

イントロダクション

このコンペは「ポルトセグロ」というブラジル保険会社の第3位が主催したものだ。

この会社は来年ドライバーが保証金請求をはじめる確率を予想したい。このノートブックではインタラクティブなチャートやplotライブラリを使ったpythonでの視覚化方法を駆使してデータ分析します。

そしていくつかの洞察と他人が真似したくなる美しいプロットを導きます。

plot.lyはplotly社が作った主要ソフトの一つです。この会社はオンラインでの統計データの可視化ツールに特化していて、さらにpython,R,Matlab, node.jsなどの広範囲な言語に対するapiも用意しています。

便利なものをかんたんに紹介します。

- シンプルな水平バープロット ターゲットの分布を調べる
- 相関係数のヒートマップ 異なる特徴間の相関係数を調べる
- 散布図 ランダムフォレストや勾配ブースティングモデルで得られた重要な変数を比べる
- 垂直バープロット 重要な特徴を降順にならべる
- 3D散布図

このノートブックのテーマはかんたんにまとめると下記になる。

- データの質をチェック 欠損値(null -1)を可視化したり、評価したり
- 特徴を調べたり、フィルタリング 相関係数や特徴間の相互情報を目的変数に対してプロットしてみる。変数を2値、カテゴリー、それ以外か調べたり
- モデル学習を通して得られた重要な特徴をランキング 特徴をランク付けするのにランダムフォレストや勾配ブースティングは便利で学習プロセスにもとづいています。

やってみましょう！

In [2]:

```
1 # 関連モジュールをロードしましょう
2 # conda install plotlyをお忘れなく
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 %matplotlib inline
8 import plotly.offline as py
9 py.init_notebook_mode(connected=True)
10 import plotly.graph_objs as go
11 import plotly.tools as tls
12 import warnings
13 from collections import Counter
14 from sklearn.feature_selection import mutual_info_classif
15 warnings.filterwarnings('ignore')
```

pandasを使って与えられたデータを読み込みましょう

In [3]:

```

1 # csvデータを読み込んで先頭から5件表示
2 data_dir = './competitions/porto-seguro-safe-driver-prediction/'
3 train = pd.read_csv(data_dir + 'train.csv')
4 train.head()

```

Out[3]:

	id	target	ps_ind_01	ps_ind_02_cat	ps_ind_03	ps_ind_04_cat	ps_ind_05_cat	ps_ind_06_bin
0	7	0	2	2	5	1	0	0
1	9	0	1	1	7	0	0	0
2	13	0	5	4	9	1	0	0
3	16	0	0	1	2	0	0	1
4	17	0	0	2	0	1	0	1

5 rows × 59 columns

In [4]:

```

1 # データの大きさを見てみましょう
2 rows = train.shape[0]
3 columns = train.shape[1]
4 print("トレーニングデータセットは{}行{}列です.".format(rows, columns))

```

トレーニングデータセットは595212行59列です。

1.データの質を確認する

nullや欠損値を確認しましょう

質チェックとして、トレーニングデータセットにnull値があるかどうかをさっと見ていきます。

In [5]:

```

1 # isnullとanyの動作確認
2 df = pd.DataFrame([[1, 2, 3],[np.nan,5, 6]])

```

In [6]:

```

1 # isnull()でNanかどうかをTrue Falseで返す
2 df.isnull()

```

Out[6]:

	0	1	2
0	False	False	False
1	True	False	False

In [7]:

```
1 # any()でTrueがある行はTrueを返す
2 df.isnull().any()
```

Out[7]:

```
0 True
1 False
2 False
dtype: bool
```

In [8]:

```
1 # 2回で全体としてNanがあるか確認できる
2 df.isnull().any().any()
```

Out[8]:

True

In [9]:

```
1 # any()2回で行、列
2 train.isnull().any().any()
```

Out[9]:

False

これよりNan値はない。nullチェックはFalseだったけど、この結果が本当にこのケースでは-1を欠損値にしているためnull値がないと早とちりしたらダメだ。だから-1のデータをnull値への置き換えを実施しましょう。もう一度確認しよう。

-1が含まれているかみて、下記のようにして-1をnullにかんたんに置換できるよ。

In [10]:

```
1 train_copy = train
2 train_copy = train_copy.replace(-1, np.NaN) # train_copyは-1をnp.NaNに置き換える。
```

次にカグラーが作ったmissingnoというパッケージがデータセットの欠損値を見るのにめっちゃ便利だからやってみよう。

In [11]:

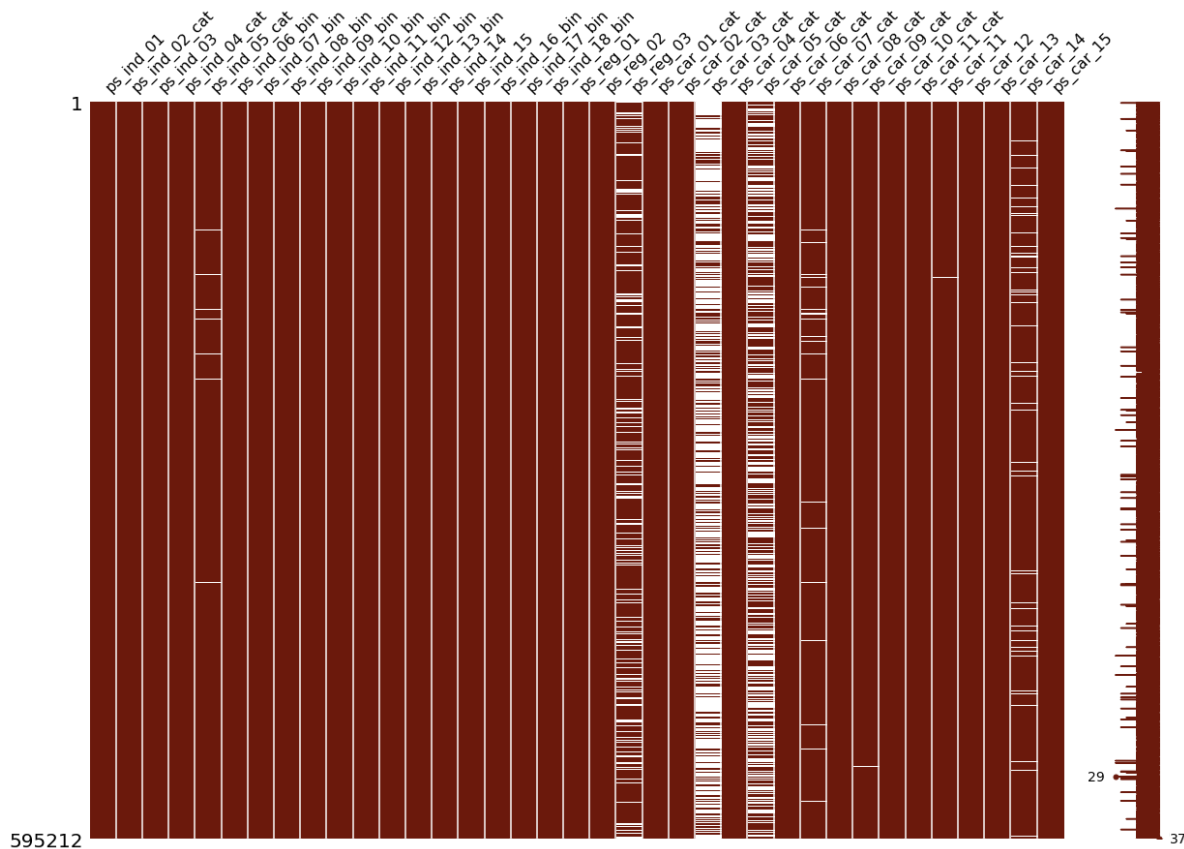
```

1 import missingno as msno
2 msno.matrix(df=train_copy.iloc[:, 2:39], figsize=(20, 14), color=(0.42, 0.1, 0.05))
3 # 欠損値の部分は白い線として表示される

```

Out[11]:

<matplotlib.axes._subplots.AxesSubplot at 0x1079a8128>



見ていると、白いところが欠損値で、縦に赤い帯になっていれば欠損値がないことがひと目で分かるね！この場合、全体59の内7つの項目で欠損値があることが分かります。だけどコメントでJustin Nafeが指摘しているようにホントは13あるようです。これはmissingnoの行列はおおよそ40奇数？の特徴を1度に表示して何個かは弾かれているかもしれない。だからこの5つのnull値を含んだものも取り除かれたんだ。すべてのnull値を見るために、画面サイズ引数をかえて引き裂かれたデータを引っ張ってみましょう。観測できた7つのnull値を含んだ項目を下に挙げておきましょう。

ps_ind_05_cat | ps_reg_03 | ps_car_03_cat | ps_car_05_cat | ps_car_07_cat | ps_car_09_cat | ps_car_14

null値のある項目の殆どの文字の後ろに_catがついています。s_reg_03, ps_car_03_cat, ps_car_05_catについて補足説明しておきましょう。暗い帯の白い部分の比率をみるとわかりますが、この3つはnull値のほうがメジャーで-1に置換するという方法があまりいい戦略ではないかもしれません。

目的変数の調査

データを導く他の普通の方法は目的変数に関することです。この場合では項目に親切に'ターゲット'ってついてますね。目的変数は「分類/ラベル/正誤」などの答えで構成され、まだ見ていないデータを上手く予測できるように予測モデルを学習するためにターゲット変数を除いたトレーニングセットの変数を与えられた教師あり学習モデルが使われる。

In [12]:

```
1 # graphオブジェクトとして、目的変数とその数とタイトルを保存
2 data = [go.Bar(
3     x = train["target"].value_counts().index.values,
4     y = train['target'].value_counts().values,
5     text = '目的変数の分布')]
6
7 # plotlyで描画する
8 layout = go.Layout(title='目的変数の分布')
9
10 fig = go.Figure(data=data, layout=layout)
11
12 py.iplot(fig, filename='basic-bar')
```

うーんデータがかなり偏っていることを心に留めておこう。偏ったターゲットはすばらしいことだと後に分かるよ。

データ型のチェック

このチェックはトレーニングデータセットの型を確認すること。整数、文字、実数など。重複なしのユニークな型を数えるのにいい方法がある。**Collections**モジュールをインポートしていれば、

In [13]:

```
1 # データ型の合計数や構成がひと目で分かる
2 Counter(train.dtypes.values)
```

Out[13]:

```
Counter({dtype('int64'): 49, dtype('float64'): 10})
```

上で分かる通り、全 59 項目構成がわかり、int64とfloat64しかないことが分かります。

他に注意すべきこととして、ポルトセグロ社は項目に_binや_cat、_regをつけてくれて、_binは2値データ、catはカテゴリーデータ、それ以外と教えてくれる。ここでさらに実数値（連続値をもつ特徴）と整数値(2値データ、カテゴリーデータ、ordinalな特徴に)とシンプルにしてみましょう。

In [14]:

```
1 # 実数型と整数型の項目をそれぞれ変数に格納
2 train_float = train.select_dtypes(include=['float64'])
3 train_int = train.select_dtypes(include=['int64'])
```

相関係数のプロット

手始めに、特徴が他のものとどう関連しているか、またもしかしたらそこから新しい気づきを得られるために相関係数をプロットしてみましょう。この際、seabornをつかって相関係数をヒットマップでぶろっとします。都合のいいことに、pandasのデータフレームにはcorr()関数が備わっているので、ピアソン相互相関係数を計算できます。またseabornの相関プロットも便利です。文字通り"heatmap"で使えます。

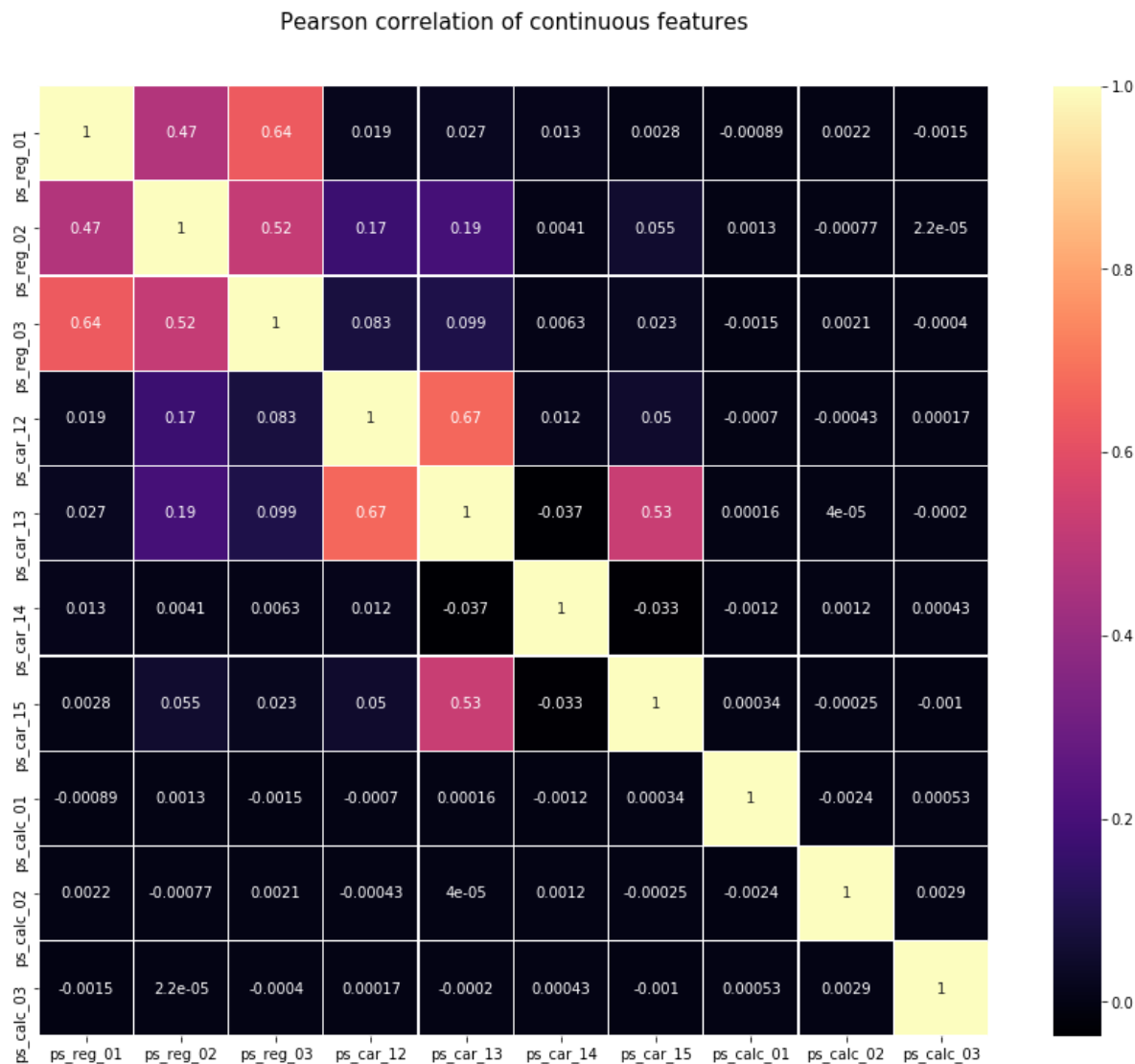
実数値の特徴の相関係数

In [15]:

```

1 # カラーマップの設定 マグマ
2 colormap = plt.cm.magma
3 plt.figure(figsize=(16, 12))
4 plt.title('Pearson correlation of continuous features', y=1.05, size=15)
5 sns.heatmap(train_float.corr(), linewidths=0.1, vmax=1.0, square=True,
6             cmap=colormap, linecolor='white', annot=True);
7 # linecolor= マス目の線を白色に、annot=True数値を表示

```



相関係数プロットから、だいたい0で他と無相関だと分かります。これはかなり興味深い観測結果で、これ以降のさらなる調査をwarrantします。さしあたって、正の相関があるペアは下記だ

(ps_reg_01, ps_reg_03)

(ps_reg_02, ps_reg_03)

(ps_car_12, ps_car_13)

(ps_car_13, ps_car_15) ※1,2はだめ？

整数値の特徴の相関関係

整数型の項目に対して、Plotlyに切り替えて相関係数のヒートマップを表示してみましょう。前のplotlyプロットほとんど同じように、ヒートマップを作るために、`go.Heatmap`を呼び出してみましょう。ここで3つの軸に値をセットする必要がある、`x,y`に項目名、`z`軸に相関係数値を。カラースケール属性は決められたワードを入れてね。（今回はGreysを使うけど、自分でPortlandやViridisも試してみて）

In [16]:

```
1 # Seabornでやるとすればこうかな?
2 #train_int = train_int.drop(["id", "target"], axis=1)
3 # colormap = plt.cm.bone
4 # plt.figure(figsize=(21,16))
5 # plt.title('Pearson correlation of categorical features', y=1.05, size=15)
6 # sns.heatmap(train_cat.corr(),linewidths=0.1,vmax=1.0, square=True, cmap=colormap, linecolor=
7
8 # plotly流のheatmapの表示方法
9 # ピポットしなくても3変数を指定するだけで表示できる
10 # まマウスカーソルを合わせると情報を表示してくれる
11 # 使い方は深入りしないが下記のサイトは参考になりそう。
12 # https://images.plot.ly/plotly-documentation/images/python\_cheat\_sheet.pdf
13 # https://qiita.com/ogrew/items/e6ba517cec9fb5161cb6
14 # https://qiita.com/inoory/items/12028af62018bf367722
15 data = [
16     go.Heatmap(
17         z = train_int.corr().values,
18         x = train_int.columns.values,
19         y = train_int.columns.values,
20         colorscale = 'Viridis',
21         reversescale = False,
22         text = True,
23         opacity = 1.0)
24 ]
25
26 layout = go.Layout(
27     title='Peason Correlation of Integer-type features',
28     xaxis = dict(ticks="", nticks=26),
29     yaxis = dict(ticks=""),
30     width = 900, height = 700)
31
32 fig = go.Figure(data=data, layout=layout)
33 py.iplot(fig, filename='labelled-heatmap')
34 # 関連のあるところにカーソルをもっていくと項目情報が表示されるので分かりやすい
```

同様に、ほとんどが線形相関がないのがわかって、かなり多くが0だと明らかになったね。これはかなり有益で、とくにもし主成分分析などはっきりと相関関係を必要とするような次元圧縮をやろうとするときにね。

負の相関関係の特徴

ps_ind_06_bin, ps_ind_07_bin, ps_ind_08_bin, ps_ind_09_bin

特筆すべき点は最初のヌル値の分析のときにpscar-3_catとps_car_05_catは多くの欠損値かヌル値を含んでいたことです。それゆえこれらの特徴が強い正の相関があっても不思議はない。たとえ一方が裏に潜んだ真実を反映していないかもしれない。

相互情報プロット

相互情報は対応する特徴やターゲット変数間の相互情報を調べる別の有益なツールです。

分類問題に対して、sklearnのmutual_info_classifを呼んでサクッと依存関係を調べられます。（0が無関係で1にいくほど強い相関です）それゆえこれは特徴内に含まれているかもしれない目的変数との多くの情報へのアイディアになります。

mutual_info_classif関数のsklearn実装はK近傍法から推測されたエントロピーに基づいたノンパラメトリックな方法になっています。くわしくは公式ホームページを見てね。

In [17]:

```
1 # 相互情報量を表示する。
2 mf = mutual_info_classif(train_float.values, train.target.values, n_neighbors=3, random_state=17)
3 print(mf)
```

```
[0.01402035 0.00431986 0.0055185  0.00778454 0.00157233 0.00197537
 0.01226   0.00553038 0.00545101 0.00562139]
```

2値特徴の調査

ほかに調査したいかもしれないデータの観点は2値しか取らない項目（たとえば0,1のどちらか）だろう。なりゆきでちょうど2値データをすべて保管しているから、垂直なplotlyバープロットをやってみよう。

In [18]:

```
1 bin_col = [col for col in train.columns if '_bin' in col] # リスト内包表記 _binのつく項目だけをbin_colに格納
2 zero_list = []
3 one_list = []
4 # 各項目ごとに0,1の数を集計
5 for col in bin_col:
6     zero_list.append((train[col]==0).sum())
7     one_list.append((train[col]==1).sum())
```

In [19]:

```
1  # 2値を視覚的にわかりやすく表示する
2  trace1 = go.Bar(
3      x=bin_col,
4      y=zero_list,
5      name='Zero count'
6  )
7  trace2 = go.Bar(
8      x=bin_col,
9      y=one_list,
10     name='One count'
11 )
12
13 data = [trace1, trace2]
14 layout = go.Layout(
15     barmode='stack',
16     title='Count of 1 and 0 in binary variables'
17 )
18
19 fig = go.Figure(data=data, layout=layout)
20 py.iplot(fig, filename='stacked-bar')
21 # すごく分かりやすい!!
```

ここで4つの特徴がみられる。ps_ind_10_bin,ps_ind_11_bin,ps_ind_12_bin,ps_ind_13_binはほぼゼロしかないね。この情報はこれらの特徴が有用なのかターゲット変数に対して大した情報がふくまれているかという疑問を招きます。

カテゴリーや順序のある特徴の調査

まず後ろに_catがついた特徴を見ていきましょう。

ランダムフォレストで得る重要な変数

まずランダムフォレストを使ってみて、学習が終わったあとの特徴のランニングを見てみましょう。これはアンサンブル学習をつかった手早い方法です。（集計されたブートストラップの下で弱意思決定木学習者のアンサンブル） 有用な特徴を見つけるのにたいへんなパラメータチューニングは必要ないし、目的変数の不均衡さに対しても堅牢です。ランダムフォレストを呼んでみよう

In [21]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 rf = RandomForestClassifier(n_estimators=150, max_depth=8, min_samples_leaf=4, max_features
3                             n_jobs=-1, random_state=0)
4 rf.fit(train.drop(['id', 'target'], axis=1), train.target)
5 features = train.drop(['id', 'target'], axis=1).columns.values
6 print("----- Training Done -----")
```

----- Training Done -----

Plot.lyの特徴の重要度の散布図

ランダムフォレストで学習して、`featureimportances`属性を呼んで特徴の重要度のリストを得て、Plotlyの散布図でプロットしましょう。

Scatterコマンドを呼んで、前のPlotlyプロットごとに、xy軸を呼び出さなければならない。しかしながら散布図で注意すべきことはマーカー属性です。定義したマーカー属性はプロット点の大きさ、色をコントロールできます。

In [22]:

```
1  # 散布図
2  trace = go.Scatter(
3      y = rf.feature_importances_,
4      x = features,
5      mode='markers',
6      marker=dict(
7          sizemode = 'diameter',
8          sizeref = 1,
9          size = 13,
10         # size = rf.feature_importances_,
11         # color = np.random.randn(500), # set color equal to a variable
12         color = rf.feature_importances_,
13         colorscale='Portland',
14         showscale=True
15     ),
16     text = features
17 )
18 data = [trace]
19
20 layout = go.Layout(
21     autosize = True,
22     title = 'Random Forest Feature Importance',
23     hovermode='closest',
24     xaxis= dict(
25         ticklen=5,
26         showgrid=False,
27         zeroline=False,
28         showline=False
29     ),
30     yaxis=dict(
31         title = 'Feature Importance',
32         showgrid=False,
33         zeroline=False,
34         ticklen=5,
35         gridwidth=2
36     ),
37     showlegend = False
38 )
39 fig = go.Figure(data=data, layout=layout)
40 py.iplot(fig, filename='scatter2010')
```

さらに重要さの順番で並べ替えて表示できます。下のように高いものから低い順でバープロットしてみます。

In [23]:

```
1 # 特徴を重要度でソートしてリストに格納する
2 x, y = (list(x) for x in zip(*sorted(zip(rf.feature_importances_, features), reverse=False)))
3
4 # 横のバースプロット
5 trace2 = go.Bar(
6     x=x,
7     y=y,
8     marker=dict(
9         color=x,
10         colorscale = 'Viridis',
11         reversescale = True
12     ),
13     name='Random Forest Feature importance',
14     orientation='h', # これで縦のバースプロットにする
15 )
16 layout = dict(
17     title='Barplot of Feature importances',
18     width = 900, height = 2000,
19     yaxis=dict(
20         showgrid=False,
21         showline=False,
22         showticklabels=True,
23         # domain=[0, 0.85],
24     ))
25
26 fig1 = go.Figure(data=[trace2])
27 fig1['layout'].update(layout)
28 py.iplot(fig1, filename='plots')
```

決定木の可視化

面白いテクニックとして決定木の枝を可視化することがよくあります。簡素化のために、深さ3で学習したので3レベルの深さの枝までしか見れませんが、sklearnのexport_graphvizのグラフ可視化属性をエクスポートしたり、このノートブックで見たイメージをインポートしたりエクスポートできます。

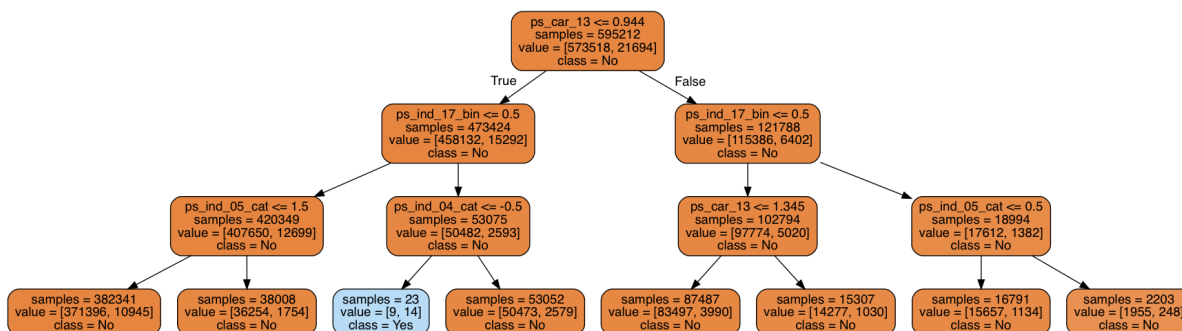
In [34]:

```

1 from sklearn import tree
2 from IPython.display import Image as PlImage
3 from subprocess import check_call
4 from PIL import Image, ImageDraw, ImageFont
5 import re # 正規表現
6
7 decision_tree = tree.DecisionTreeClassifier(max_depth = 3)
8 decision_tree.fit(train.drop(['id', 'target'], axis=1), train.target)
9
10 # 学習モデルを.dotファイルにエクスポート
11 with open("tree1.dot", 'w') as f:
12     f = tree.export_graphviz(decision_tree,
13                             out_file=f,
14                             max_depth = 4,
15                             impurity = False,
16                             feature_names = train.drop(['id', 'target'], axis=1).columns.values,
17                             class_names = ['No', 'Yes'],
18                             rounded = True,
19                             filled = True)
20
21 # .dotを.pngファイルに変換
22 # brew install graphvizする必要あり
23 # https://qiita.com/umasoya/items/c1c4e583f393b27ebb7a
24 check_call(['dot', '-Tpng', 'tree1.dot', '-o', 'tree1.png'])
25
26 # Annotating chart with PIL
27 img = Image.open("tree1.png")
28 draw = ImageDraw.Draw(img)
29 img.save('sample-out.png')
30 PlImage("sample-out.png", )
31 # カレントディレクトリに.dotと.pngファイルが生成されることを確認

```

Out[34]:



勾配ブースティングを使った特徴の重要度

好奇心のために、別のモデルで学習してみましょう。ここでは勾配ブースティングの分類器を使います。勾配ブースティングは段階を踏んで進み、それぞれの段階で回帰木は損失関数の勾配が与えられます。(sklearn実装では デビアンスがデフォルトです。)

In [24]:

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 # 勾配ブースティング
3 gb = GradientBoostingClassifier(n_estimators=100, max_depth=3, min_samples_leaf=4,
4                               max_features=0.2, random_state=0)
5 # モデル学習
6 gb.fit(train.drop(['id', 'target'], axis=1), train.target)
7 # idと目的変数を除いた項目名を特徴として保持
8 features = train.drop(['id', 'target'], axis=1).columns.values
9 print('-----Training Done -----')
```

-----Training Done -----

In [25]:

```
1  # 特徴とその重要度を散布図としてプロットする
2  trace = go.Scatter(
3      y = gb.feature_importances_,
4      x = features,
5      mode='markers',
6      marker=dict(
7          sizemode = 'diameter',
8          sizeref = 1,
9          size = 13,
10         # size = rf.feature_importances_,
11         # color = np.random.randn(500), # set color equal to a variable
12         color = gb.feature_importances_,
13         colorscale='Portland',
14         showscale=True
15     ),
16     text = features
17 )
18
19 data = [trace]
20
21 layout = go.Layout(
22     autosize = True,
23     title='Gradient Boosting Machine Feature Importance',
24     hovermode='closest',
25     xaxis=dict(
26         ticklen=5,
27         showgrid=False,
28         zeroline=False,
29         showline=False
30     ),
31     yaxis=dict(
32         title = 'Feature Importance',
33         showgrid = False,
34         zeroline=False,
35         ticklen=5,
36         gridwidth=2
37     ),
38     showlegend=False
39 )
40 fig = go.Figure(data=data, layout=layout)
41 py.iplot(fig, filename='scatter2010')
```


In [26]:

```
1  # さきほどと同じように重要度順にバープロットしてみる。
2  x, y = (list(x) for x in zip(*sorted(zip(gb.feature_importances_, features),
3                                     reverse = False)))
4  trace2 = go.Bar(
5      x=x,
6      y=y,
7      marker=dict(
8          color=x,
9          colorscale = 'Viridis',
10         reversescale = True
11     ),
12     name="Gradient Boosting Classifier Feature importance",
13     orientation='h',
14 )
15
16 layout =dict(
17     title='Barplot of Feature importances',
18     width = 900, height =2000,
19     yaxis = dict(
20         showgrid=False,
21         showline=False,
22         showticklabels=True,
23     ))
24 fig1 = go.Figure(data=[trace2])
25 fig1['layout'].update(layout)
26 py.iplot(fig1, filename='plots')
```

興味深いことに、ランダムフォレストと勾配ブースティングどちらも最も重要な特徴は同じで**ps_car_13**のようです。

この特徴はかなりさらなる調査が必要なのでより深く調べることを検討してください。

結論

これまで広範囲に渡る調査をやってきました。ヌル値とデータの質の調査、特徴間の相関関係を調べ、ラ特徴の分布の幾つかを調べ、ランダムフォレストと勾配ブースティングという2つの学習モデルを行ってみた。そしてモデルが重要だと示す特徴を特定できた！

続く..