

CASE Study-1 - Don't not Overfit ||

Data Source : <https://www.kaggle.com/c/dont-overfit-ii/data> (<https://www.kaggle.com/c/dont-overfit-ii/data>)

The primary goal of this case study to avoid overfitting problem of given small amount of training sample

We have 250 samples data and 300 features and 1 id and 1 class label(0/1) We have 19750 features and 300 features and 1 id,

So, with small amount of training data set, We must to do task carefully to avoid overfitting.

Import Necessary Libraries

In []:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Read the training data with the help of pandas library

In []:

```
train_data = pd.read_csv('train.csv')
train_data.head()
```

Out[]:

	id	target	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	
0	0	1.0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174	-0.
1	1	0.0	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483	0.5
2	2	0.0	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587	1.4
3	3	1.0	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361	0.9
4	4	0.0	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697	-1.

5 rows × 302 columns

In []:

```
train_data.columns
```

Out[]:

```
Index(['id', 'target', '0', '1', '2', '3', '4', '5', '6', '7',
       '8', '9', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '20', '21', '22', '23', '24', '25', '26', '27', '28', '29', '30', '31', '32', '33', '34', '35', '36', '37', '38', '39', '40', '41', '42', '43', '44', '45', '46', '47', '48', '49', '50', '51', '52', '53', '54', '55', '56', '57', '58', '59', '60', '61', '62', '63', '64', '65', '66', '67', '68', '69', '70', '71', '72', '73', '74', '75', '76', '77', '78', '79', '80', '81', '82', '83', '84', '85', '86', '87', '88', '89', '90', '91', '92', '93', '94', '95', '96', '97', '98', '99', '100', '101', '102', '103', '104', '105', '106', '107', '108', '109', '110', '111', '112', '113', '114', '115', '116', '117', '118', '119', '120', '121', '122', '123', '124', '125', '126', '127', '128', '129', '130', '131', '132', '133', '134', '135', '136', '137', '138', '139', '140', '141', '142', '143', '144', '145', '146', '147', '148', '149', '150', '151', '152', '153', '154', '155', '156', '157', '158', '159', '160', '161', '162', '163', '164', '165', '166', '167', '168', '169', '170', '171', '172', '173', '174', '175', '176', '177', '178', '179', '180', '181', '182', '183', '184', '185', '186', '187', '188', '189', '190', '191', '192', '193', '194', '195', '196', '197', '198', '199', '200', '201', '202', '203', '204', '205', '206', '207', '208', '209', '210', '211', '212', '213', '214', '215', '216', '217', '218', '219', '220', '221', '222', '223', '224', '225', '226', '227', '228', '229', '230', '231', '232', '233', '234', '235', '236', '237', '238', '239', '240', '241', '242', '243', '244', '245', '246', '247', '248', '249', '250'], dtype='object', length=302)
```

In []:

```
train_data.shape
```

Out[]:

```
(250, 302)
```

In []:

```
train_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 250 entries, 0 to 249
Columns: 302 entries, id to 299
dtypes: float64(301), int64(1)
memory usage: 590.0 KB
```

In []:

```
train_data.describe()
```

Out[]:

	id	target	0	1	2	3	4	5	6	
count	250.000000	250.000000	250.000000	250.000000	250.000000	250.000000	250.000000	250.000000	250.000000	250.000000
mean	124.500000	0.268000	-0.098064	0.001208	0.090680	-0.122248	0.011500	-0.116624	0.006932	0.100988
std	72.312977	0.443806	0.996063	0.955117	0.968065	0.933001	0.945662	1.081705	1.014091	1.028042
min	0.000000	0.000000	-3.181000	-3.041000	-2.967000	-2.898000	-2.837000	-3.831000	-2.873000	-2.489000
25%	62.250000	0.000000	-0.756250	-0.624750	-0.515750	-0.695500	-0.678000	-0.758500	-0.646250	-0.589000
50%	124.500000	0.000000	-0.064500	-0.008000	0.067500	-0.090000	0.028000	-0.073500	-0.076500	0.104500
75%	186.750000	1.000000	0.647750	0.493250	0.716000	0.436250	0.625250	0.554250	0.676500	0.717000
max	249.000000	1.000000	2.347000	3.138000	2.609000	2.590000	2.413000	2.687000	2.793000	3.766000

8 rows × 302 columns

Observation

Train data has 300 independent features(0-299) and one dependent feature(target) along with unique id. All the values of feature are continuous and have same distribution from distribution table.

In []:

```
print(train_data["target"].value_counts())
```

```
0.0    183  
1.0     67  
Name: target, dtype: int64
```

We found 183 count for target value 0 and 67 for target value 1. This is imbalanced dataset(Not highly imbalanced but decent).

In []:

```
nul=0  
for col in train_data.columns:  
    if(train_data[col].isnull().any()):  
        print(col,"has null value")  
        nul=1  
if(nul==0):  
    print("There is no null value")
```

There is no null value

There are no null values found and all are real values

Exploratory Data Analysis(EDA)

In []:

```
sns.set(rc={'figure.figsize':(30,10)})
plt.subplot(1,2,1)
ax = sns.countplot(train_data['target'])
for p in ax.patches:
    ax.annotate('{:.1f}%'.format(100*p.get_height()/len(train_data)), (p.get_x()+0.35, p.get_height()+1))
ax.set_title('count plot of target distribution')
ax.set_xlabel('Target')
ax.set_ylabel('count')

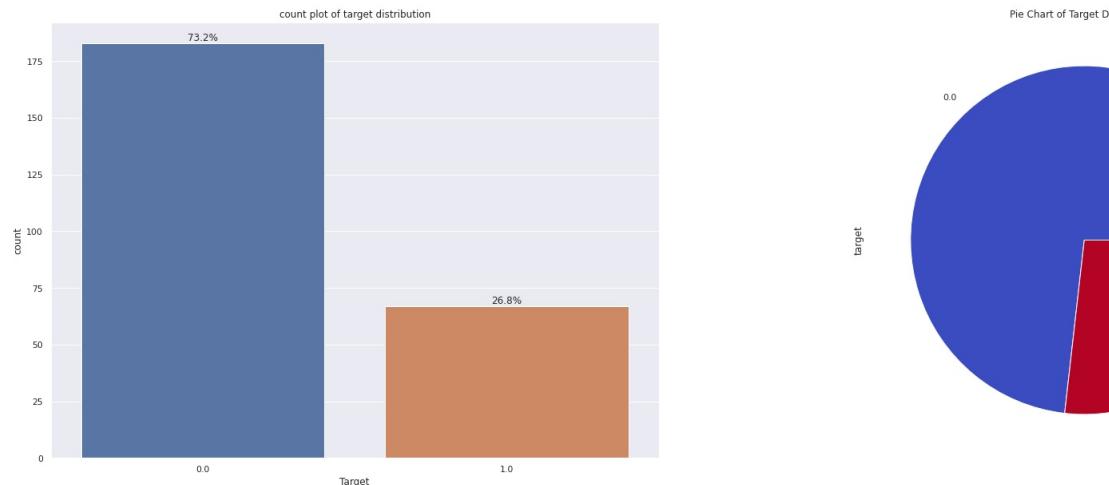
plt.subplot(1,2,2)
ax = train_data['target'].value_counts().plot(kind='pie', colormap='coolwarm')
ax.set_title("Pie Chart of Target Distribution")
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation

FutureWarning

Out[]:

Text(0.5, 1.0, 'Pie Chart of Target Distribution')



Observation

It is clearly said that data is imbalanced, where data points belonging to target 0 is 73.2% and 1 is 26.8%. Thus it is imbalanced data.

New Thing : I have tried to get mean and standard deviation

In []:

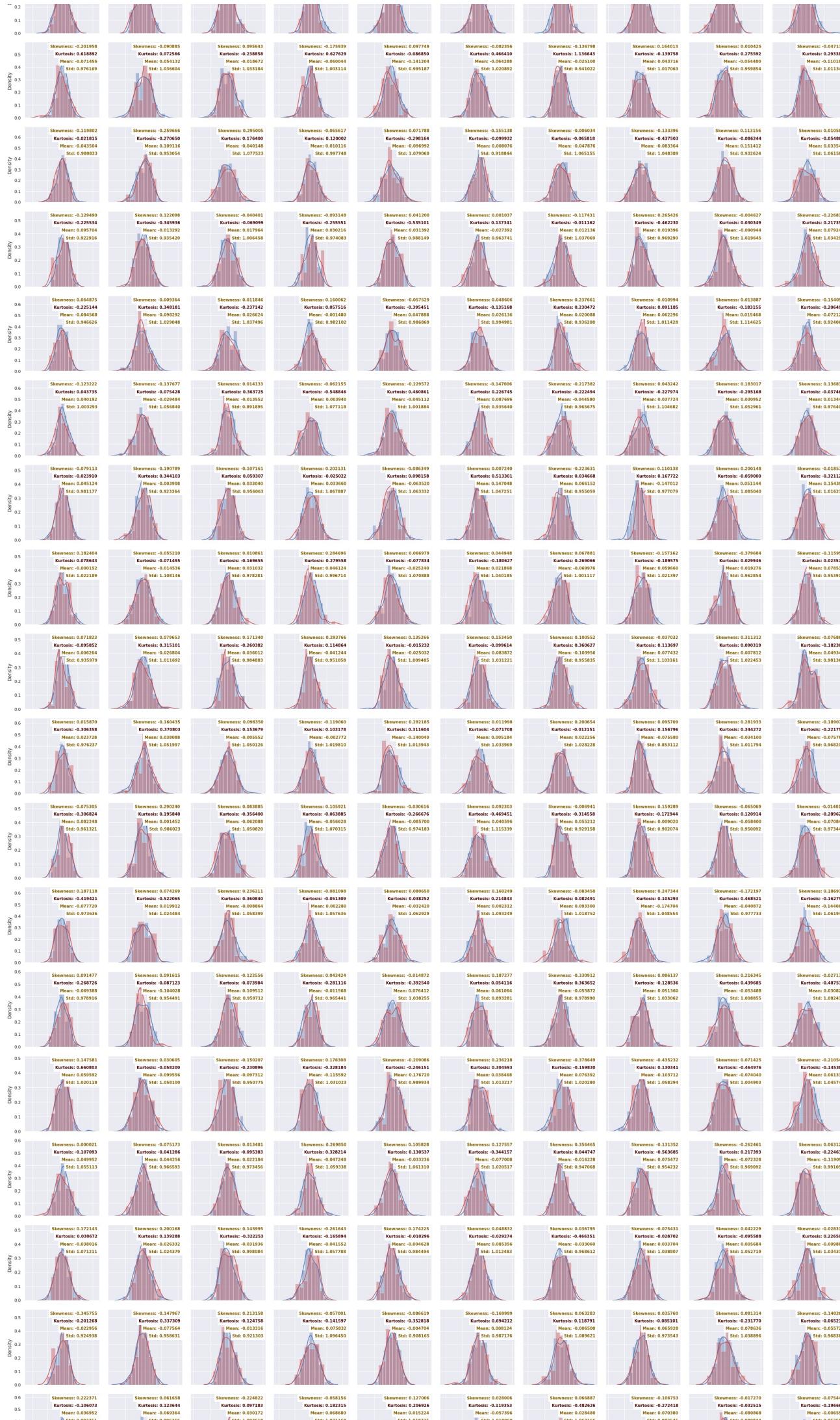
#<https://www.analyticsvidhya.com/blog/2021/05/shape-of-data-skewness-and-kurtosis/>

```
# https://stackoverflow.com/questions/50940283/show-metrics-like-kurtosis-skewness-on-distribution-plot-using-seaborn-in-python
import warnings
warnings.filterwarnings("ignore")
sns.set(rc={'figure.figsize':(30,90)})
fig, ax = plt.subplots(30, 10,sharex='col', sharey='row')
ax = ax.reshape(-1)
for i, col in enumerate(train_data.columns[2:]):
    g0 = train_data[train_data['target']==0.0][col]
    g1 = train_data[train_data['target']==1.0][col]
    sns.distplot(g0, label = 'target 0', ax=ax[i], color='b')
    sns.distplot(g1, label = 'target 1', ax=ax[i], color='r')
min_skew, max_skew = 100, -100
min_kurt, max_kurt = 100, -100
for i, x in enumerate(ax):
    skew = train_data.iloc[:,i+2].skew()
    min_skew = min(min_skew, skew)
    max_skew = max(max_skew, skew)
    x.text(x=0.97, y=0.97, transform=x.transAxes, s="Skewness: %f" % skew,\n        fontweight='demibold', fontsize=10, verticalalignment='top', horizontalalignment='right',\n        backgroundcolor='white', color='xkcd:poo brown')
    kurt = train_data.iloc[:,i+2].kurt()
    min_kurt = min(min_kurt, kurt)
    max_kurt = max(max_kurt, kurt)
    x.text(x=0.97, y=0.87, transform=x.transAxes, s="Kurtosis: %f" % kurt,\n        fontweight='demibold', fontsize=10, verticalalignment='top', horizontalalignment='right',\n        backgroundcolor='white', color='xkcd:dried blood')
for i, x in enumerate(ax):
    mean = train_data.iloc[:,i+2].mean()
    x.text(x=0.97, y=0.77, transform=x.transAxes, s="Mean: %f" % mean,\n        fontweight='demibold', fontsize=10, verticalalignment='top', horizontalalignment='right',\n        backgroundcolor='white', color='xkcd:poo brown')
for i, x in enumerate(ax):
    std = train_data.iloc[:,i+2].std()
    x.text(x=0.97, y=0.67, transform=x.transAxes, s="Std: %f" % std,\n        fontweight='demibold', fontsize=10, verticalalignment='top', horizontalalignment='right',\n        backgroundcolor='white', color='xkcd:poo brown')
plt.tight_layout()
plt.show()
```

print('Skewness Range:[{},{}]\'.format(min_skew,max_skew))

print('Kurtosis Range:[{},{}]\'.format(min_kurt,max_kurt))







Skewness Range: [-0.4889114906320757, 0.4864285153563513]
Kurtosis Range: [-0.5636848359050028, 1.5206991488929726]

In []:

```
print('Skewness Range:[{},{}]\t'.format(min_skew,max_skew))
print('Kurtosis Range:[{},{}]\t'.format(min_kurt,max_kurt))
```

Skewness Range: [-0.4889114906320757, 0.4864285153563513]
Kurtosis Range: [-0.5636848359050028, 1.5206991488929726]

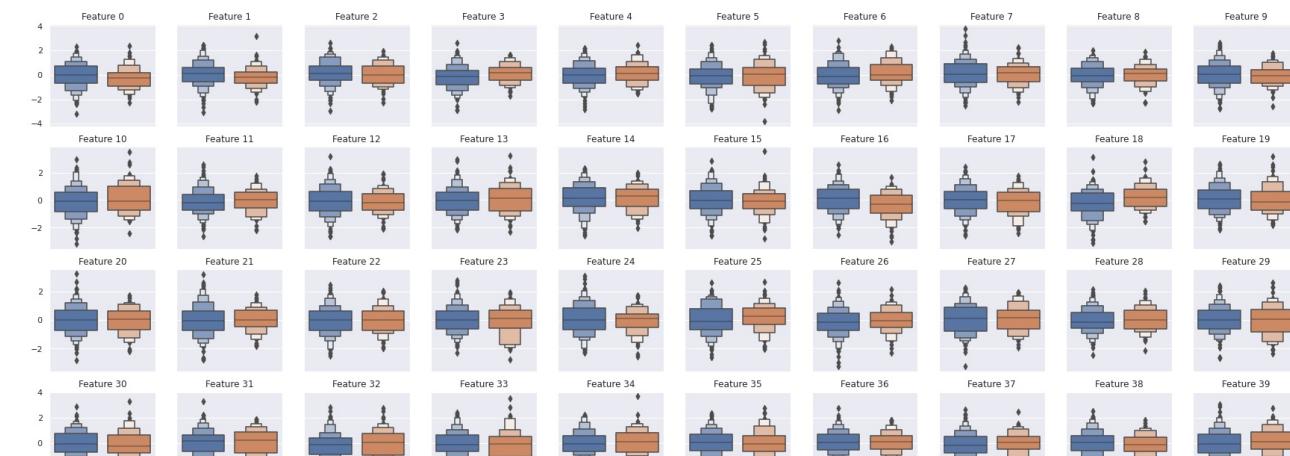
Based on above distribution, We have found skewness range between -0.5 to 0.5 and kurtosis range is less than 3. The distribution for all 300 features are showing almost Symetric distribution.

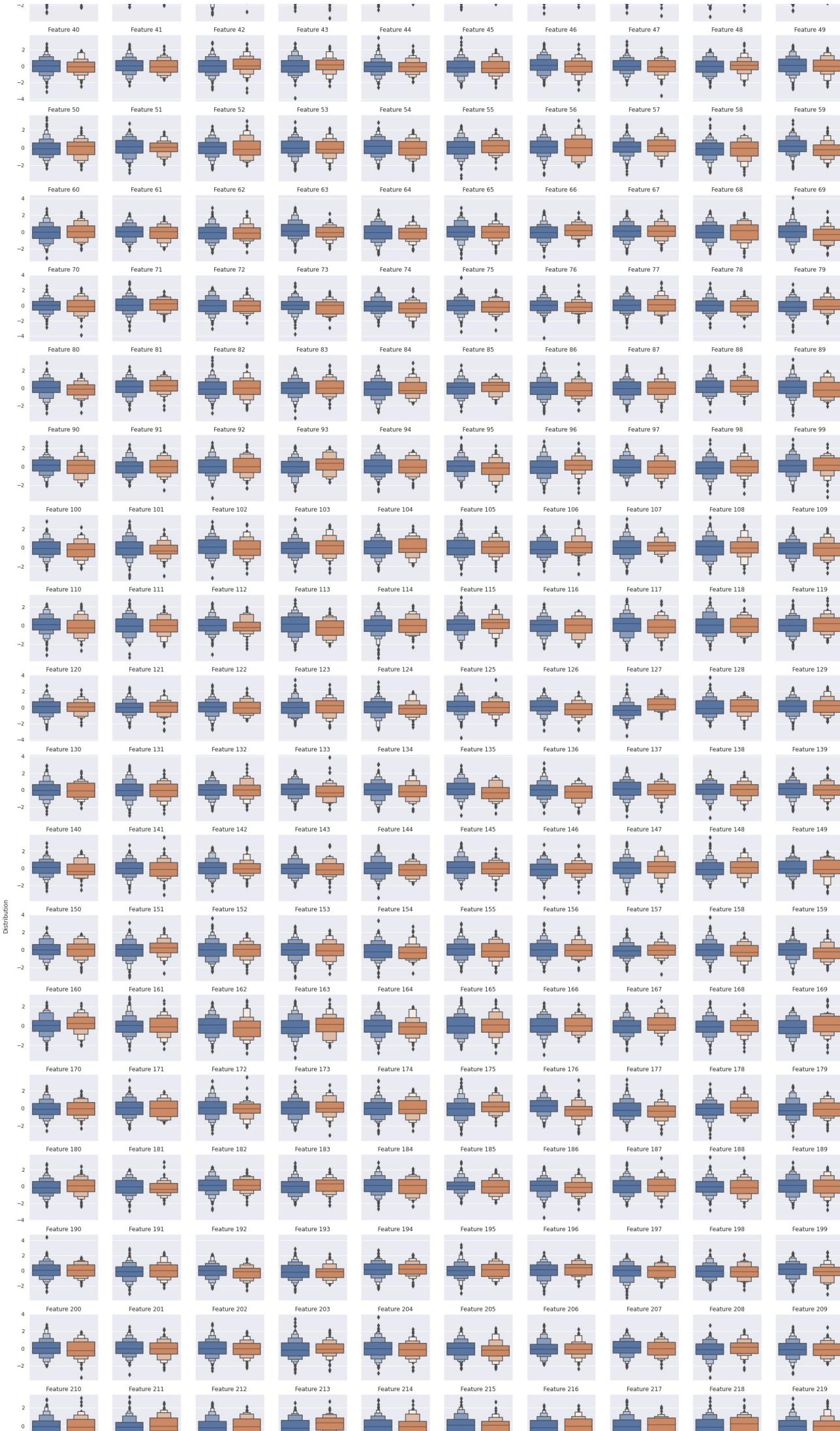
The mean for all 300 features are almost close to 0. And The mean for all 300 features are almost close to 1. The Distributions are very close to Standard Gaussian Distribution

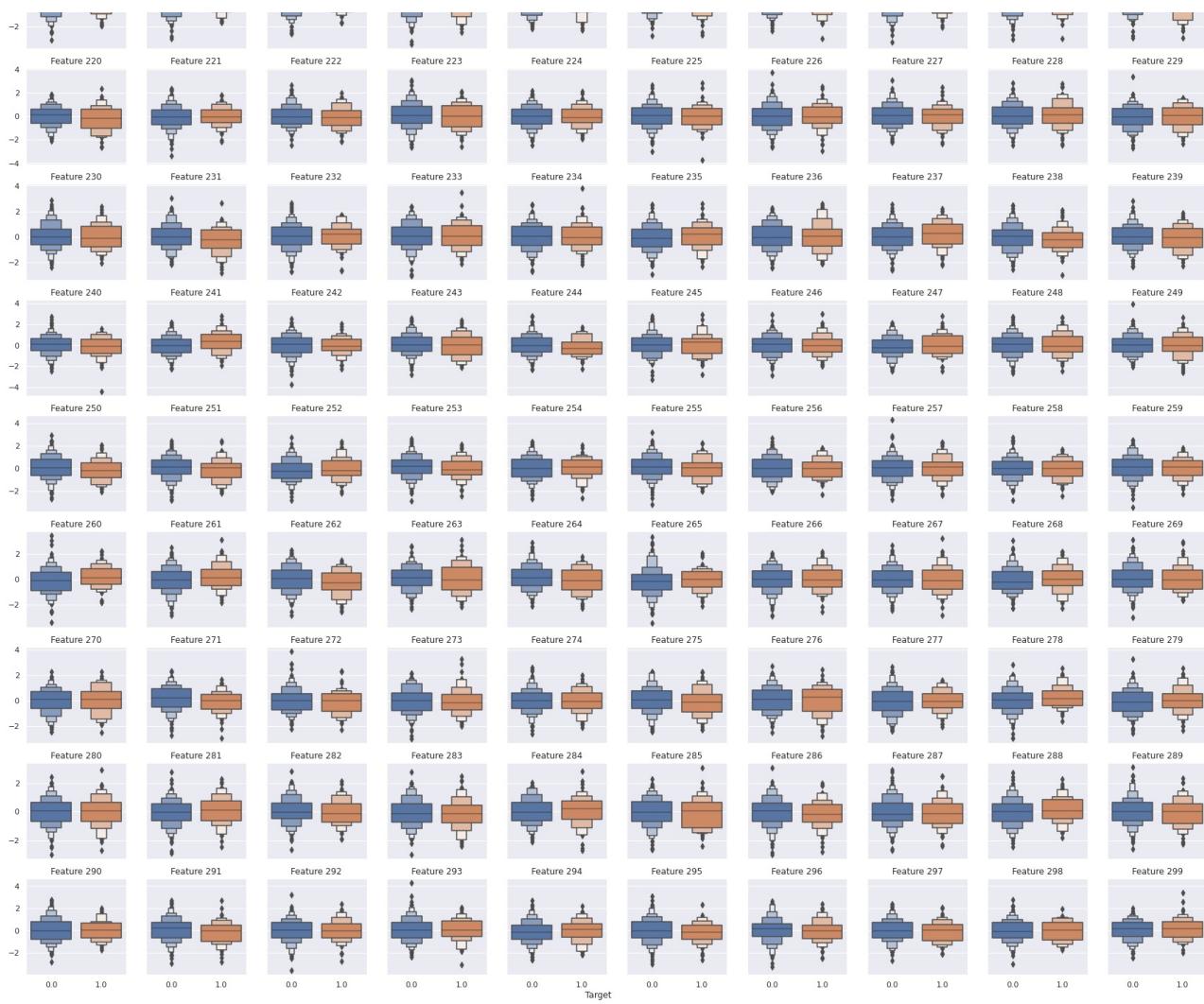
Box Plot of all 300 features in single frame

In []:

```
sns.set(rc={'figure.figsize':(30,90)})
fig, ax = plt.subplots(30, 10,sharex='col', sharey='row')
fig.text(0.5, 0.12, s='Target', ha='center',)
fig.text(0.1, 0.5, s='Distribution', va='center', rotation='vertical')
ax = ax.reshape(-1)
for i, col in enumerate(train_data.columns[2:]):
    axe= sns.boxenplot(np.array(train_data['target']),np.array(train_data[col]), ax=ax[i])
    ax[i].set_title('Feature '+col)
plt.show()
```







Observation

From the above two plot it is clearly mentioned that distribution of all 300 features are almost identical. All the features have almost same means and standard deviations

Seem like , Median of both target 0 and 1 have somewhat same. Boxplot detecting few points as outlier

In []:

```
###Scatter plot
r = np.random.randint(0,299,2,dtype=int)

print('***** scatter plot between {} and {}'.format(r[0],r[1]))

fig = plt.figure(figsize=(7,7))

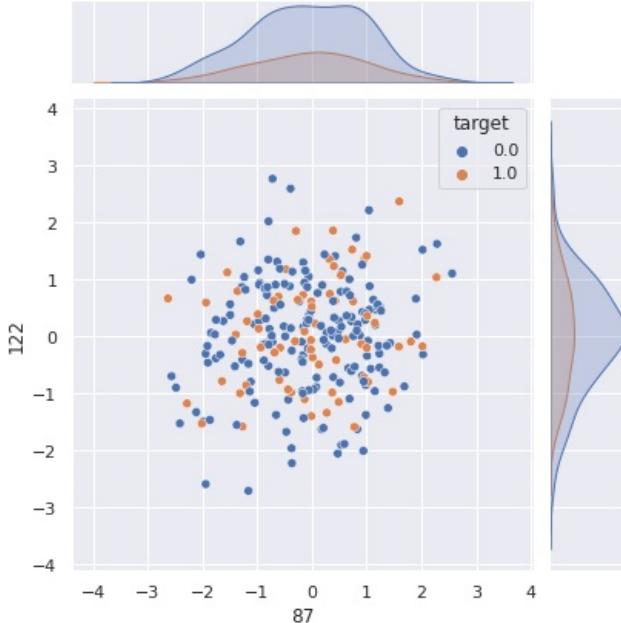
sns.jointplot(data=train_data,x = str(r[0]),y = str(r[1]),hue=train_data['target'])

***** scatter plot between 87 and 122 *****
```

Out[]:

<seaborn.axisgrid.JointGrid at 0x7f8f65615e50>

<Figure size 504x504 with 0 Axes>



We are unable to interpret using scatter plot

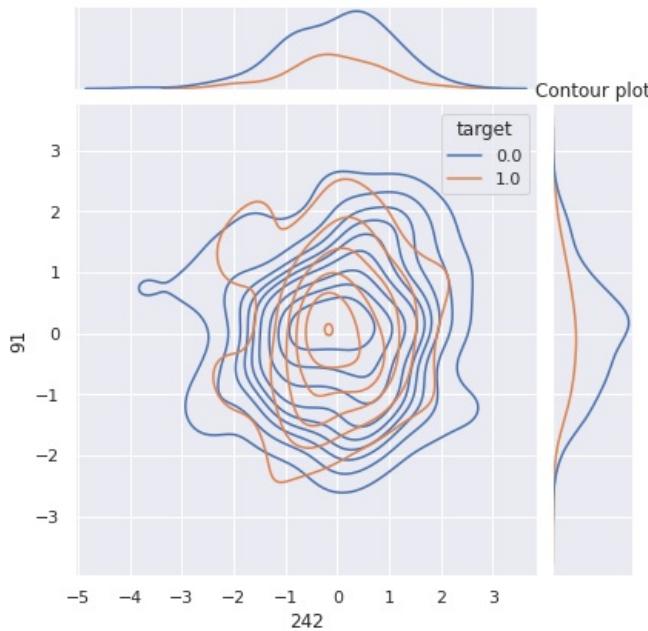
Let's observe the density of region rather than points. One of the seaborn toolkit called Contour

In []:

```
print('***** Contour plot between {} and {}'.format(r[0],r[1]))  
fig = plt.figure(figsize=(7,7))  
r = np.random.randint(0,299,2,dtype=int)  
sns.jointplot(x = str(r[0]),y = str(r[1]),data = train_data, hue=train_data['target'],kind = "kde")  
plt.title("Contour plot")  
plt.show()
```

***** Contour plot between 87 and 122 *****

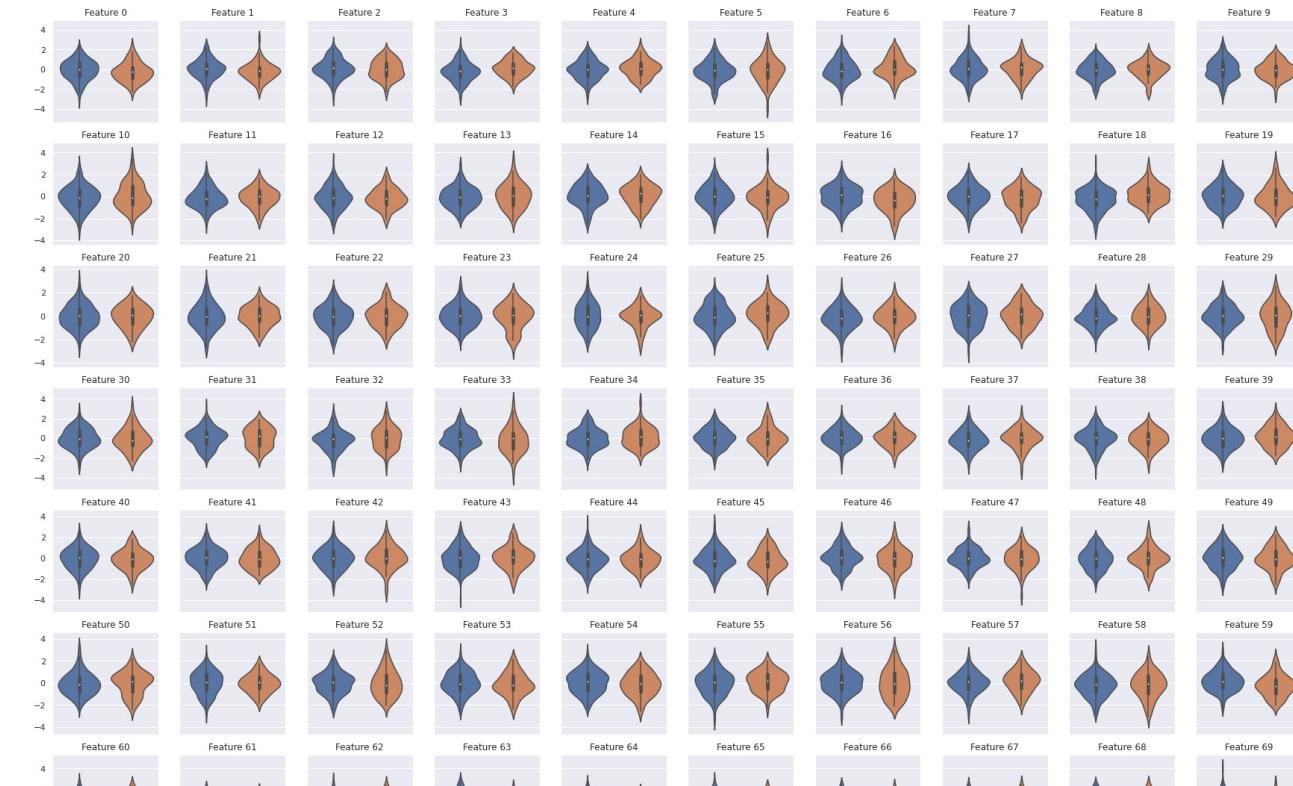
<Figure size 504x504 with 0 Axes>

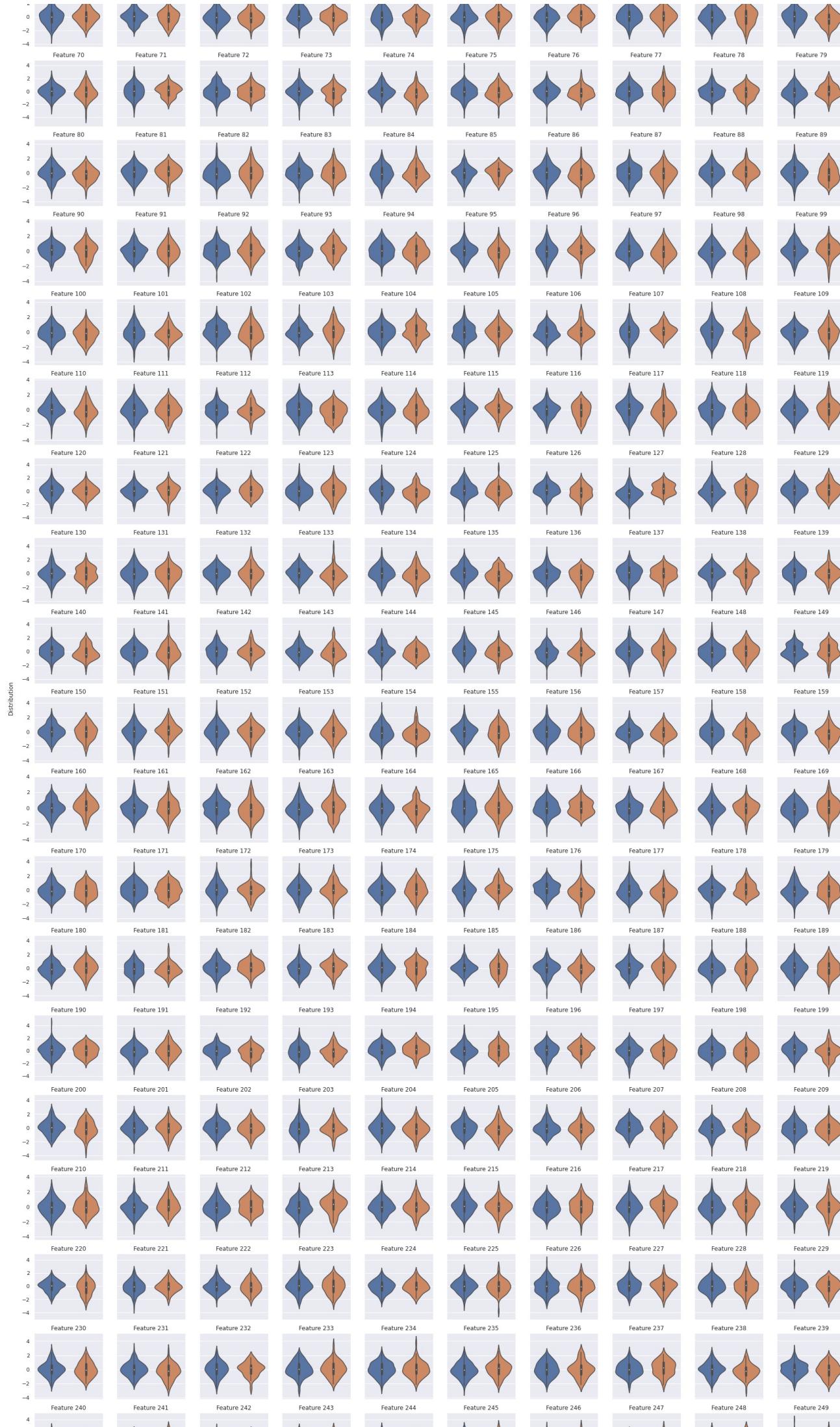


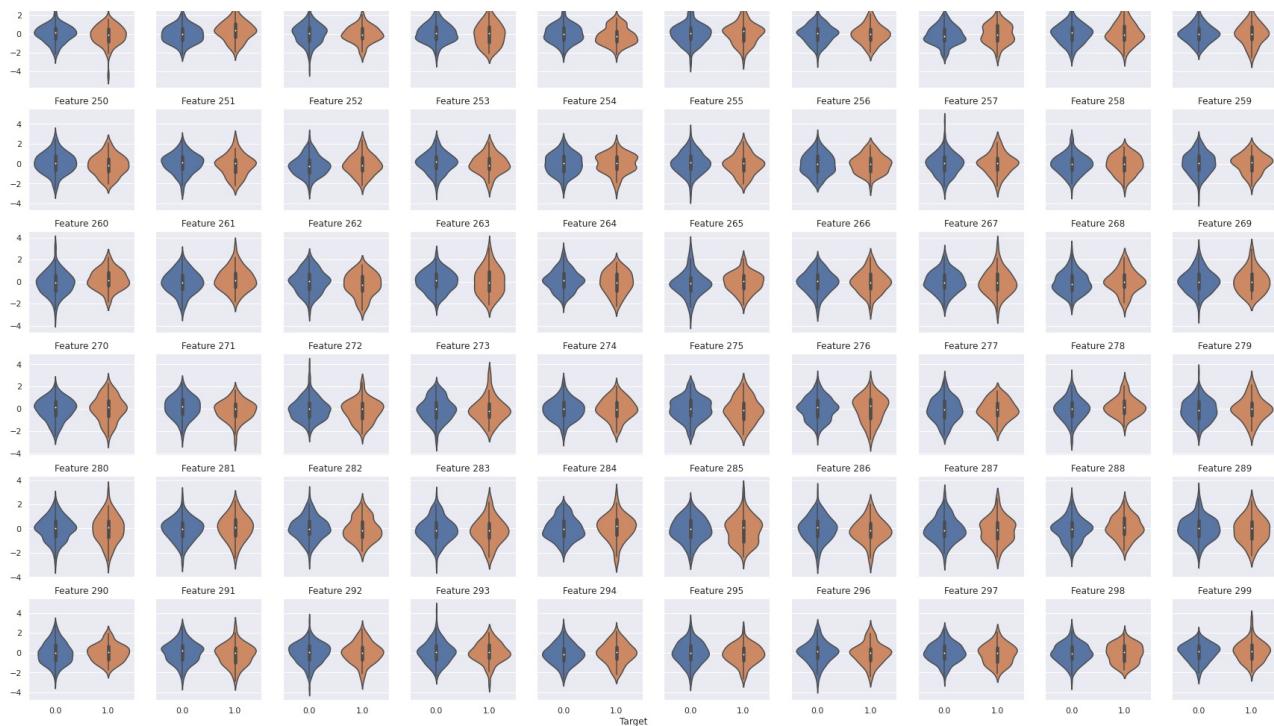
More density region is in range between -1 to 1.

In []:

```
sns.set(rc={'figure.figsize':(30,90)})  
fig, ax = plt.subplots(10, 10, sharex='col', sharey='row')  
fig.text(0.5, 0.12, s='Target', ha='center',)  
fig.text(0.1, 0.5, s='Distribution', va='center', rotation='vertical')  
ax = ax.reshape(-1)  
for i, col in enumerate(train_data.columns[2:]):  
    axe= sns.violinplot(np.array(train_data['target']),np.array(train_data[col]), ax=ax[i])  
    ax[i].set_title('Feature '+col)  
plt.show()
```







It seems like Median for both target value 0 and 1 have somewhat same.

New thing : I have tried Correlation Among Features till Feature Engineering

In []:

```
train_data.drop(['id'],axis=1).corr()
```

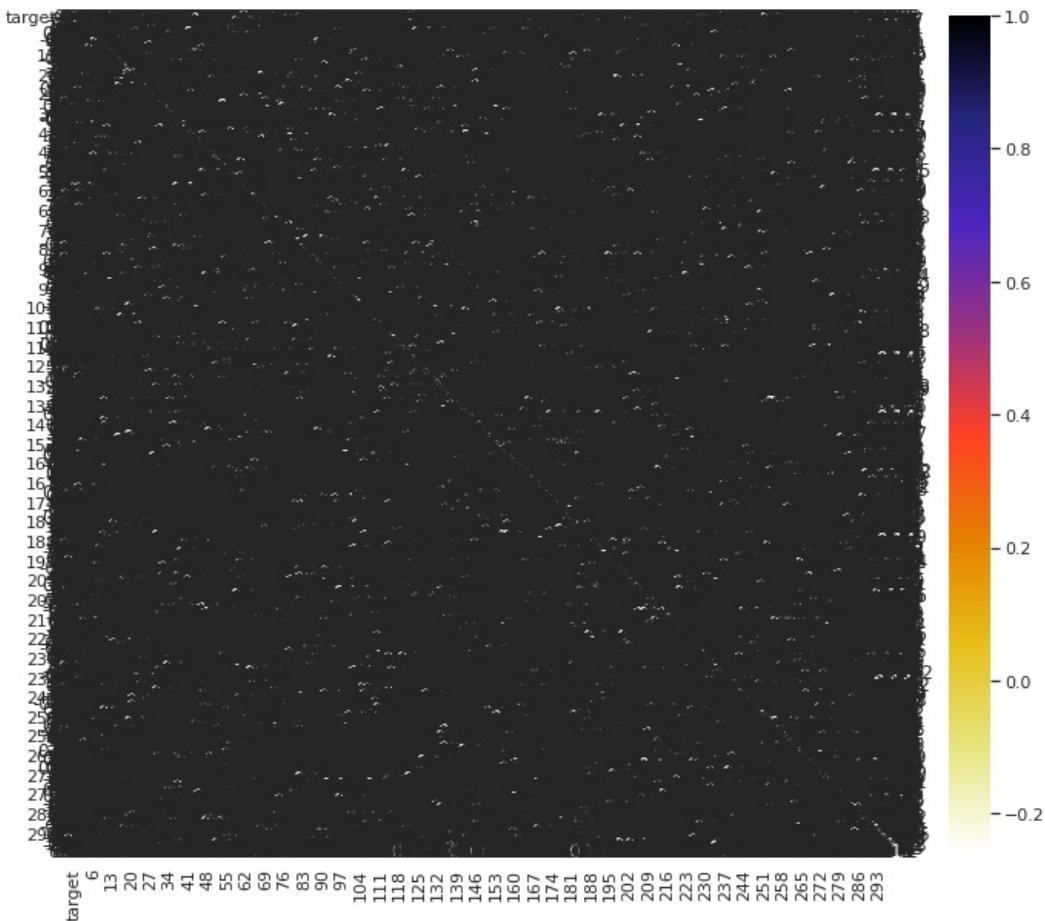
Out[]:

	target	0	1	2	3	4	5	6	7	8	9
target	1.000000	-0.085268	-0.111172	-0.036050	0.153317	0.077830	0.006206	0.090943	-0.003443	0.020330	-0.037787
0	-0.085268	1.000000	0.039939	0.069846	-0.139829	0.079360	-0.063259	-0.071990	0.016923	0.043081	0.051027
1	-0.111172	0.039939	1.000000	0.023237	-0.096873	-0.052313	0.048742	-0.086880	-0.013418	-0.005268	0.007490
2	-0.036050	0.069846	0.023237	1.000000	-0.107534	0.077051	-0.044687	-0.086411	0.000560	0.028217	-0.115959
3	0.153317	-0.139829	-0.096873	-0.107534	1.000000	-0.068026	-0.009757	0.049916	-0.063416	0.069039	0.047272
...
295	-0.071524	0.041362	-0.060408	0.151670	0.019118	-0.014687	0.070082	-0.069957	0.048203	0.113801	0.024786
296	-0.020375	0.009961	-0.192962	-0.030104	-0.049607	-0.002116	0.016213	-0.004527	-0.087313	-0.042409	-0.041123
297	-0.065313	-0.002739	0.025713	-0.025195	-0.008971	-0.031669	0.074662	0.007120	-0.047059	0.013366	-0.080101
298	-0.012973	0.092893	-0.049503	-0.086761	-0.052998	-0.044680	0.032147	-0.020821	-0.061363	0.019664	-0.045190
299	0.050805	0.065807	0.010546	-0.169051	0.076426	0.054352	0.064655	0.037515	-0.034875	0.060992	0.082712

301 rows x 301 columns

In []:

```
#using pearson correlation
plt.figure(figsize=(12,10))
cor = train_data.drop(['id'],axis=1).corr()
sns.heatmap(cor,annot=True,cmap = plt.cm.CMRmap_r)
plt.show()
```



From above heatmap, it is really difficult to get exact correlation value between the features

Top 20 correlated features

In []:

```
correlations = train_data.drop(['id'],axis=1).corr().drop_duplicates()
print("Top 20 positive correlated features with target")
print(correlations['target'].sort_values(ascending = False)[:20])

correlations = train_data.drop(['id'],axis=1).corr().drop_duplicates()
print("Top 20 Negative correlated features with target")
print(correlations['target'].sort_values(ascending = True)[:20])
```

Top 20 positive correlated features with target

```
target    1.000000
127     0.337540
18      0.206452
241     0.173879
3       0.153317
66      0.140056
93      0.136455
260     0.134119
213     0.121495
167     0.121217
175     0.120121
261     0.118305
278     0.118187
211     0.116670
151     0.110165
178     0.108901
169     0.106084
208     0.103583
196     0.101347
183     0.098502
```

Name: target, dtype: float64

Top 20 Negative correlated features with target

```
176    -0.217100
59     -0.203166
135    -0.179960
16     -0.179796
126    -0.167064
69     -0.164571
74     -0.157756
133    -0.147518
113    -0.146098
199    -0.145666
136    -0.140761
231    -0.137473
271    -0.132168
159    -0.128942
63     -0.128926
177    -0.123470
262    -0.119571
220    -0.117128
95     -0.116410
192    -0.114118
```

Name: target, dtype: float64

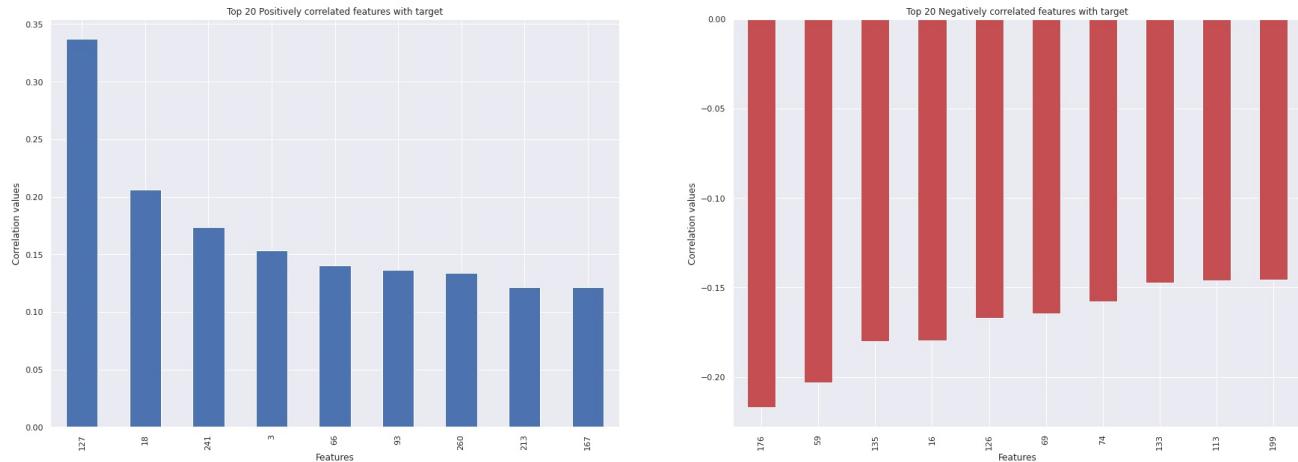
In []:

```
correlations = train_data.drop(['id'], axis=1).corr().unstack().drop_duplicates()
sns.set(rc={'figure.figsize':(30,10)})
plt.subplot(1,2,1)
x = correlations["target"].sort_values(ascending=False)[1:10].plot(kind='bar', title='Top 20 Positively correlated features with target')
x.set_xlabel('Features')
x.set_ylabel('Correlation values')

plt.subplot(1,2,2)
x = correlations["target"].sort_values(ascending=True)[:10].plot(kind='bar', title='Top 20 Negatively correlated features with target', color='r')
x.set_xlabel('Features')
x.set_ylabel('Correlation values')
```

Out[]:

Text(0, 0.5, 'Correlation values')

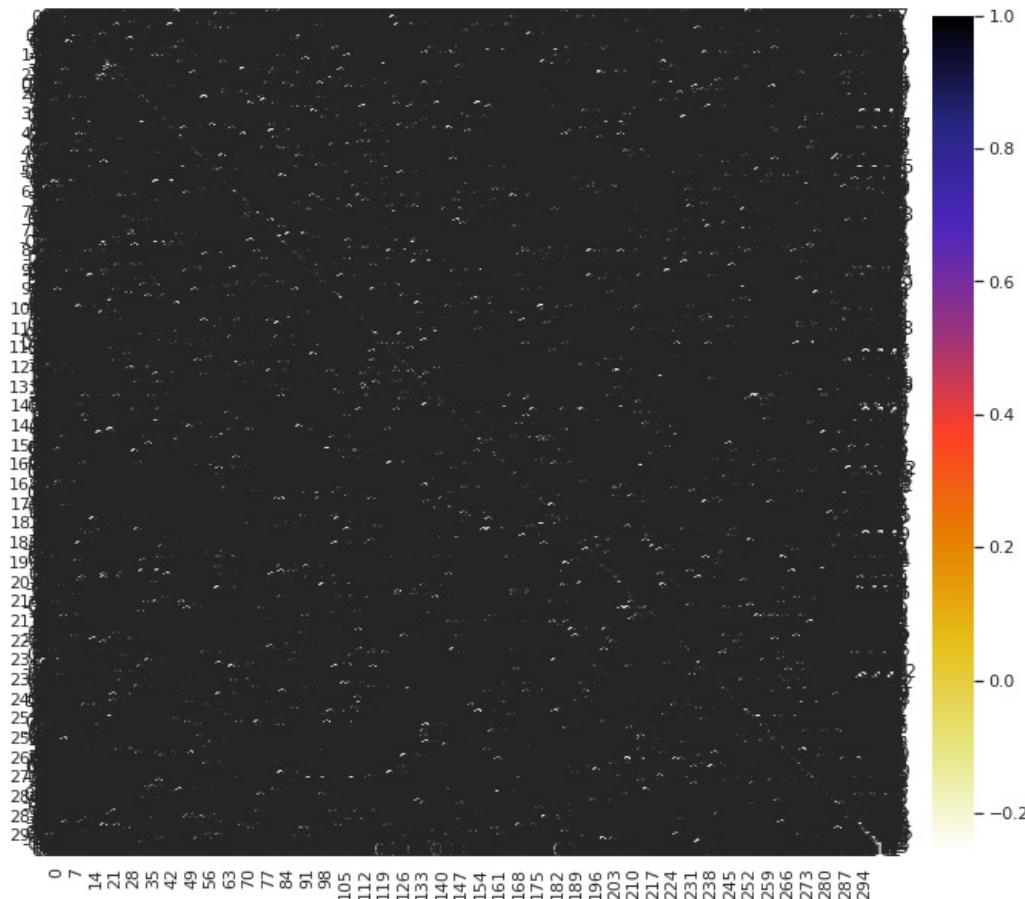


The Blue color is showing positive feature and red color negative features. Based on correlation values , we are see that features are not highly correlated with target. Feature 127 and 176 most correlated with target.

Correlation among 300 features

In []:

```
plt.figure(figsize=(12,10))
cor = train_data.drop(['id','target'],axis=1).corr()
sns.heatmap(cor,annot=True,cmap = plt.cm.CMRmap_r)
plt.show()
```



From above heatmap, it is really difficult to get exact correlation value between the features

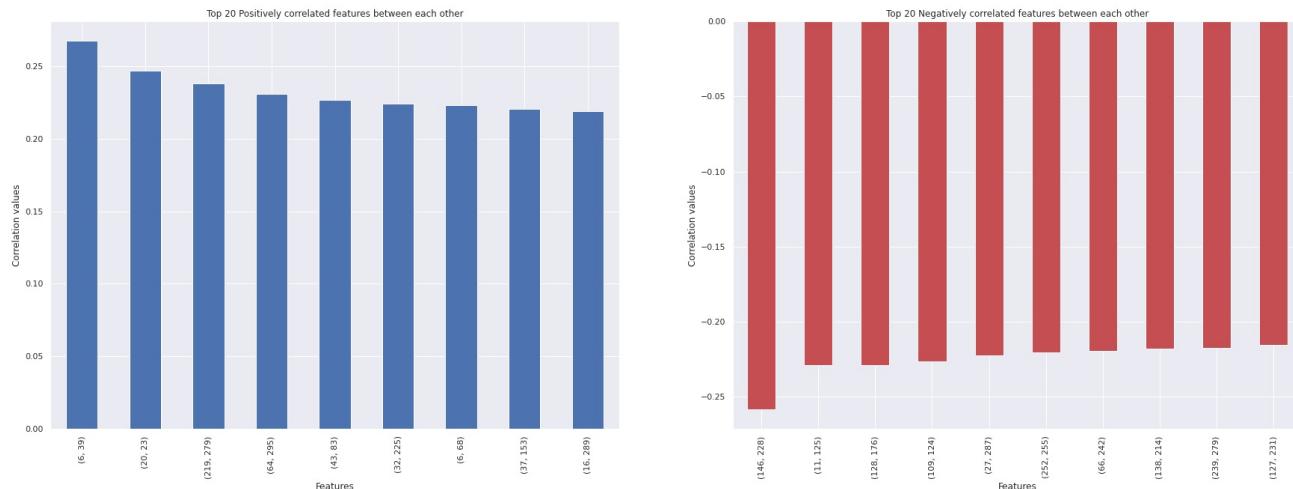
In []:

```
correlations = train_data.drop(['id','target'], axis=1).corr().unstack().drop_duplicates()
sns.set(rc={'figure.figsize':(30,10)})
plt.subplot(1,2,1)
x = correlations.sort_values(ascending=False)[1:10].plot(kind='bar', title='Top 20 Positively correlated features between each other')
x.set_xlabel('Features')
x.set_ylabel('Correlation values')

plt.subplot(1,2,2)
x = correlations.sort_values(ascending=True)[:10].plot(kind='bar', title='Top 20 Negatively correlated features between each other', color='r')
x.set_xlabel('Features')
x.set_ylabel('Correlation values')
```

Out[]:

Text(0, 0.5, 'Correlation values')



In []:

```
def correlation(dataset,threshold):
    col_corr = set()#Set of all name of correlated columns
    corr_matrix = dataset.corr()
    for i in range(len(corr_matrix.columns)):
        for j in range(i):
            if abs(corr_matrix.iloc[i,j])>threshold:
                # we are interested in absolute coeff value
                colname = corr_matrix.columns[i] #getting the name of column
                col_corr.add(colname)
    return col_corr
```

In []:

```
corr_features = correlation(train_data,0.5)
len(set(corr_features))
```

Out[]:

0

We have set threshold value for correlation value between feature is 0.5 but features are not highly correlated to each other. So we are not deleting any features here.

Feature Engineering

In []:

```
train_data = pd.read_csv('train.csv')
train_data.head()
```

Out[]:

	id	target	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	1
0	0	1.0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174	-0.24
1	1	0.0	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483	0.564
2	2	0.0	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587	1.464
3	3	1.0	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361	0.962
4	4	0.0	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697	-1.17

5 rows × 302 columns

In []:

```
#create duplicate
data = train_data.drop(['id','target'],axis=1)
data.head()
```

Out[]:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174	-0.245	-1.070
1	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483	0.564	0.047
2	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587	1.464	-1.550
3	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361	0.962	0.021
4	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697	-1.175	0.720

5 rows × 300 columns

Mean and Standard deviation value of each row

In []:

```
train_data['mean'] = np.mean(data, axis=1)
train_data['std'] = np.std(data, axis=1)
train_data.head()
```

Out[]:

	id	target	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	1.0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174	-0.
1	1	0.0	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483	0.564
2	2	0.0	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587	1.464
3	3	1.0	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361	0.962
4	4	0.0	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697	-1.

5 rows × 301 columns

In []:

```
#ref : : https://docs.scipy.org/doc/numpy/reference/routines.math.html)
sin_data = np.sin(data)
cos_data = np.cos(data)
tan_data = np.tan(data)

print("sin data:",sin_data)
print("cos_data:",cos_data)
print("tan_data:", tan_data)
```

```
sin data:      0      1      2 ...    297     298     299
0   -0.875756 -0.897470 -0.577775 ...  0.422317  0.818043 -0.733869
1   -0.738606  0.267695  0.989476 ... -0.028996 -0.935472 -0.336314
2    0.098838  0.983701 -0.668359 ... -0.147460 -0.601997  0.663135
3   -0.835477 -0.793172 -0.974166 ... -0.882429  0.720831 -0.256508
4    0.724976 -0.998091  0.498615 ...  0.991717 -0.382961 -0.599599
...
245  -0.067948 -0.182964 -0.913985 ...  0.928740  0.858327 -0.391260
246  -0.231870 -0.980502 -0.887362 ... -0.852631  0.679441  0.127651
247  -0.727446 -0.965561 -0.690370 ...  0.625678 -0.939099 -0.754571
248  -0.435866 -0.202588 -0.690370 ... -0.854503 -0.519273 -0.092866
249   0.663135  0.874304  0.971859 ... -0.992904 -0.196709  0.599599
```

[250 rows x 300 columns]

```
cos_data:      0      1      2 ...    297     298     299
0    0.482754  0.441075  0.816196 ...  0.906448  0.575157  0.679291
1    0.674137  0.963504 -0.144694 ...  0.999580 -0.353400  0.941750
2    0.995104  0.179813  0.743839 ...  0.989068  0.798498  0.748499
3    0.549526  0.608998  0.225831 ...  0.470446  0.693111 -0.966542
4    0.688774  0.061757  0.866824 ...  0.128441  0.923764  0.800301
...
245   0.997689  0.983120  0.405747 ...  0.370731  0.513103  0.920280
246   0.972747  0.196509 -0.461073 ...  0.522514  0.733730  0.991819
247  -0.686165 -0.260175  0.723457 ...  0.780082  0.343646  0.656219
248   0.900012  0.979264  0.723457 ... -0.519447  0.854608  0.995679
249   0.748499  0.485379  0.235562 ... -0.118922  0.980462  0.800301
```

[250 rows x 300 columns]

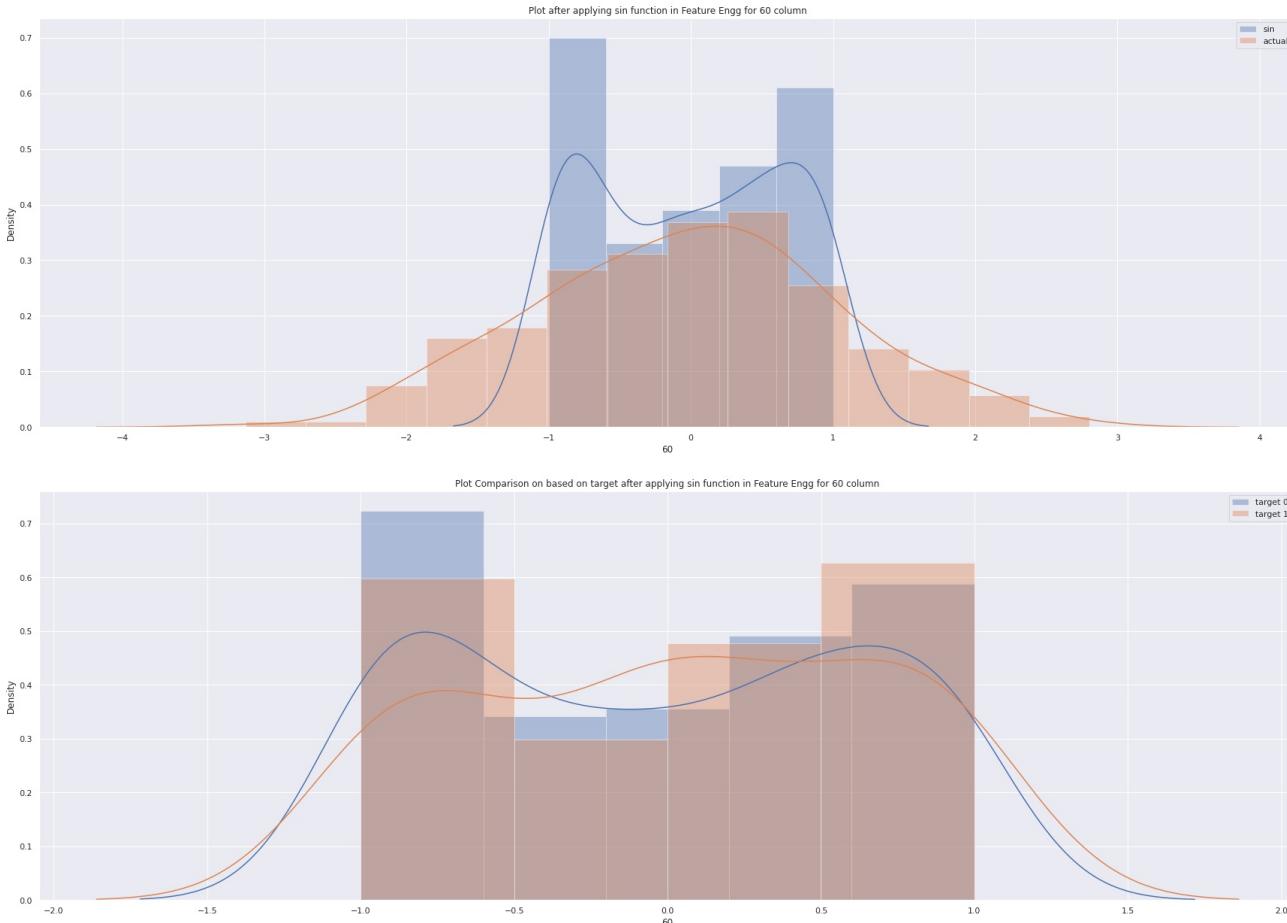
```
tan_data:      0      1      2 ...    297     298     299
0   -1.814085 -2.034733 -0.707888 ...  0.465903  1.422294 -1.080345
1   -1.095631  0.277835 -6.838409 ... -0.029008  2.647059 -0.357116
2    0.099325  5.470689 -0.898526 ... -0.149090 -0.753912  0.885953
3   -1.520360 -1.302422 -4.313690 ... -1.875727  1.039993  0.265387
4    1.052561 -16.161588  0.575221 ...  7.721217 -0.414566 -0.749217
...
245  -0.068105 -0.186105 -2.252597 ...  2.505157  1.672815 -0.425153
246  -0.238367 -4.989601  1.924561 ... -1.631787  0.926008  0.128704
247   1.060162  3.711197 -0.954265 ...  0.802067 -2.732754 -1.149878
248  -0.484289 -0.206878 -0.954265 ...  1.645024 -0.607615 -0.093269
249   0.885953  1.801282  4.125714 ...  8.349231 -0.200629  0.749217
```

[250 rows x 300 columns]

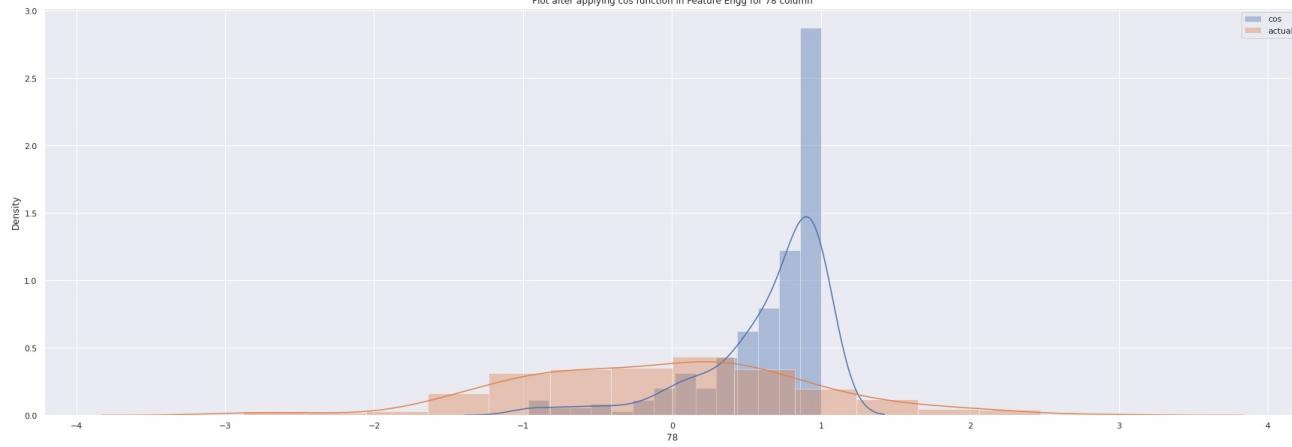
In []:

```
def feature_function(fe_name, fe_var):
    ...
    Parameter:
    fe_name: name of transformation Feature engg (string)
    fe_var: data after applying feature engineering
    Return:
    Plot 2 graphs.
    First plot for transformation of actual plot and after applying feature_engg
    ,>engg plot for any 'r' column
    Second plot for showing the plot of applying feature engg on the based on target
    ,>target value for any 'r' column
    ...
    r = str(np.random.randint(0,300))
    plt.figure(1)
    sns.distplot(fe_var[r], label=fe_name)
    sns.distplot(train_data[r], label='actual')
    plt.title('Plot after applying {} function in Feature Engg for {} column'.format(fe_name,r))
    plt.legend()
    plt.figure(2)
    sns.distplot(fe_var[train_data['target']==0][r], label='target 0')
    sns.distplot(fe_var[train_data['target']==1][r], label='target 1')
    plt.title('Plot Comparison on based on target after applying {} function in Feature Engg for {} column'.format(fe_name,r))
    plt.legend()
    plt.show()

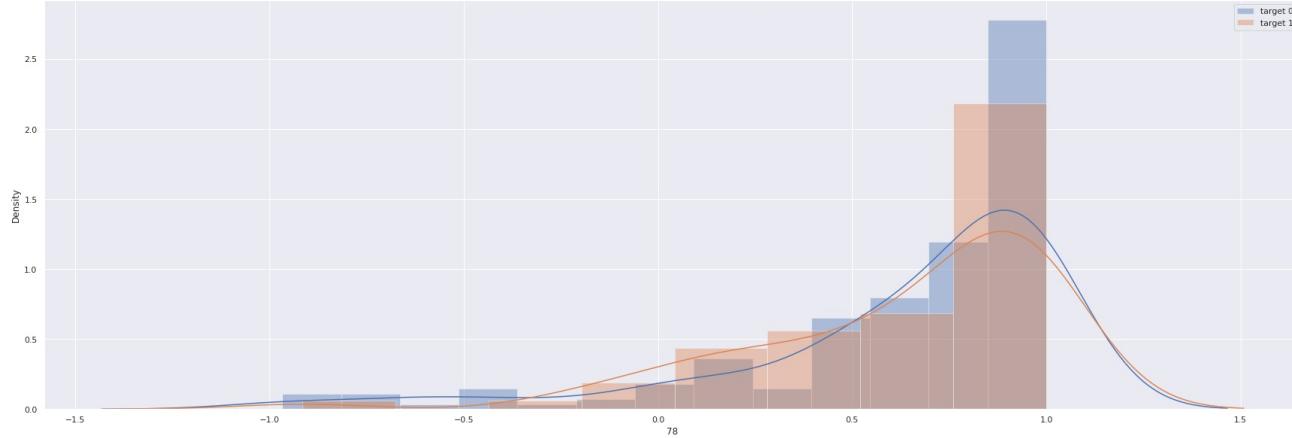
feature_function('sin',sin_data)
feature_function('cos',cos_data)
feature_function('tan',tan_data)
```



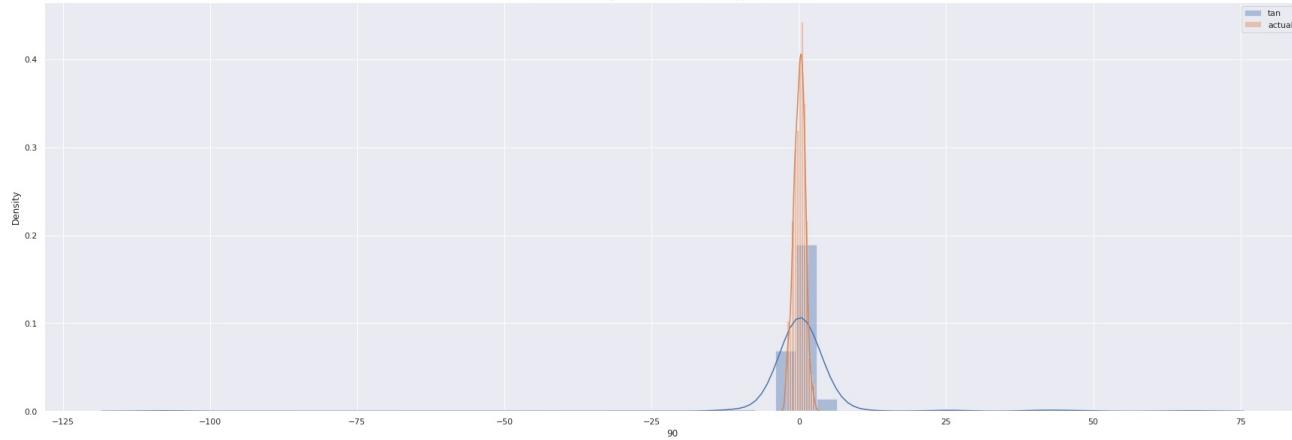
Plot after applying cos function in Feature Engg for 78 column



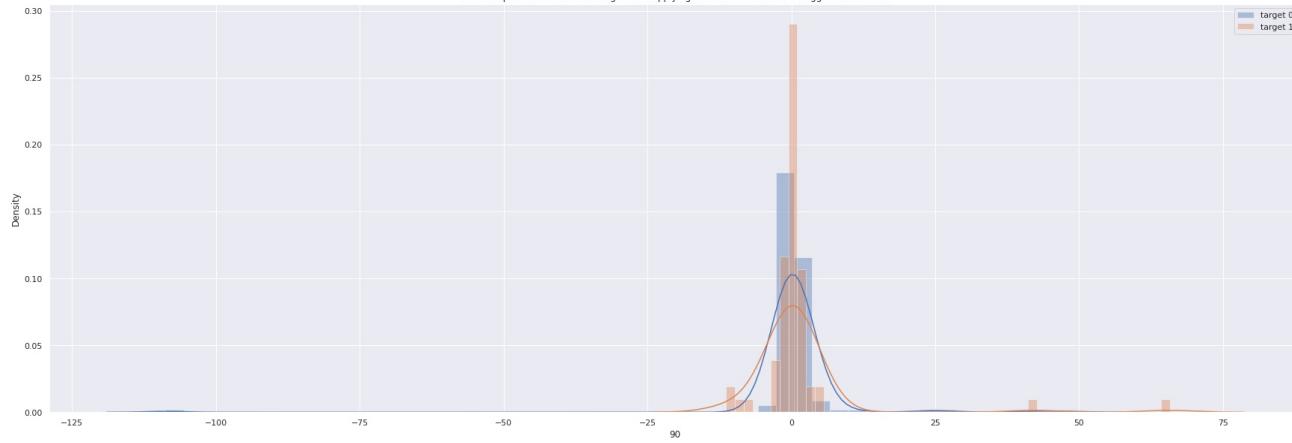
Plot Comparison on based on target after applying cos function in Feature Engg for 78 column



Plot after applying tan function in Feature Engg for 90 column



Plot Comparison on based on target after applying tan function in Feature Engg for 90 column



In []:

```
train_data['mean_sin'] = np.mean(sin_data, axis=1)
train_data['cos_sin'] = np.mean(cos_data, axis=1)
train_data['mean_tan'] = np.mean(tan_data, axis=1)
```

In []:

```
train_data.head()
```

Out[]:

	id	target	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	1
0	0	1.0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174	-0.24
1	1	0.0	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483	0.564
2	2	0.0	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587	1.464
3	3	1.0	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361	0.962
4	4	0.0	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697	-1.17

5 rows × 307 columns

We apply hyperbolic function

In []:

```
sinh_data = np.sinh(data)
cosh_data = np.cosh(data)
tanh_data = np.tanh(data)
print(sinh_data)
print(cosh_data)
print(tanh_data)
```

```
      0         1         2   ...     297        298        299
0 -1.281304 -1.359138 -0.655703   ...  0.449946  1.111409 -0.920463
1 -0.930000  0.274329  2.691226   ... -0.029004 -3.379222 -0.349765
2  0.099162  1.882887 -0.799144   ... -0.148541 -0.691878  0.790203
3 -1.158298 -1.049578 -1.784728   ... -1.304185  0.894804 14.980370
4  0.902872 -2.150538  0.546031   ...  1.996346 -0.403195 -0.688233
...   ...
245 -0.068052 -0.185040 -1.425997   ...  1.493226  1.225190 -0.412915
246 -0.236141 -1.846914 -3.819583   ... -1.207867  0.818436  0.128350
247 -5.074783 -3.049549 -0.837913   ...  0.728675 -1.545979 -0.963046
248 -0.466445 -0.205418 -0.837913   ... -4.092895 -0.573536 -0.093134
249  0.790203  1.276433  1.764359   ... -2.617481 -0.199296  0.688233
```

[250 rows × 300 columns]

```
      0         1         2   ...     297        298        299
0  1.625343  1.687382  1.195804   ...  1.096563  1.495069  1.359137
1  1.365613  1.036946  2.871009   ...  1.000421  3.524081  1.059403
2  1.004905  2.131963  1.280091   ...  1.010972  1.216016  1.274528
3  1.530247  1.449695  2.045790   ...  1.643441  1.341892 15.013710
4  1.347285  2.371669  1.139364   ...  2.232800  1.078224  1.213946
...   ...
245  1.002313  1.016976  1.741685   ...  1.797144  1.581484  1.081896
246  1.027503  2.100260  3.948318   ...  1.568102  1.292222  1.008203
247  5.172371  3.209323  1.304645   ...  1.237323  1.841209  1.388329
248  1.103436  1.020880  1.304645   ...  4.213287  1.152798  1.004328
249  1.274528  1.621506  2.028044   ...  2.802000  1.019666  1.213946
```

[250 rows × 300 columns]

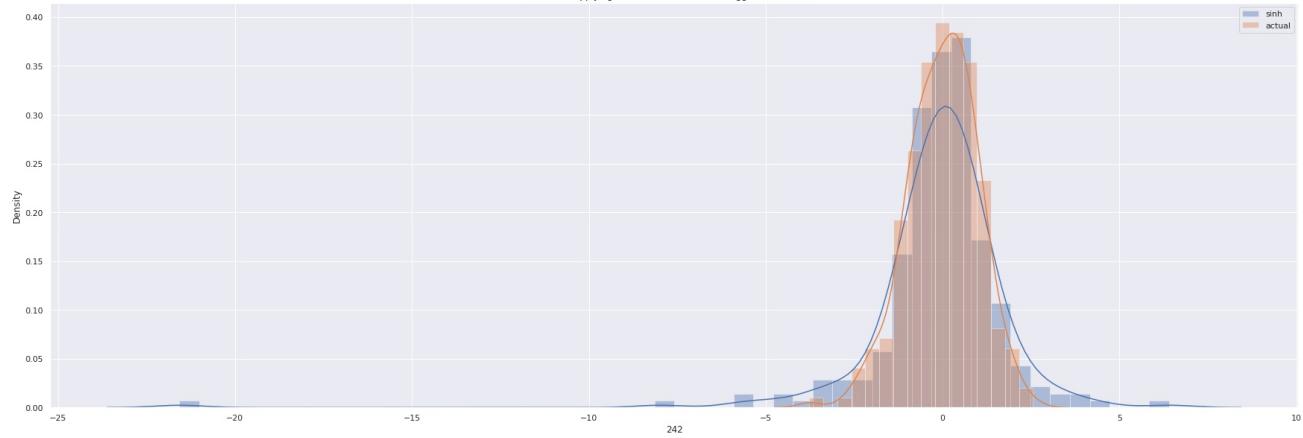
```
      0         1         2   ...     297        298        299
0 -0.788328 -0.805472 -0.548337   ...  0.410323  0.743383 -0.677241
1 -0.681013  0.264555  0.937380   ... -0.028992 -0.958895 -0.330153
2  0.098678  0.883171 -0.624287   ... -0.146929 -0.568971  0.619997
3 -0.756936 -0.724000 -0.872391   ... -0.793570  0.666823  0.997779
4  0.670142 -0.906761  0.479242   ...  0.894100 -0.373944 -0.566939
...   ...
245 -0.067895 -0.181951 -0.818745   ...  0.830889  0.774709 -0.381659
246 -0.229821 -0.879374 -0.967395   ... -0.770274  0.633356  0.127306
247 -0.981133 -0.950216 -0.642253   ...  0.588913 -0.839654 -0.693673
248 -0.422721 -0.201216 -0.642253   ... -0.971426 -0.497516 -0.092733
249  0.619997  0.787190  0.869981   ... -0.934147 -0.195452  0.566939
```

[250 rows × 300 columns]

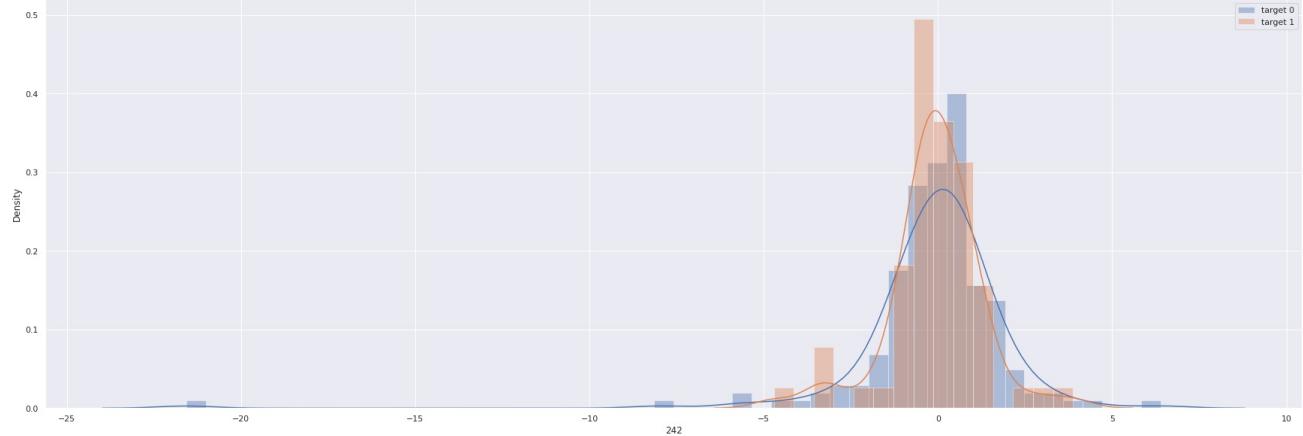
In []:

```
feature_function('sinh',sinh_data)
feature_function('cosh',cosh_data)
feature_function('tanh',tanh_data)
```

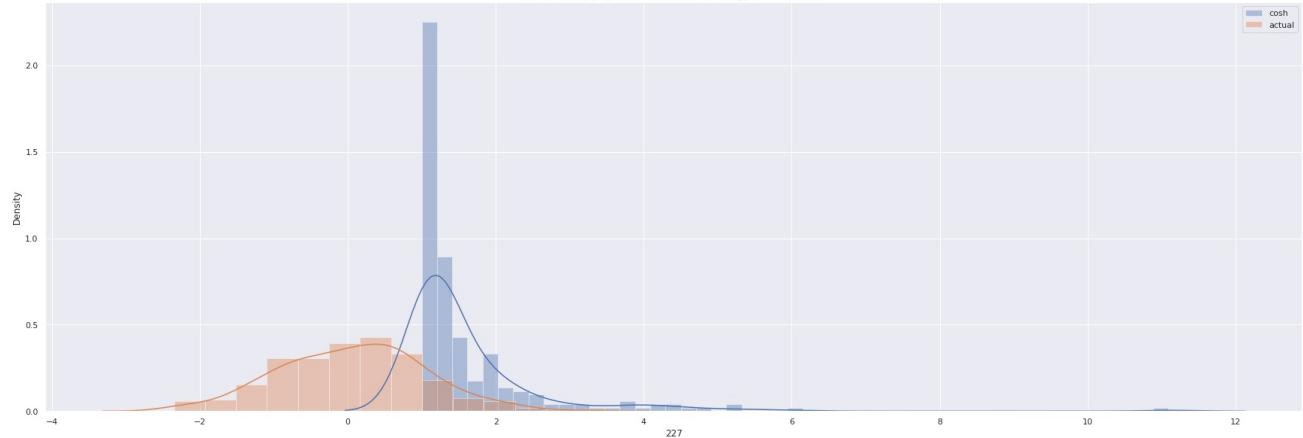
Plot after applying sinh function in Feature Engg for 242 column



Plot Comparison on based on target after applying sinh function in Feature Engg for 242 column

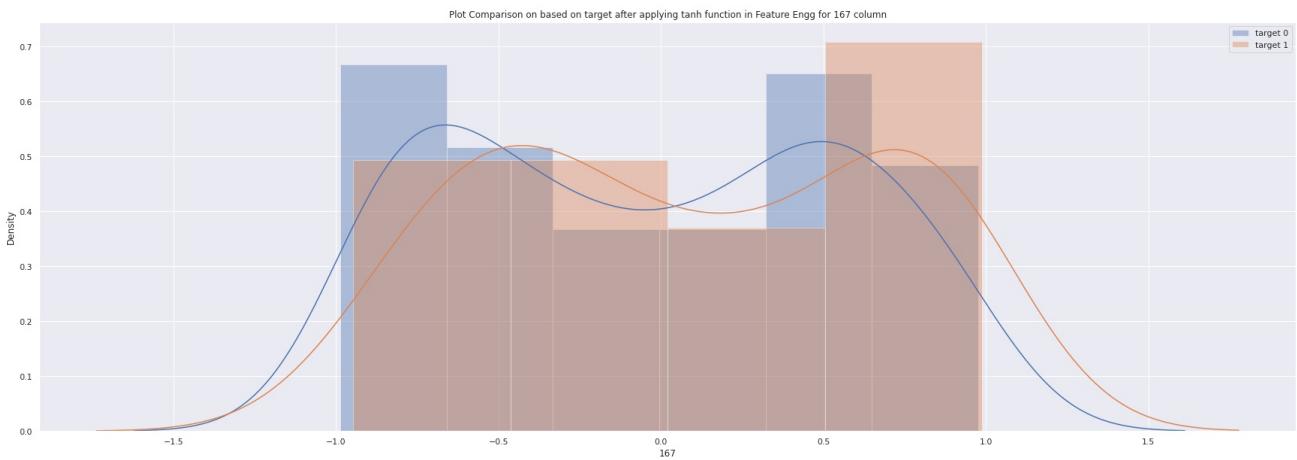
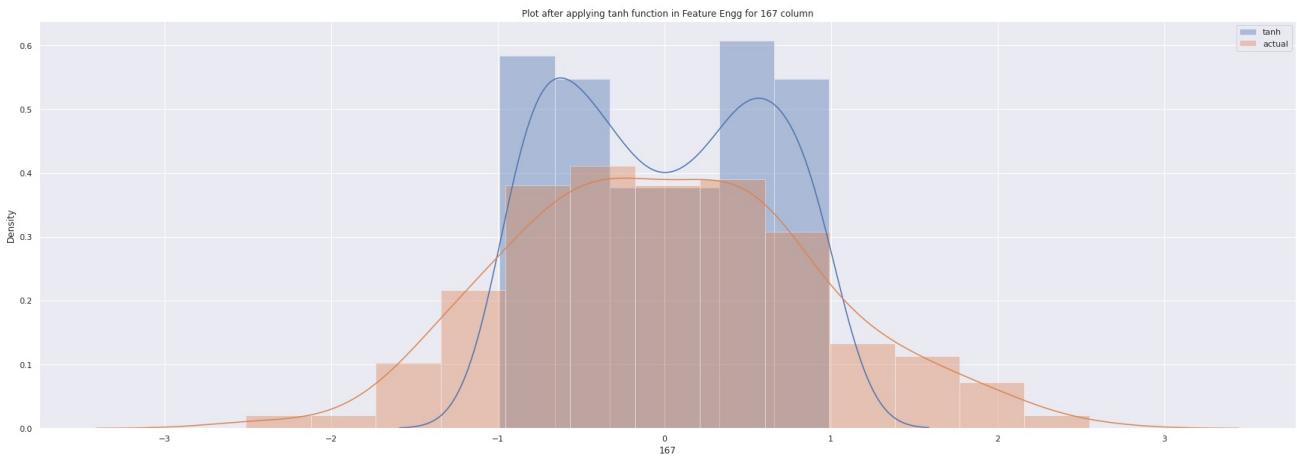


Plot after applying cosh function in Feature Engg for 227 column



Plot Comparison on based on target after applying cosh function in Feature Engg for 227 column





In []:

```
train_data['mean_sinh'] = np.mean(sinh_data, axis=1)
train_data['mean_cosh'] = np.mean(cosh_data, axis=1)
train_data['mean_tanh'] = np.mean(tanh_data, axis=1)
```

In []:

```
train_data.head()
```

Out[]:

	id	target	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	
0	0	1.0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174	-0.
1	1	0.0	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483	0.5
2	2	0.0	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587	1.4
3	3	1.0	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361	0.9
4	4	0.0	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697	-1.

5 rows x 310 columns

In []:

```
# We have taken ref https://numpy.org/doc/stable/reference/routines.math.html
```

```
exp_data = np.exp(data)
expm1_data = np.expm1(data)
exp2_data = np.exp2(data)
```

In []:

```
print(exp_data)
print(expm1_data)
print(exp2_data)
```

```
          0         1         2   ...      297        298        299
0  0.344039  0.328243  0.540101   ...  1.546509  2.606478  0.438673
1  0.435613  1.311275  5.562235   ...  0.971416  0.144858  0.709638
2  1.104066  4.014850  0.480946   ...  0.862431  0.524138  2.064731
3  0.371948  0.400116  0.261061   ...  0.339256  2.236696  29.994079
4  2.250157  0.221131  1.685395   ...  4.229146  0.675029  0.525713
..  ...  ...  ...  ...  ...  ...  ...
245 0.934260  0.831936  0.315688   ...  3.290370  2.806674  0.668981
246 0.791362  0.253346  0.128735   ...  0.360235  2.110659  1.136553
247 0.097588  0.159773  0.466732   ...  1.965998  0.295230  0.425283
248 0.636991  0.815462  0.466732   ...  0.120392  0.579262  0.911194
249 2.064731  2.897940  3.792404   ...  0.184520  0.820370  1.902179
```

[250 rows x 300 columns]

```
          0         1         2   ...      297        298        299
0  -0.655961 -0.671757 -0.459899   ...  0.546509  1.606478 -0.561327
1  -0.564387  0.311275  4.562235   ... -0.028584 -0.855142 -0.290362
2  0.104066  3.014850 -0.519054   ... -0.137569 -0.475862  1.064731
3  -0.628052 -0.599884 -0.738939   ... -0.660744  1.236696  28.994079
4  1.250157 -0.778869  0.685395   ...  3.229146 -0.324971 -0.474287
..  ...  ...  ...  ...  ...  ...
245 -0.065740 -0.168064 -0.684312   ...  2.290370  1.806674 -0.331019
246 -0.208638 -0.746654 -0.871265   ... -0.639765  1.110659  0.136553
247 -0.902412 -0.840227 -0.533268   ...  0.965998 -0.704770 -0.574717
248 -0.363009 -0.184538 -0.533268   ... -0.879608 -0.420738 -0.088806
249 1.064731  1.897940  2.792404   ... -0.815480 -0.179630  0.902179
```

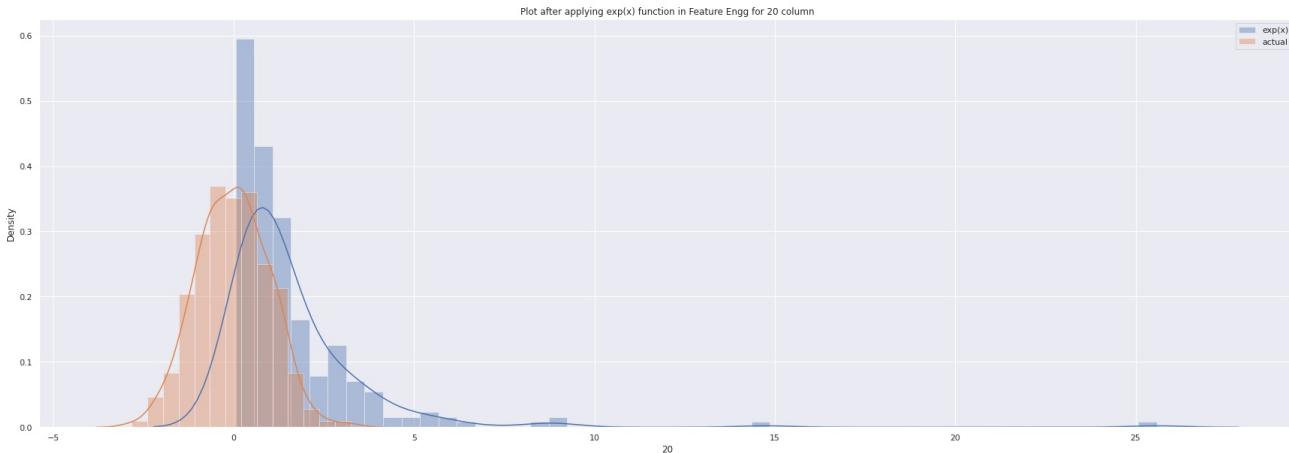
[250 rows x 300 columns]

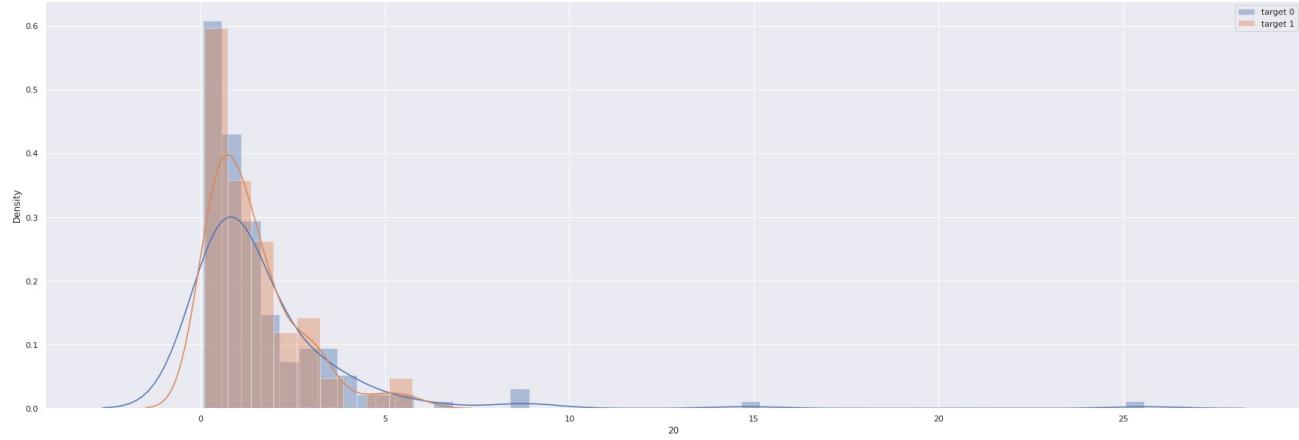
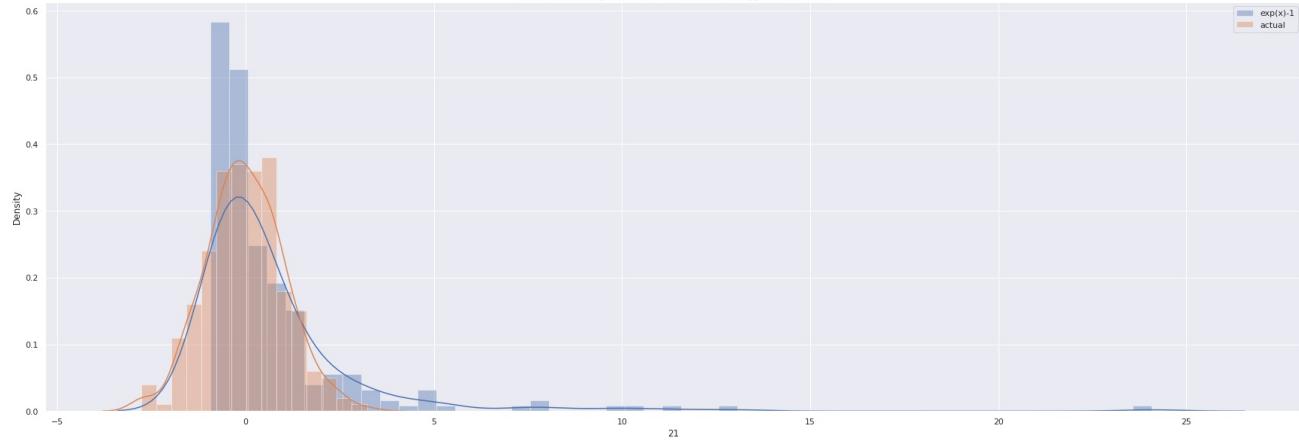
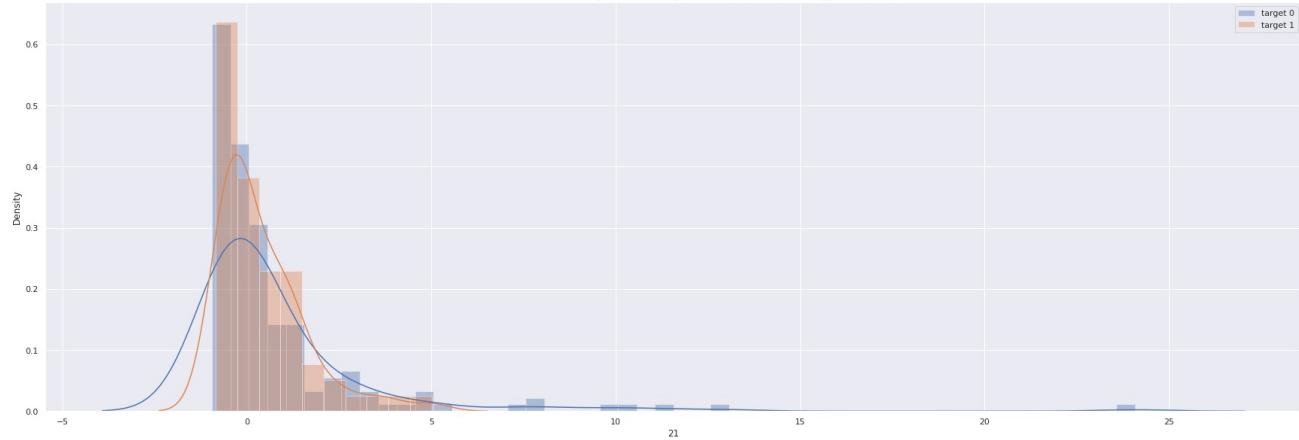
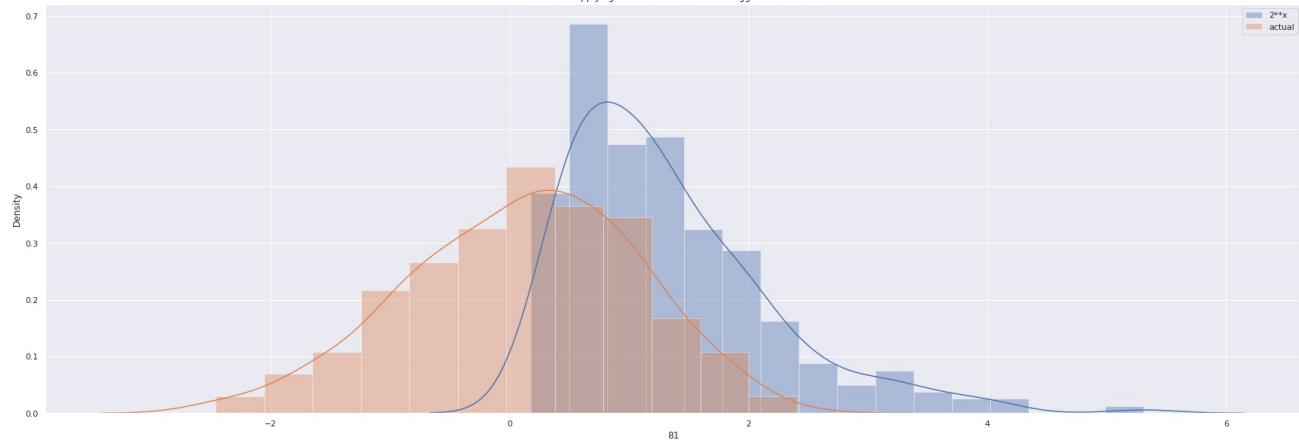
```
          0         1         2   ...      297        298        299
0  0.477311  0.462011  0.652477   ...  1.352848  1.942615  0.564874
1  0.562139  1.206644  3.285243   ...  0.980099  0.262066  0.788400
2  1.071031  2.620787  0.602069   ...  0.902501  0.639050  1.652901
3  0.503827  0.529976  0.394200   ...  0.472701  1.747146  10.563383
4  1.754427  0.351355  1.435945   ...  2.716973  0.761544  0.640380
..  ...  ...  ...  ...  ...  ...
245 0.953960  0.880259  0.449689   ...  2.283109  2.044857  0.756808
246 0.850274  0.386088  0.241484   ...  0.492775  1.678299  1.092778
247 0.199298  0.280486  0.589678   ...  1.597704  0.429283  0.552865
248 0.731536  0.868140  0.589678   ...  0.230526  0.684916  0.937571
249 1.652901  2.090720  2.519260   ...  0.309927  0.871758  1.561573
```

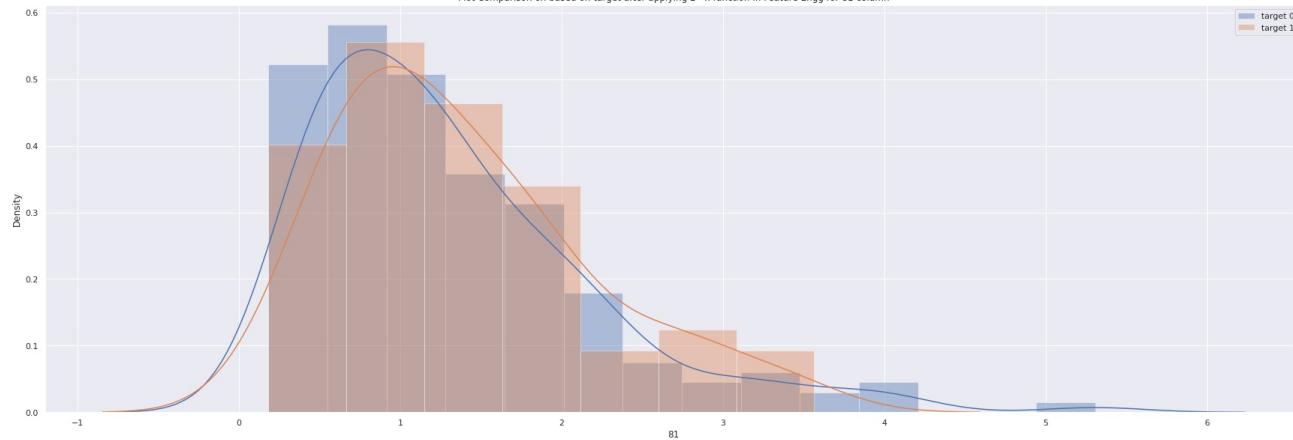
[250 rows x 300 columns]

In []:

```
feature_function('exp(x)',exp_data)
feature_function('exp(x)-1',expm1_data)
feature_function('2**x',exp2_data)
```



Plot Comparison on based on target after applying $\exp(x)$ in Feature Engg for 20 columnPlot after applying $\exp(x)-1$ function in Feature Engg for 21 columnPlot Comparison on based on target after applying $\exp(x)-1$ function in Feature Engg for 21 columnPlot after applying 2^{**x} function in Feature Engg for 81 column

Plot Comparison on based on target after applying 2^{xx} function in Feature Engg for 81 column

In []:

```
train_data['mean_exp'] = np.mean(exp_data, axis=1)
train_data['mean_expm1'] = np.mean(expm1_data, axis=1)
train_data['mean_exp2'] = np.mean(exp2_data, axis=1)
```

In []:

```
train_data.head()
```

Out[]:

	id	target	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
0	0	1.0	-1.067	-1.114	-0.616	0.376	1.090	0.467	-0.422	0.460	-0.443	-0.338	0.416	-2.177	-0.326	0.340	1.174
1	1	0.0	-0.831	0.271	1.716	1.096	1.731	-0.197	1.904	-0.265	0.557	1.202	0.542	0.424	-1.572	-0.968	-1.483
2	2	0.0	0.099	1.390	-0.732	-1.065	0.005	-0.081	-1.450	0.317	-0.624	-0.017	-0.665	1.905	0.376	-1.373	1.587
3	3	1.0	-0.989	-0.916	-1.343	0.145	0.543	0.636	1.127	0.189	-0.118	-0.638	0.760	-0.360	-2.048	-0.996	-0.361
4	4	0.0	0.811	-1.509	0.522	-0.360	-0.220	-0.959	0.334	-0.566	-0.656	-0.499	-0.653	-0.058	-0.046	0.654	-0.697

5 rows x 313 columns