# 数据摘要

## 1.标称属性

此数据集中的标称属性包括
['Date', 'GameID', 'Drive', 'qtr', 'down', 'time', 'TimeUnder', 'TimeSecs', 'SideofField', 'GoalToGo', 'FirstDown', 'posteam', 'DefensiveTeam', 'desc', 'PlayAttempted', 'sp', 'Touchdown', 'ExPointResult', 'TwoPointConv', 'DefTwoPoint', 'Safety', 'Onsidekick', 'PuntResult', 'PlayType', 'Passer', 'Passer_ID', 'PassAttempt', 'PassOutcome', 'PassLength', 'QBHit', 'PassLocation', 'InterceptionThrown', 'Interceptor','Rusher', 'Rusher_ID', 'RushAttempt', 'RunLocation', 'RunGap', 'Receiver', 'Receiver_ID', 'Reception', 'ReturnResult', 'Returner', 'BlockingPlayer', 'Tackler1', 'Tackler2', 'FieldGoalResult', 'Fumble', 'RecFumbTeam', 'RecFumbPlayer', 'Sack', 'Challenge.Replay', 'ChalReplayResult', 'Accepted.Penalty', 'PenalizedTeam', 'PenaltyType', 'PenalizedPlayer', 'HomeTeam', 'AwayTeam', 'Timeout_Indicator', 'Timeout_Team', 'Season']

数据读取：

```python
path = "/Users/margo/Downloads/NFL Play by Play 2009-2017 (v4).csv"
data = pd.read_csv(path)
```

对于标称属性，求出每个可能值的频数并保存在txt文件中，代码如下：

```python
def frequency(NominalAttribute, data):
    for n in NominalAttribute:
        doc = open('frequencyOfNominalAttribute.txt', 'a')
        print('*************** %s ***************' % n, file=doc)
        print(data[n].value_counts(), file=doc)
        doc.close()
```

求出的部分结果如下：

## 2.数值属性

此数据集中的数值属性包括:

['PlayTimeDiff', 'yrdln', 'yrdline100', 'ydstogo', 'ydsnet', 'Yards.Gained', 'AirYards', 'YardsAfterCatch', 'FieldGoalDistance', 'Penalty.Yards', 'PosTeamScore', 'DefTeamScore', 'ScoreDiff', 'AbsScoreDiff', 'posteam_timeouts_pre', 'HomeTimeouts_Remaining_Pre', 'AwayTimeouts_Remaining_Pre', 'HomeTimeouts_Remaining_Post','AwayTimeouts_Remaining_Post', 'No_Score_Prob', 'Opp_Field_Goal_Prob', 'Opp_Safety_Prob', 'Opp_Touchdown_Prob', 'Field_Goal_Prob', 'Safety_Prob', 'Touchdown_Prob', 'ExPoint_Prob', 'TwoPoint_Prob', 'ExpPts', 'EPA', 'airEPA', 'yacEPA', 'Home_WP_pre', 'Away_WP_pre', 'Home_WP_post', 'Away_WP_post', 'Win_Prob', 'WPA', 'airWPA', 'yacWPA']

对于数值属性,求出了每个数值属性的最大值、最小值、均值、中位数、四分位数及缺失值的个数,代码如下:

```python
def describeNumericalAttribute(NumericalAttribute, data):
    for n in NumericalAttribute:
        doc = open("describeNumericalAttribute.txt", 'a')
        print('*************** %s ***************' % n, file=doc)
        print("max:%s" % (data[n].max()), file=doc)
        print("min:%s" % (data[n].min()), file=doc)
        print("mean:%s" % (data[n].mean()), file=doc)
        print("median:%s" % (data[n].median()), file=doc)
        print("quantile1:%s" % (data[n].quantile(0.25)), file=doc)
        print("quantile2:%s" % (data[n].quantile(0.5)), file=doc)
        print("quantile3:%s" % (data[n].quantile(0.75)), file=doc)
        print("theNumberOfNull:%s" % (data[n].count()), file=doc)
        doc.close
```

求出的部分结果如下所示:

```
*************** PlayTimeDiff ***************
max:943.0
min:0.0
mean:20.576762334128926
median:17.0
quantile1:5.0
quantile2:17.0
quantile3:37.0
theNumberOfNull:407244
*************** yrdln ***************
max:50.0
min:1.0
mean:28.48832733600755
median:30.0
quantile1:20.0
quantile2:30.0
quantile3:39.0
theNumberOfNull:406848
*************** yrdline100 ***************
max:99.0
min:1.0
mean:48.644080836086204
median:49.0
quantile1:30.0
quantile2:49.0
quantile3:70.0
theNumberOfNull:406848
```
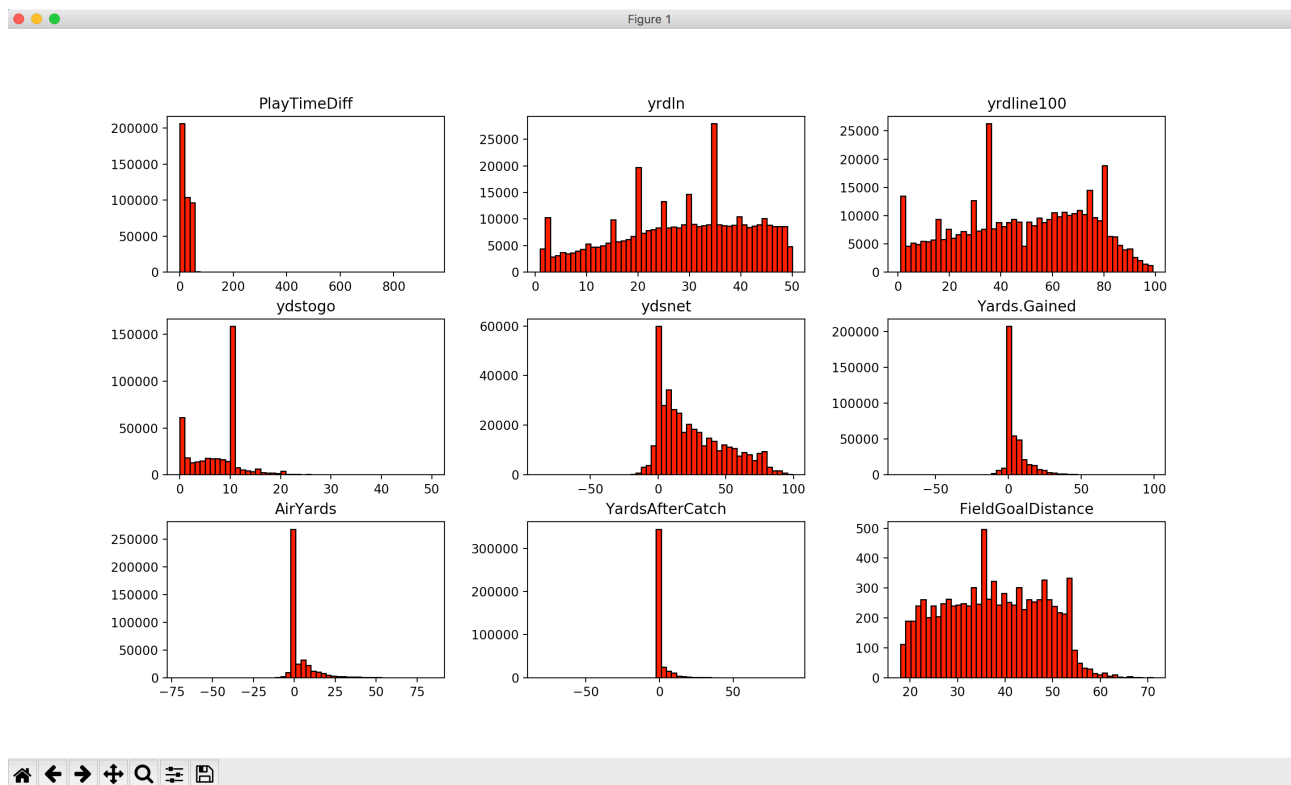
# 数据的可视化

## 1.直方图

对每种数值属性绘制直方图，代码如下

```python
def histogram(NumericalAttribute, data):
    for i, col in enumerate(NumericalAttribute):
        if i % cellSize == 0:    # 一页有cellSize张子图，如果第i列是第cellSize+1张图，就另起一页
            fig = plt.figure(FIGSIZE=(15, 15))
        ax = fig.add_subplot(rowSize, colSize, (i % cellSize) + 1)
        data[col].hist(ax=ax, grid=False,
                       bins=50, facecolor='red', edgecolor='black')
        plt.title(col)
        if (i + 1) % cellSize == 0 or i + 1 == len(NumericalAttribute):
            plt.subplots_adjust(wspace=0.3, hspace=0.3)
            plt.show()
```
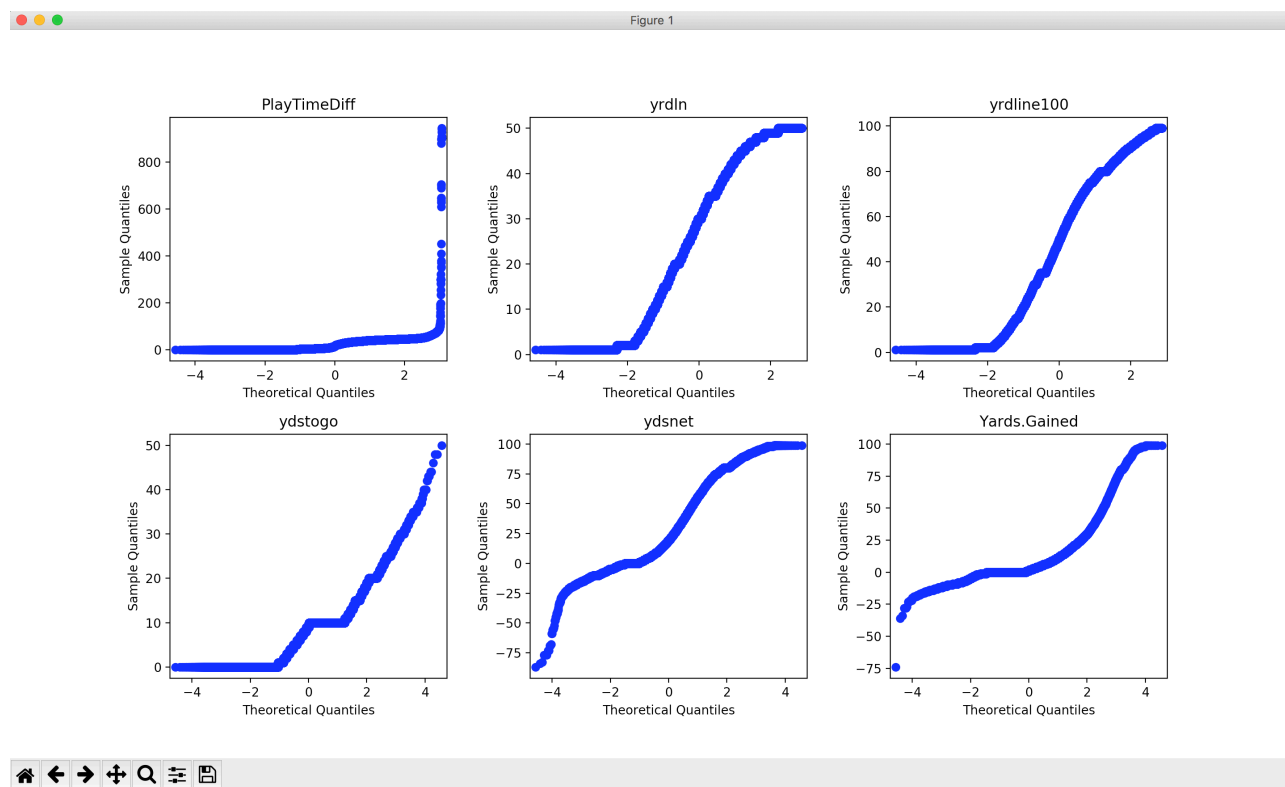
绘制出的部分直方图如下所示：



## 2.qq图

绘制qq图的代码如下所示：

```python
def qq(NumericalAttribute, data):
    for i, col in enumerate(NumericalAttribute):
        if i % 6 == 0:
            fig = plt.figure()
        ax = fig.add_subplot(2, 3, (i % 6) + 1)
        sm.qqplot(data[col], ax=ax)
        ax.set_title(col)
        if (i + 1) % 6 == 0 or i + 1 == len(NumericalAttribute):
            plt.subplots_adjust(wspace=0.3, hspace=0.3)
            plt.show()
```
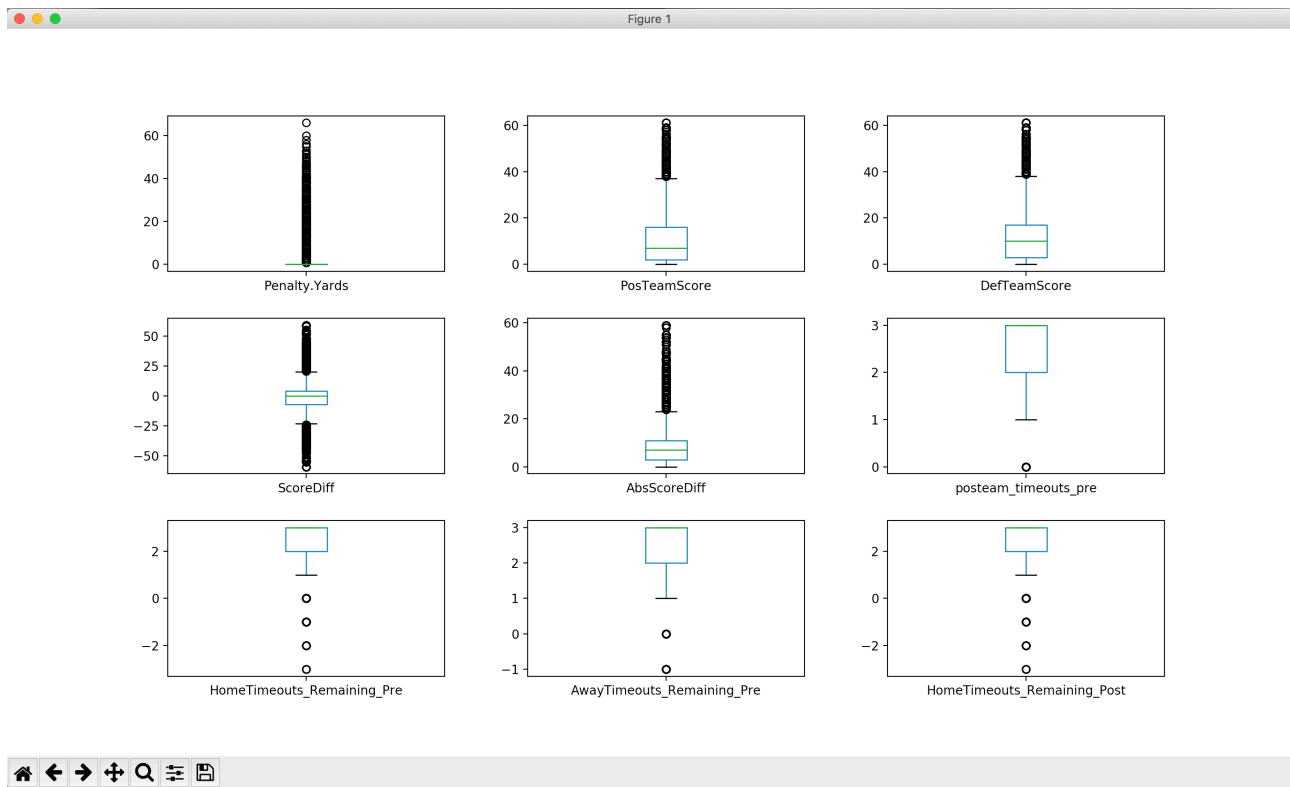
绘制出的部分qq图如下，可以直观地利用qq图检验分布是否为正态分布



## 3.盒图

绘制盒图的代码如下：

```python
def boxplot(NumericalAttribute, data):
    for i, col in enumerate(NumericalAttribute):
        if i % cellSize == 0:
            fig = plt.figure()
        ax = fig.add_subplot(colSize, rowSize, (i % cellSize) + 1)
        data[col].plot.box(ax=ax)
        if (i + 1) % cellSize == 0 or i + 1 == len(NumericalAttribute):
            plt.subplots_adjust(wspace=0.3, hspace=0.3)
            plt.show()
```

绘制出的部分盒图如下，可以直观的观察到离群值的分布

# 缺失数据的处理

## 1.将缺失部分剔除

可以进行填充的字段有如下：

```
NullValue = ['No_Score_Prob', 'Opp_Field_Goal_Prob', 'Opp_Safety_Prob', 'Opp_Touchdown_Prob', 'Field_Goal_Prob',
             'Safety_Prob', 'Touchdown_Prob', 'ExpPts', 'EPA', 'airEPA', 'yacEPA', 'Home_WP_pre', 'Away_WP_pre',
             'Home_WP_post', 'Away_WP_post', 'Win_Prob', 'WPA', 'airWPA', 'yacWPA']
```
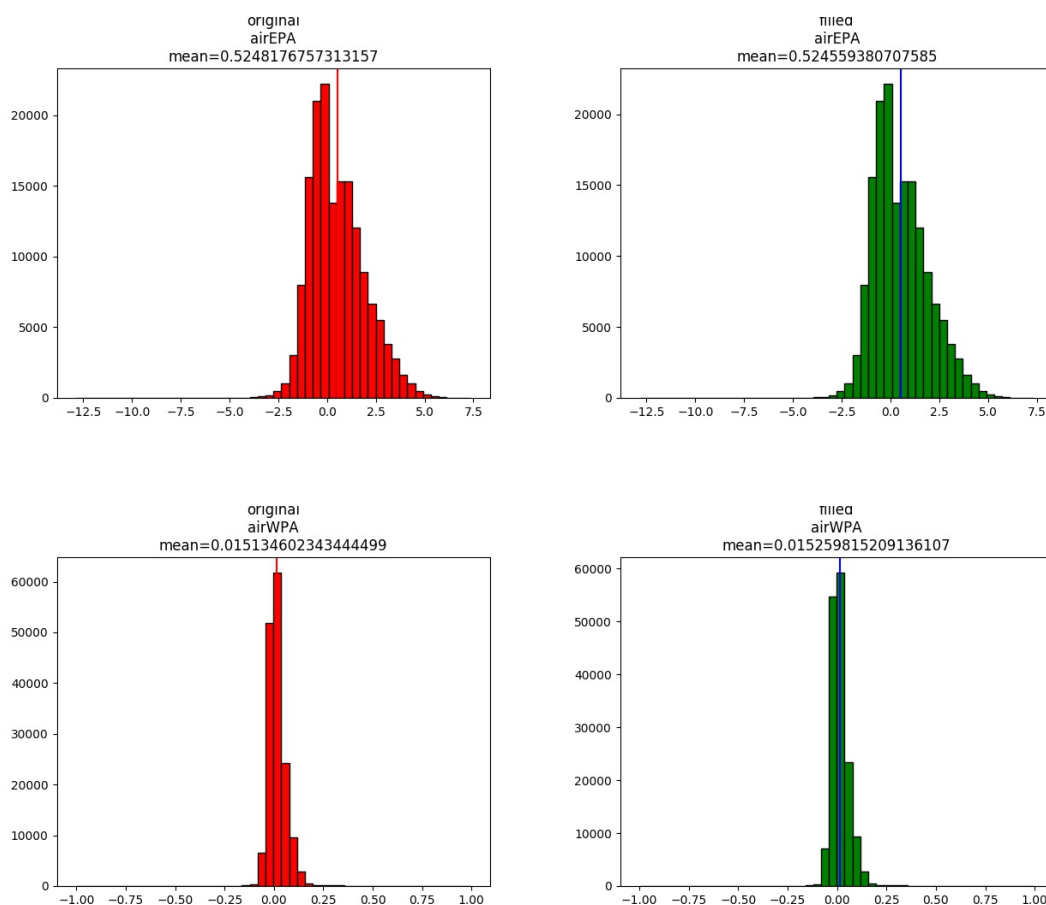
剔除缺失数据的代码如下：

```
index = data[NullValue].isnull().sum(axis=1) == 0
df_fillna = data[index]
compare(data, df_fillna, NullValue)
```

将剔除之后的数据与剔除前的数据进行直方图的比较，代码如下：

```python
def compare(df1, df2, columns, bins=50):
    for col in columns:
        mean1 = df1[col].mean()
        mean2 = df2[col].mean()

        fig = plt.figure()
        ax1 = fig.add_subplot(121)
        df1[col].hist(ax=ax1, grid=False, figsize=(15, 5),
                      bins=bins, edgecolor='black', facecolor='red')
        ax1.axvline(mean1, color='r')
        plt.title('original\n{}\nmean={}'.format(col, str(mean1)))
        ax2 = fig.add_subplot(122)
        df2[col].hist(ax=ax2, grid=False, figsize=(15, 5),
                      bins=bins, edgecolor='black', facecolor='green')
        ax2.axvline(mean2, color='b')
        plt.title('filled\n{}\nmean={}'.format(col, str(mean2)))
        plt.subplots_adjust(wspace=0.3, hspace=10)
        plt.savefig('/Users/margo/Desktop/ScreenShot/%s.jpg' %
                    col, format='jpg')
```

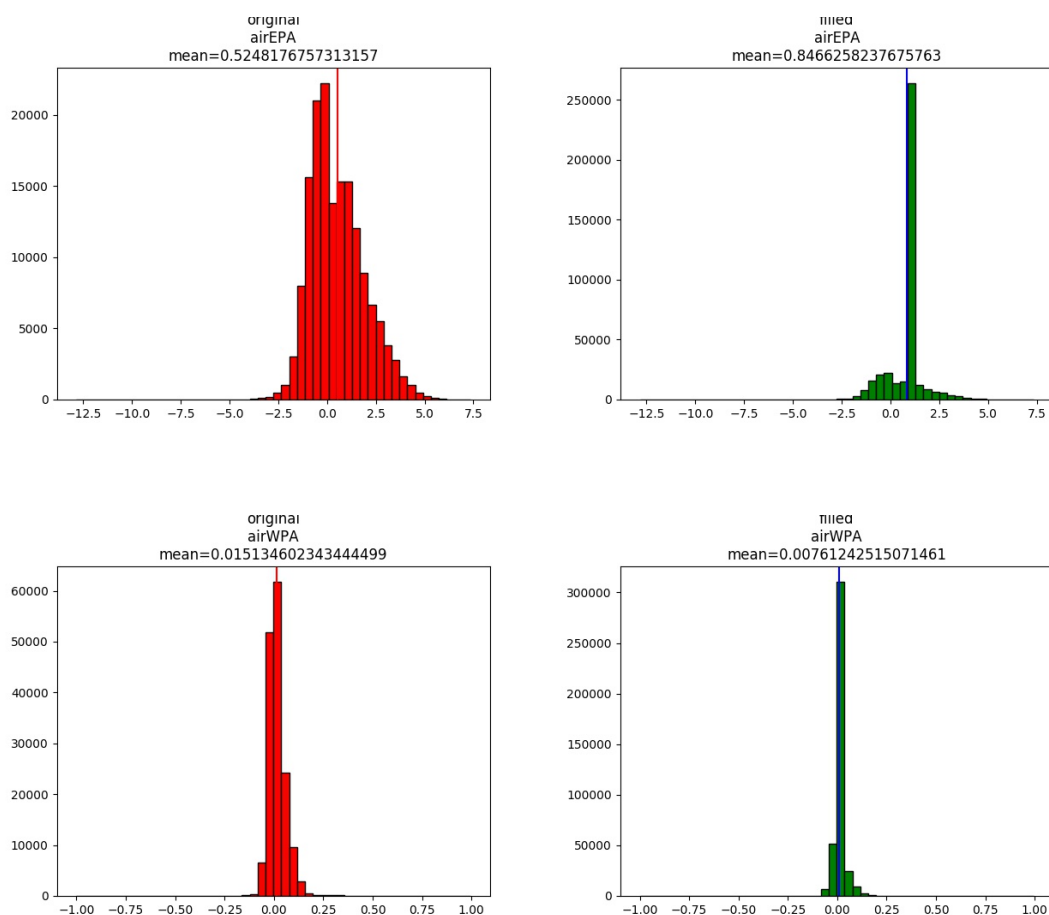比较的部分直方图如下，其中左边的红线为原有数据集的均值，右边的蓝线为剔除后的均值

## 2.用最高频率值填补缺失数据

用最高频率值填补缺失数据的代码如下：

```
# 最高频率值填补，比较
df_filled = data.copy()
for col in NullValue:
    most_frequent_value = df_filled[col].value_counts().idxmax()
    df_filled[col].fillna(value=most_frequent_value, inplace=True)
compare(data, df_filled, NullValue)
```

比较的部分直方图如下，其中左边的红线为原有数据集的均值，右边的蓝线为新数据的均值



## 3.通过属性的相关关系来填补缺失值

通过属性的相关关系来填补缺失值的代码如下：

```python
df_filled_inter = data.copy()
for col in NullValue:
    df_filled_inter[col].interpolate(inplace=True)

compare(data, df_filled_inter, NullValue)
```

得到的对比直方图如下所示，其中左边的红线为原有数据集的均值，右边的蓝线为新数据的均值