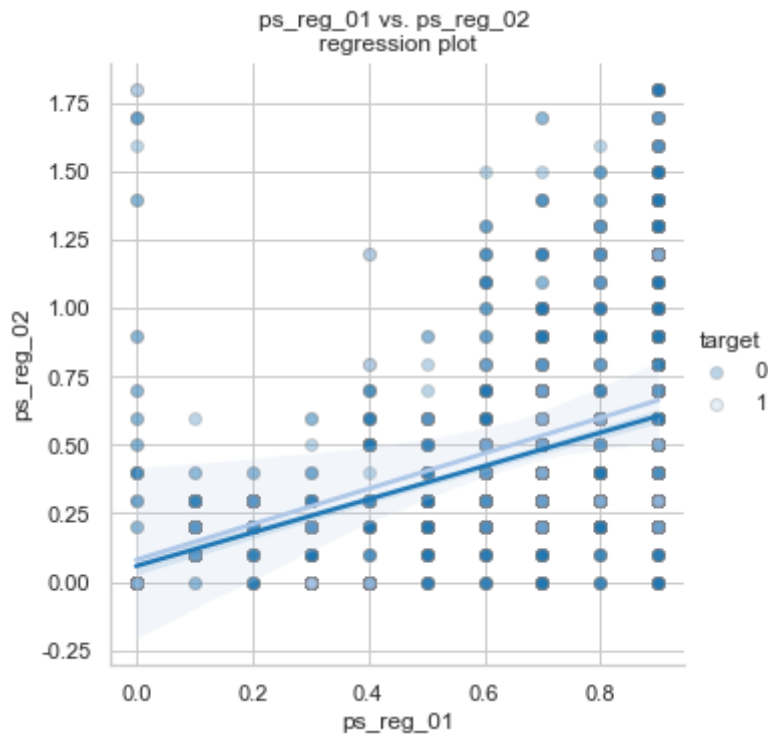

```
In [27]: sns.lmplot(x='ps_car_13', y='ps_car_15', data=trn_sample, hue='target',  
                  palette='tab10', height = 5, scatter_kws={'alpha':0.3,'edgecolor'
```

Out[27]: <seaborn.axisgrid.FacetGrid at 0x19f0f8fa8c8>



```
In [28]: sns.lmplot(x='ps_reg_01', y='ps_reg_02', data=trn_sample, hue='target',  
                  palette='tab20', height = 5, scatter_kws={'alpha':0.3,'edgecolor'
```

Out[28]: <seaborn.axisgrid.FacetGrid at 0x19f0fe41888>



```
In [29]: from xgboost import XGBClassifier
from xgboost import plot_importance
plt.figure(figsize = [100,20])

X = train_imp.drop(['id', 'target'], axis=1)
y = train_imp.target

model = XGBClassifier()

model.fit(X, y)
# plot feature importance
```

```
Out[29]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                      importance_type='gain', interaction_constraints='',
                      learning_rate=0.300000012, max_delta_step=0, max_depth=6,
                      min_child_weight=1, missing=nan, monotone_constraints='()',
                      n_estimators=100, n_jobs=0, num_parallel_tree=1, random_state=0,
                      reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                      tree_method='exact', validate_parameters=1, verbosity=None)

<Figure size 7200x1440 with 0 Axes>
```

```

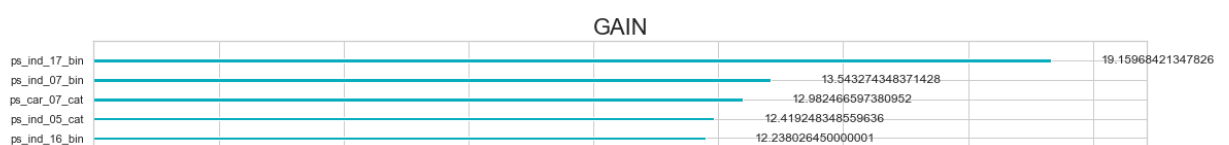
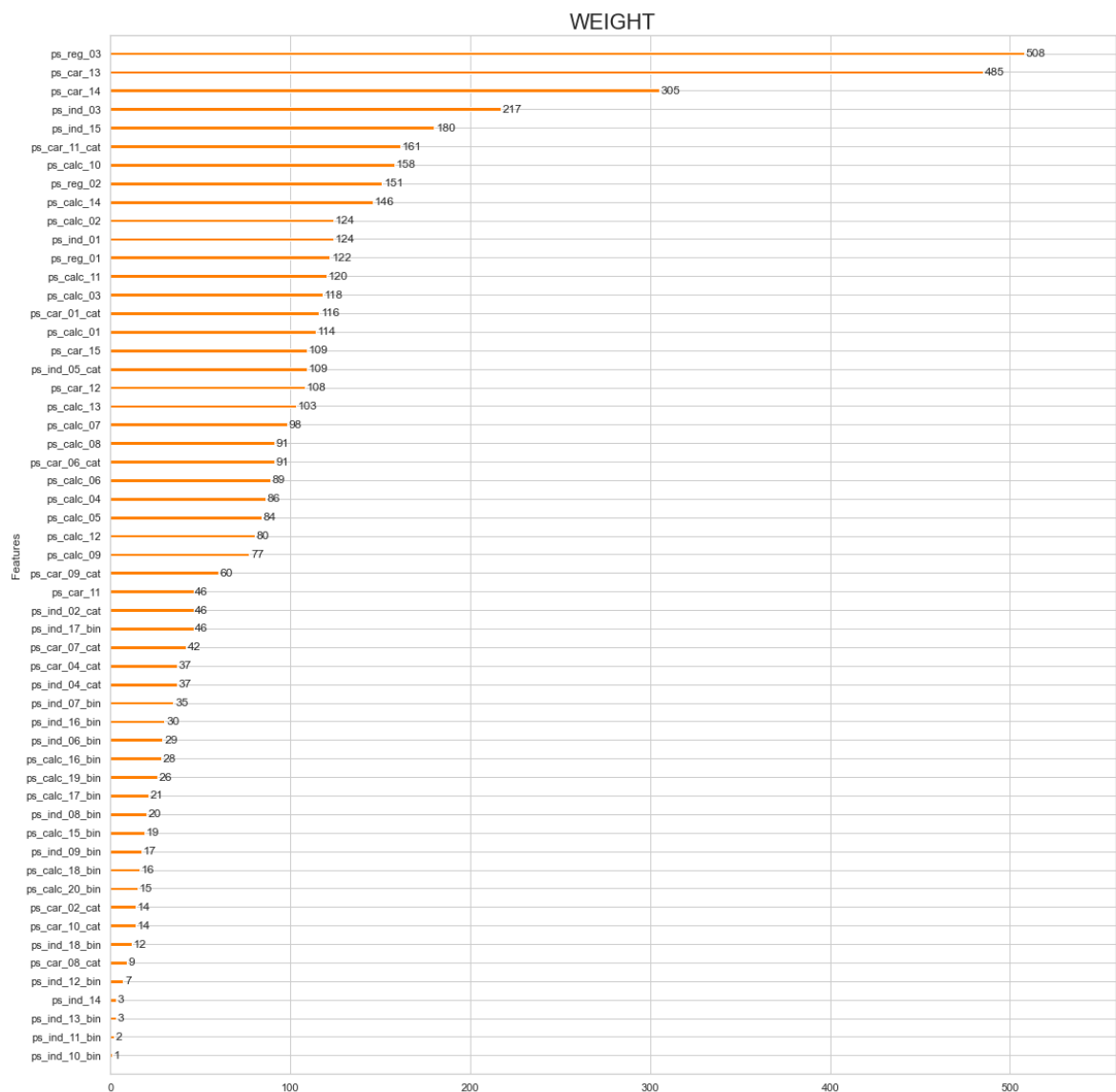
In [30]: # define subplot grid
fig, axs = plt.subplots(nrows=5, ncols=1, figsize=(15, 80))
plt.tight_layout()
plt.subplots_adjust(hspace=0.1)

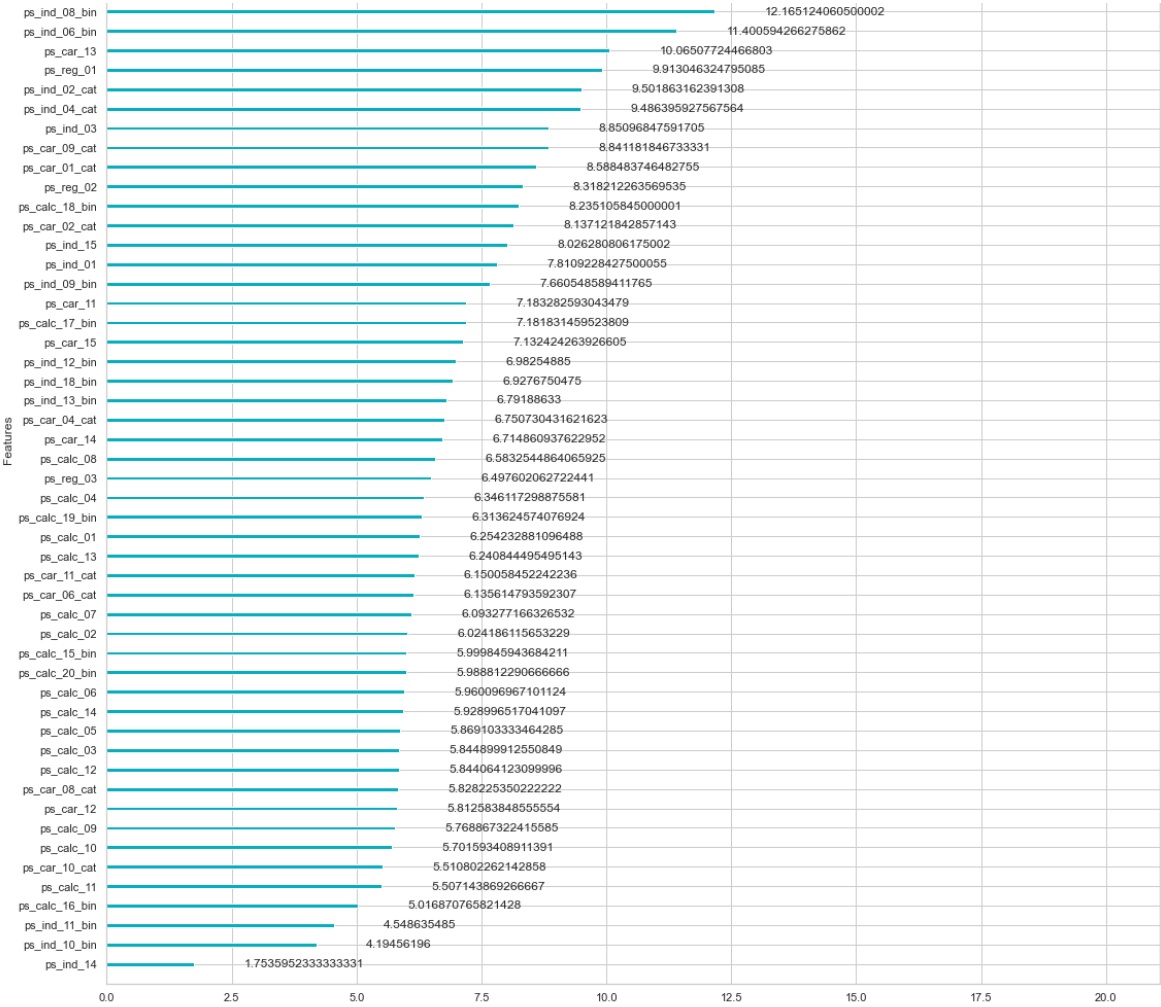
types = ['weight', 'gain', 'cover', 'total_gain', 'total_cover']
# loop through tickers and axes
colors = ['#ff7f01', '#08aebd', '#fc5531', '#139948', '#8950fe']
for ty, ax, color in zip(types, axs.ravel(), colors):
    # filter df for ticker and plot on specified axes
    plot_importance(ax = ax, booster = model, importance_type=ty, color = color)

    # chart formatting
    ax.set_title(ty.upper(), fontsize = 22)
    ax.set_xlabel("")

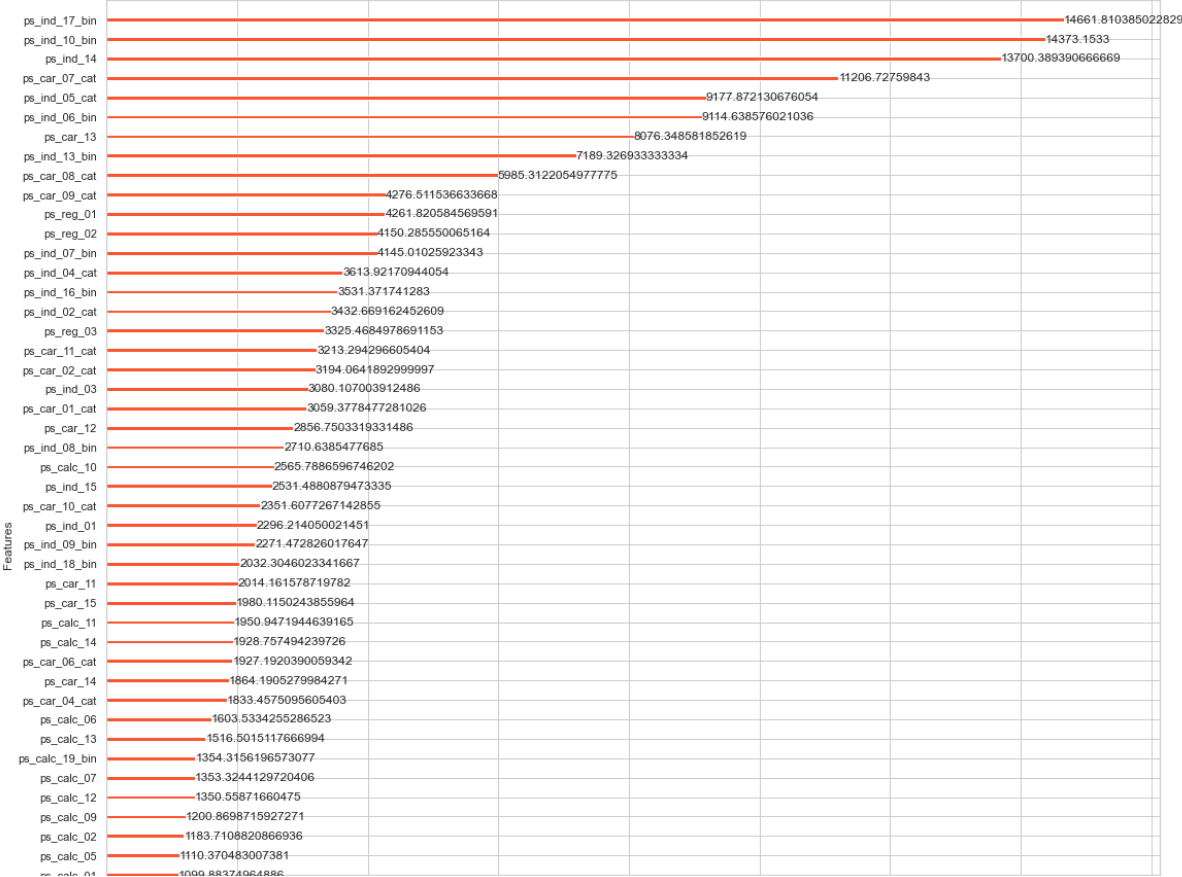
plt.show()

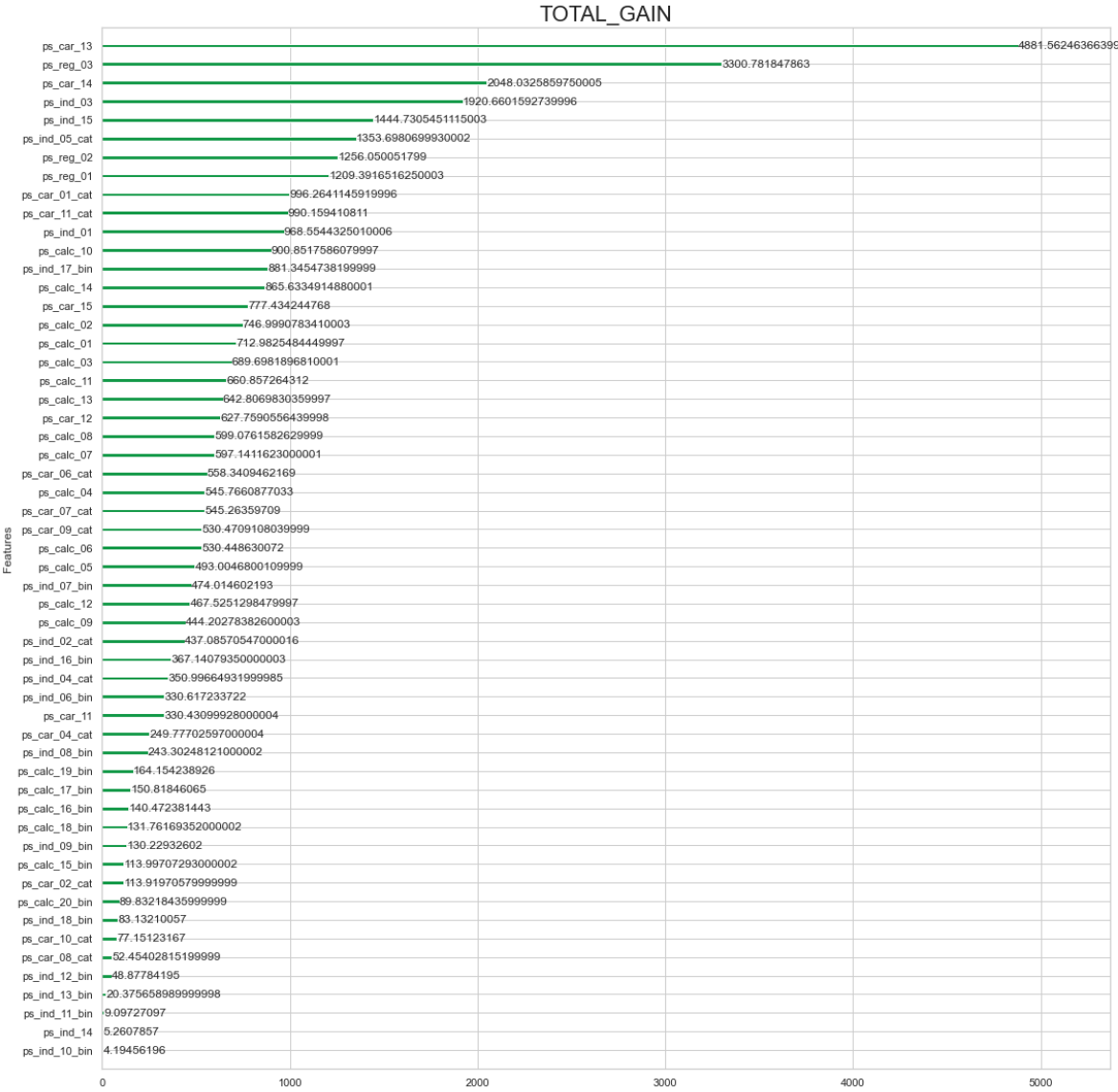
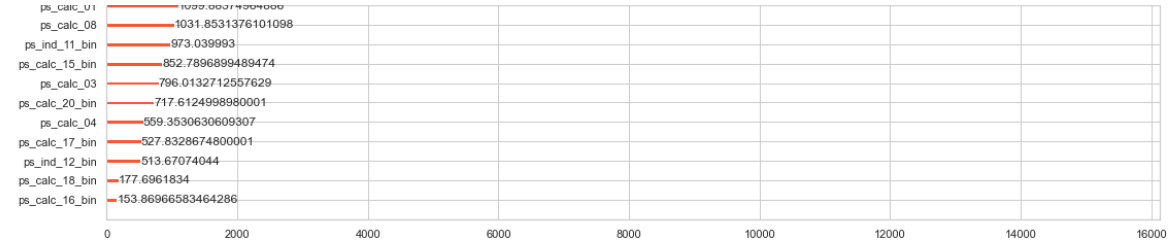
```





COVER







```
In [32]: %%time
mf = mutual_info_classif(train_imp.drop(['id', 'target'], axis=1).values,
                          train_imp.target.values,
                          n_neighbors=30, random_state=2022)
print(mf)
```

```
[1.13688315e-03 4.22158074e-03 1.35599884e-03 1.68686989e-03
 8.57707895e-04 1.86228440e-03 1.12560973e-03 2.77409873e-04
 4.06129362e-04 0.00000000e+00 1.14051723e-05 7.66921627e-05
 0.00000000e+00 1.12701163e-05 1.24501323e-03 4.49947066e-03
 5.92141362e-04 2.36087711e-04 1.34802166e-03 9.28019750e-04
 5.18843733e-04 2.80119096e-03 6.64504597e-03 7.93398089e-04
 1.53250383e-03 8.38246069e-03 6.60925418e-03 5.06933218e-03
 9.00388241e-03 1.44670132e-03 3.91695305e-03 1.48370371e-03
 1.84021357e-03 9.77655896e-04 1.26728143e-03 3.82313257e-04
 3.90869701e-04 2.84671759e-04 1.96475285e-03 2.05495176e-03
 1.96014168e-03 1.60294717e-03 1.88351358e-03 1.79755088e-03
 9.60430682e-04 9.10600275e-04 1.74713733e-03 1.38065000e-03
 8.71062322e-04 1.20886246e-04 3.56145372e-03 2.70045648e-03
 7.07345176e-04 1.10377568e-03 9.43819366e-05]
Wall time: 6min 54s
```

```
In [33]: threshold = 10
high_score_features = []
high_score = []
for score, f_name in sorted(zip(mf, train_imp.drop(['id', 'target'], axis=1)).
    print(f_name, score)
    high_score_features.append(f_name)
    high_score.append(score)
```

```
ps_car_10_cat 0.009003882411455333
ps_car_07_cat 0.008382460694422278
ps_car_02_cat 0.006645045972727637
ps_car_08_cat 0.006609254183864266
ps_car_09_cat 0.005069332176625085
ps_ind_16_bin 0.0044994706649483796
ps_ind_02_cat 0.00422158074290202
ps_car_11 0.003916953051236849
ps_calc_16_bin 0.00356145371805372
ps_car_01_cat 0.0028011909644849453
```

```
In [34]: sns.barplot(y = high_score_features, x = high_score, palette = 'crest')
```

Out[34]: <AxesSubplot:>

