

## ▼ Predicting interest rates from Federal Reserve documents

### Exploratory Data Analysis/ Feature Engineering (Vol. 4)

FE 690: Machine Learning in Finance

Author: Theo Dimitrasopoulos

Advisor: Zachary Feinstein

## ▼ Environment

```
import sys
import os
IN_COLAB = 'google.colab' in sys.modules
IN_COLAB

    True

if IN_COLAB:
    from google.colab import drive
    drive.mount('/content/drive', force_remount=True)

    Mounted at /content/drive
```

```
#if IN_COLAB:
# # Uninstall existing versions:
# !pip uninstall bs4 -y
# !pip uninstall textract -y
# !pip uninstall numpy -y
# !pip uninstall pandas -y
# !pip uninstall requests -y
# !pip uninstall tqdm -y
# !pip uninstall nltk -y
# !pip uninstall quandl -y
# !pip uninstall scikit-plot -y
# !pip uninstall seaborn -y
# !pip uninstall sklearn -y
# !pip uninstall torch -y
# !pip uninstall transformers -y
# !pip uninstall wordcloud -y
# !pip uninstall xgboost -y
#
# # Install packages:
# !pip install bs4==0.0.1
```

```
# !pip install texttract==1.6.3
# !pip install numpy==1.19.4
# !pip install pandas==1.1.4
# !pip install requests==2.24.0
# !pip install tqdm==4.51.0
# !pip install nltk==3.5
# !pip install quandl==3.5.3
# !pip install scikit-plot==0.3.7
# !pip install seaborn==0.11.0
# !pip install sklearn==0.0
# !pip install torch==1.7.1+cu101 torchvision==0.8.2+cu101 -f https://download.pytorch.org/whl/torch\_stable.html
# !pip install transformers==3.5.0
# !pip install wordcloud==1.8.0
# !pip install xgboost==1.2.1
# os.kill(os.getpid(), 9)
#
```

```
if IN_COLAB:
```

```
    employment_data_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/Employment/'
    cpi_data_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/CPI/'
    fed_rates_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/FEDRates/'
    fx_rates_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/FXRates/'
    gdp_data_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/GDP/'
    ism_data_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/ISM/'
    sales_data_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/Sales/'
    treasury_data_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/MarketData/Treasury/'
    fomc_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/FOMC/'
    preprocessed_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/'
    train_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/train_data/'
    output_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/result/'
    keyword_lm_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/LoughranMcDonald/'
    glove_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/GloVe/'
    model_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/models/'
    graph_dir = '/content/drive/My Drive/Colab Notebooks/proj2/src/data/graphs/'
```

```
else:
```

```
    employment_data_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/Employment/'
    cpi_data_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/CPI/'
    fed_rates_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/FEDRates/'
    fx_rates_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/FXRates/'
    gdp_data_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/GDP/'
    ism_data_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/ISM/'
    sales_data_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/Sales/'
    treasury_data_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/MarketData/Treasury/'
    fomc_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/FOMC/'
    preprocessed_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/preprocessed/'
    train_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/train_data/'
    output_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/result/'
    keyword_lm_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/LoughranMcDonald/'
    glove_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/GloVe/'
    model_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/models/'
```

```
graph_dir = 'C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/graphs/'
```

```
# Python libraries
```

```
import pprint
import datetime as dt
import re
import pickle
from tqdm.notebook import tqdm
import time
import logging
import random
from collections import defaultdict, Counter
import xgboost as xgb
```

```
# Data Science modules
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
plt.style.use('ggplot')
import seaborn as sns; sns.set(style='white', context='notebook', palette='deep')
```

```
# Import Scikit-learn models
```

```
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.metrics import accuracy_score, f1_score, plot_confusion_matrix
from sklearn.pipeline import Pipeline, FeatureUnion
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier, ExtraTreesClassifier, VotingClassifier
from sklearn.linear_model import LogisticRegression, Perceptron, SGDClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.neural_network import MLPClassifier
from sklearn.svm import SVC, LinearSVC
from sklearn import model_selection
from sklearn.model_selection import GridSearchCV, cross_val_score, cross_validate, StratifiedKFold, learning_curve, RandomizedSearchCV
import scikitplot as skplt
```

```
# Import nltk modules and download dataset
```

```
import nltk
from nltk.corpus import stopwords
from nltk.util import ngrams
from nltk.tokenize import word_tokenize, sent_tokenize
```

```
# Import Pytorch modules
```

```
import torch
from torch import nn, optim
import torch.nn.functional as F
from torch.utils.data import (DataLoader, RandomSampler, SequentialSampler, TensorDataset)
from torch.autograd import Variable
from torch.optim import Adam, AdamW
```

from torch.optim import Adam, AdamW

plt.style.use('ggplot')

sns.set()

# Fiinalize nltk setup:

nltk.download('stopwords')

nltk.download('punkt')

nltk.download('wordnet')

stop = set(stopwords.words('english'))

# Test pprint

pprint.pprint(sys.path)

```
['',
 '/env/python',
 '/usr/lib/python36.zip',
 '/usr/lib/python3.6',
 '/usr/lib/python3.6/lib-dynload',
 '/usr/local/lib/python3.6/dist-packages',
 '/usr/lib/python3/dist-packages',
 '/usr/local/lib/python3.6/dist-packages/IPython/extensions',
 '/root/.ipython']
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

## Use TPU Runtime:

#if IN\_COLAB:

# assert os.environ['COLAB\_TPU\_ADDR'], 'Make sure to select TPU from Edit > Notebook setting > Hardware accelerator'

# VERSION = "20200220"

# !curl https://raw.githubusercontent.com/pytorch/xla/master/contrib/scripts/env-setup.py -o pytorch-xla-env-setup.py

# !python pytorch-xla-env-setup.py --version \$VERSION

## Use GPU Runtime:

if IN\_COLAB:

if torch.cuda.is\_available():

torch.cuda.get\_device\_name(0)

gpu\_info = !nvidia-smi

gpu\_info = '\n'.join(gpu\_info)

print(gpu\_info)

else:

print('Select the Runtime > "Change runtime type" menu to enable a GPU accelerator, and then re-execute this cell.')

os.kill(os.getpid(), 9)

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											
NVIDIA-SMI		460.32.03		Driver Version: 418.67				CUDA Version: 10.1																			
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											
GPU	Name		Persistence-M		Bus-Id		Disp.A		Volatile Uncorr. ECC																		
Fan	Temp	Perf	Pwr:Usage/Cap				Memory-Usage		GPU-Util		Compute M.																
												MIG M.															
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											
0	Tesla	V100-SXM2...	Off	00000000:00:04.0		Off				0																	
N/A	35C	P0	25W / 300W		10MiB / 16130MiB				0%		Default																
												ERR!															
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											
Processes:																											
GPU	GI	CI	PID		Type	Process name				GPU Memory																	
		ID	ID							Usage																	
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											
No running processes found																											
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+																											

```
# Load nontext data
if IN_COLAB:
    file = open('/content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_data.pickle', 'rb')
    nontext_data = pickle.load(file)
    file.close()
else:
    file = open('C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/preprocessed/nontext_data.pickle', 'rb')

print(nontext_data.shape)
nontext_data.head()
```

(390, 57)

	unscheduled	forecast	confcall	ChairPerson	Rate	RateDiff	RateDecision	RateChanged	GDP_date	GDP_value	GDP_diff_prev	GDP_diff_year	GDPPOT_
date													
1982-10-05	False	False	False	Paul Volcker	9.5	-0.5	-1	1	1982-04-01	6825.876	0.456197	-1.010549	1982-10-05
1982-11-16	False	False	False	Paul Volcker	9.0	-0.5	-1	1	1982-07-01	6799.781	-0.382295	-2.555898	1982-11-16
1982-12-21	False	False	False	Paul Volcker	8.5	0.0	0	0	1982-07-01	6799.781	-0.382295	-2.555898	1982-12-21
1983-01-14	False	False	True	Paul Volcker	8.5	0.0	0	0	1982-07-01	6799.781	-0.382295	-2.555898	1983-01-14
1983-01-21	False	False	True	Paul Volcker	8.5	0.0	0	0	1982-07-01	6799.781	-0.382295	-2.555898	1983-01-21

```
if IN_COLAB:
    file = open('/content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma2.pickle', 'rb')
    nontext_ma2 = pickle.load(file)
    file.close()
    file = open('/content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma3.pickle', 'rb')
    nontext_ma3 = pickle.load(file)
    file.close()
    file = open('/content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma6.pickle', 'rb')
    nontext_ma6 = pickle.load(file)
    file.close()
    file = open('/content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma12.pickle', 'rb')
    nontext_ma12 = pickle.load(file)
    file.close()
else:
    file = open('C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma2.pickle', 'rb')
    nontext_ma2 = pickle.load(file)
    file.close()
    file = open('C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma3.pickle', 'rb')
    nontext_ma3 = pickle.load(file)
    file.close()
    file = open('C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma6.pickle', 'rb')
    nontext_ma6 = pickle.load(file)
    file.close()
    file = open('C:/Users/theon/GDrive/Colab Notebooks/proj2/src/data/preprocessed/nontext_ma12.pickle', 'rb')
    nontext_ma12 = pickle.load(file)
    file.close()
```

## Non-text dataset analysis

```
nontext_data['prev_decision'] = nontext_data['RateDecision'].shift(1)
nontext_data['next_decision'] = nontext_data['RateDecision'].shift(-1)
nontext_data[['RateDecision', 'prev_decision', 'next_decision']].head()
```

	RateDecision	prev_decision	next_decision
date			
1982-10-05	-1	<NA>	-1
1982-11-16	-1	-1	0
1982-12-21	0	-1	0
1983-01-14	0	0	0
1983-01-21	0	0	0

```
nontext_data.describe()
```

	Rate	RateDiff	RateDecision	RateChanged	GDP_value	GDP_diff_prev	GDP_diff_year	GDPPOT_value	GDPPOT_diff_prev	GDPPOT_diff_year	PCI
count	390.000000	390.000000	390.000000	390.000000	390.00000	390.000000	390.000000	390.000000	390.000000	390.000000	390
mean	4.141346	-0.025641	-0.020513	0.348718	12450.45391	0.647576	2.594788	12647.034520	0.657539	2.675490	79
std	3.056928	0.239839	0.590925	0.477177	3542.17116	0.623227	1.954636	3563.685301	0.204059	0.829742	17
min	0.000000	-1.000000	-1.000000	0.000000	6799.78100	-2.163811	-3.924447	7224.140335	0.263026	1.080299	47
25%	1.250000	0.000000	0.000000	0.000000	9341.64200	0.396135	1.722931	9578.876710	0.484126	1.970276	66
50%	4.375000	0.000000	0.000000	0.000000	12345.14650	0.631956	2.694679	12142.911975	0.648491	2.672563	79
75%	6.187500	0.000000	0.000000	1.000000	15645.57475	0.970216	3.908121	15938.819440	0.822396	3.316956	94
max	11.500000	1.125000	1.000000	1.000000	19221.97000	2.275605	8.578274	19099.880000	1.064642	4.300945	112

nontext\_data.isnull().sum()

```

forecast          0
confcalls         0
ChairPerson       0
Rate              0
RateDiff          0
RateDecision      0
RateChanged       0
GDP_date          0
GDP_value         0
GDP_diff_prev     0
GDP_diff_year     0
GDPPOT_date       0
GDPPOT_value      0
GDPPOT_diff_prev  0
GDPPOT_diff_year  0
PCE_date          0
PCE_value         0
PCE_diff_prev     0
PCE_diff_year     0
CPI_date          0
CPI_value         0
CPI_diff_prev     0
CPI_diff_year     0
Unemp_date        0
Unemp_value       0
Unemp_diff_prev   0
Unemp_diff_year   0
Employ_date       0
Employ_value      0
Employ_diff_prev  0
Employ_diff_year  0

PMI_date          0

```

PMI_value	0
PMI_diff_prev	0
PMI_diff_year	0
NMI_date	281
NMI_value	281
NMI_diff_prev	281
NMI_diff_year	295
Rsales_date	116
Rsales_value	116
Rsales_diff_prev	117
Rsales_diff_year	129
Hsales_date	0
Hsales_value	0
Hsales_diff_prev	0
Hsales_diff_year	0
Taylor	0
Balanced	0
Inertia	0
Taylor-Rate	0
Balanced-Rate	0
Inertia-Rate	0
Taylor_diff	1
Balanced_diff	1
Inertia_diff	1
prev_decision	1
next_decision	1
dtvne: int64	

```
x = nontext_data['RateDecision'].value_counts()
print("Count: ")
print(x)
print("Percent: ")
print(round(x/sum(x) * 100))
plt.figure(figsize=(8,5))
ax = sns.countplot(x='RateDecision', data=nontext_data)
ax.set_title('nontext_data')
```



```

Count:
0      254
-1      72
1       64
Name: RateDecision, dtype: Int64
Percent:
0      65.0
-1     18.0
1     16.0
Name: RateDecision, dtype: float64
Text(0.5, 1.0, 'nontext_data')

```



## Correlation



```
nontext_data.columns.values
```

```

array(['unscheduled', 'forecast', 'confcall', 'ChairPerson', 'Rate',
      'RateDiff', 'RateDecision', 'RateChanged', 'GDP_date', 'GDP_value',
      'GDP_diff_prev', 'GDP_diff_year', 'GDPPOT_date', 'GDPPOT_value',
      'GDPPOT_diff_prev', 'GDPPOT_diff_year', 'PCE_date', 'PCE_value',
      'PCE_diff_prev', 'PCE_diff_year', 'CPI_date', 'CPI_value',
      'CPI_diff_prev', 'CPI_diff_year', 'Unemp_date', 'Unemp_value',
      'Unemp_diff_prev', 'Unemp_diff_year', 'Employ_date',
      'Employ_value', 'Employ_diff_prev', 'Employ_diff_year', 'PMI_date',
      'PMI_value', 'PMI_diff_prev', 'PMI_diff_year', 'NMI_date',
      'NMI_value', 'NMI_diff_prev', 'NMI_diff_year', 'Rsales_date',
      'Rsales_value', 'Rsales_diff_prev', 'Rsales_diff_year',
      'Hsales_date', 'Hsales_value', 'Hsales_diff_prev',
      'Hsales_diff_year', 'Taylor', 'Balanced', 'Inertia', 'Taylor-Rate',
      'Balanced-Rate', 'Inertia-Rate', 'Taylor_diff', 'Balanced_diff',
      'Inertia_diff', 'prev_decision', 'next_decision'], dtype=object)

```

```

corr_columns = ['RateDecision', 'next_decision', 'prev_decision', 'unscheduled', 'forecast', 'confcall',
               'GDP_diff_prev', 'GDP_diff_year', 'GDPPOT_diff_prev', 'GDPPOT_diff_year',
               'PCE_diff_prev', 'PCE_diff_year', 'CPI_diff_prev', 'CPI_diff_year',
               'Unemp_value', 'Unemp_diff_prev', 'Unemp_diff_year',
               'Employ_value', 'Employ_diff_prev', 'Employ_diff_year',
               'PMI_value', 'PMI_diff_prev', 'PMI_diff_year',
               'Rsales_diff_prev', 'Rsales_diff_year', 'Hsales_diff_prev', 'Hsales_diff_year',
               'Taylor_diff', 'Balanced_diff', 'Inertia_diff', 'Rate', 'RateDiff', 'RateChanged']

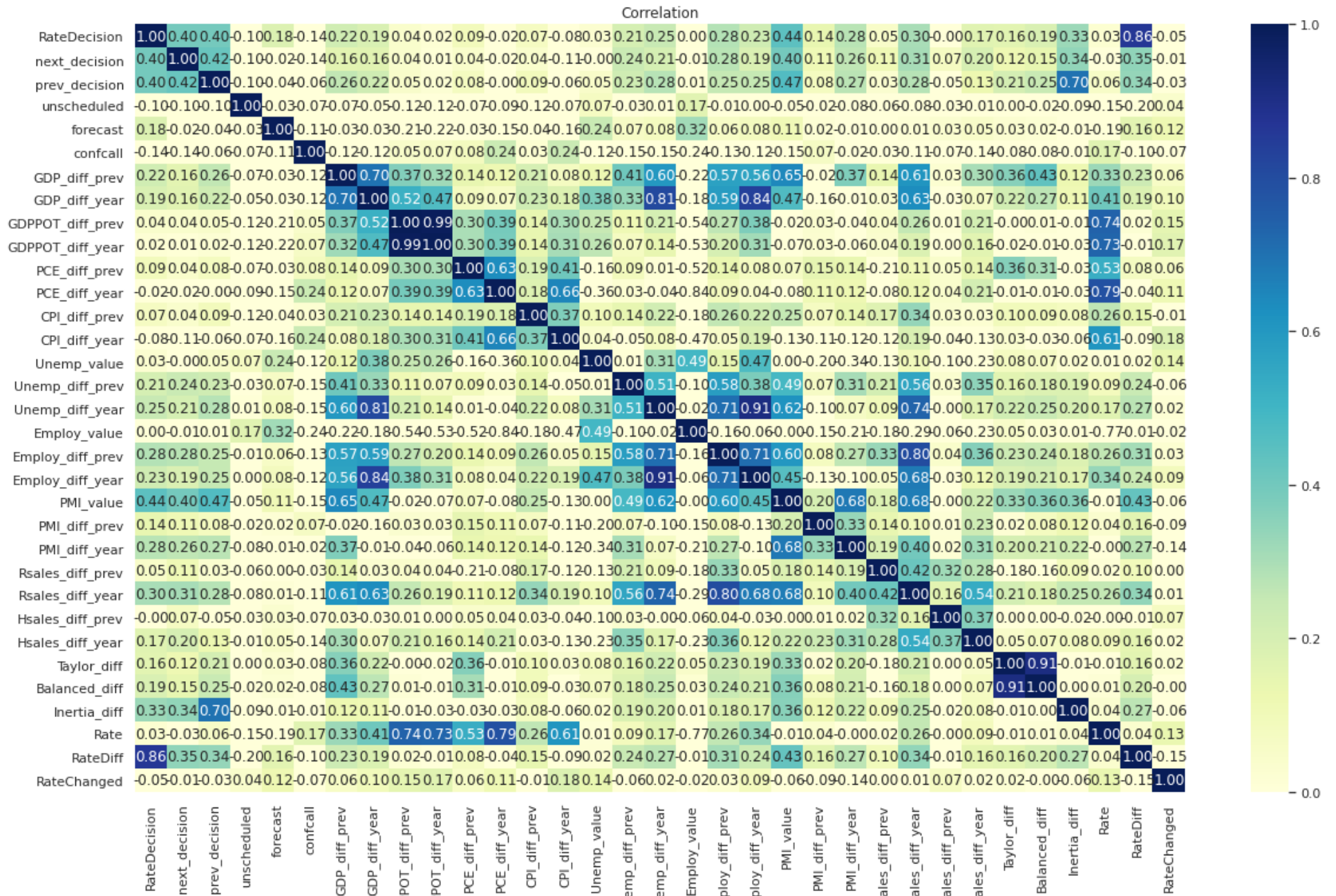
```

```
fig, ax = plt.subplots(1, 1, figsize=(20, 12))
```

```

sns.heatmap(nontext_data[corr_columns].astype(float).corr(), cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax)
ax.set_title("Correlation")
plt.show()

```



## Moving average

```
corr_columns = ['RateDecision',
                'GDP_diff_prev', 'GDP_diff_year', 'GDPPOT_diff_prev', 'GDPPOT_diff_year',
                'PCE_diff_prev', 'PCE_diff_year', 'CPI_diff_prev', 'CPI_diff_year',
                'Unemp_value', 'Unemp_diff_prev', 'Unemp_diff_year',
                'Employ_value', 'Employ_diff_prev', 'Employ_diff_year',
                'PMI_value', 'PMI_diff_prev', 'PMI_diff_year',
                'Rsales_diff_prev', 'Rsales_diff_year', 'Hsales_diff_prev', 'Hsales_diff_year',
                'Taylor_diff', 'Balanced_diff', 'Inertia_diff', 'Rate', 'RateDiff', 'RateChanged']
```

```
PMI_value , PMI_diff_prev , PMI_diff_year ,  
'Rsales_diff_prev', 'Rsales_diff_year', 'Hsales_diff_prev', 'Hsales_diff_year',  
'Taylor_diff', 'Balanced_diff', 'Inertia_diff', 'Taylor-Rate', 'Balanced-Rate', 'Inertia-Rate']
```

```
fig, (ax1, ax2, ax3, ax4, ax5) = plt.subplots(5, 1, figsize=(17,8))
```

```
sns.heatmap(nontext_data[corr_columns].astype(float).corr().iloc[:1], cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax1)  
ax1.set_title("Correlation: Original")  
ax1.set_xticks([])  
ax1.set_yticks([])  
sns.heatmap(nontext_ma2[corr_columns].astype(float).corr().iloc[:1], cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax2)  
ax2.set_title("Correlation: Moving average of 2 periods")  
ax2.set_xticks([])  
ax2.set_yticks([])  
sns.heatmap(nontext_ma3[corr_columns].astype(float).corr().iloc[:1], cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax3)  
ax3.set_title("Correlation: Moving average of 3 periods")  
ax3.set_xticks([])  
ax3.set_yticks([])  
sns.heatmap(nontext_ma6[corr_columns].astype(float).corr().iloc[:1], cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax4)  
ax4.set_title("Correlation: Moving average of 6 periods")  
ax4.set_xticks([])  
ax4.set_yticks([])  
sns.heatmap(nontext_ma12[corr_columns].astype(float).corr().iloc[:1], cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax5)  
ax5.set_title("Correlation: Moving average of 12 periods")  
ax5.set_yticks([])  
  
fig.tight_layout(pad=1.0)  
plt.show()
```

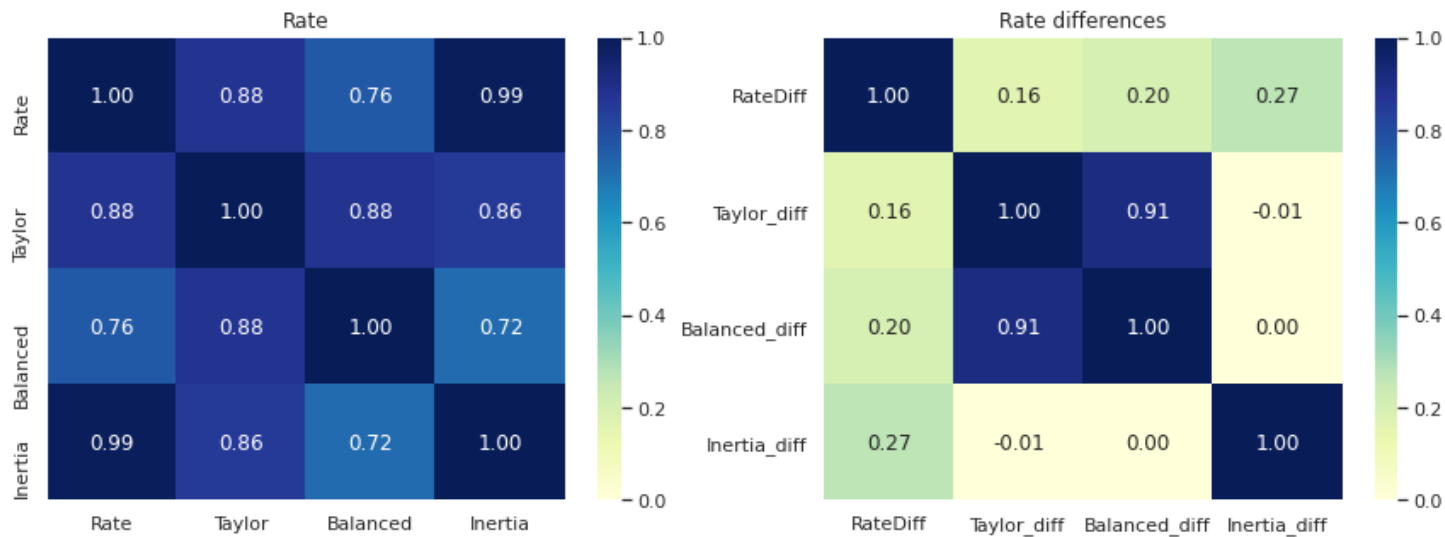


```
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 5))
```

```
corr_columns = ['Rate', 'Taylor', 'Balanced', 'Inertia']
sns.heatmap(nontext_data[corr_columns].astype(float).corr(), cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax1)
ax1.set_title("Rate")
```

```
corr_columns = ['RateDiff', 'Taylor_diff', 'Balanced_diff', 'Inertia_diff']
sns.heatmap(nontext_data[corr_columns].astype(float).corr(), cmap="YlGnBu", annot=True, fmt=".2f", vmin=0, vmax=1, ax=ax2)
ax2.set_title("Rate differences")
```

```
plt.show()
```



```
fig, ax = plt.subplots(figsize=(18, 10))
sns.lineplot(data=nontext_data[corr_columns], ax=ax)
```

```
decision_raise = nontext_data.loc[nontext_data['RateDecision'] == 1]
decision_hold = nontext_data.loc[nontext_data['RateDecision'] == 0]
decision_lower = nontext_data.loc[nontext_data['RateDecision'] == -1]
```

```
ax.plot(decision_raise.index.values, decision_raise['Rate'], 'o', color="g", label="Raise")
ax.plot(decision_hold.index.values, decision_hold['Rate'], 'o', color="grey", label="Hold")
ax.plot(decision_lower.index.values, decision_lower['Rate'], 'o', color="r", label="Lower")
```

```
[<matplotlib.lines.Line2D at 0x7fe77167a9b0>]
```



## ▼ Check individual columns

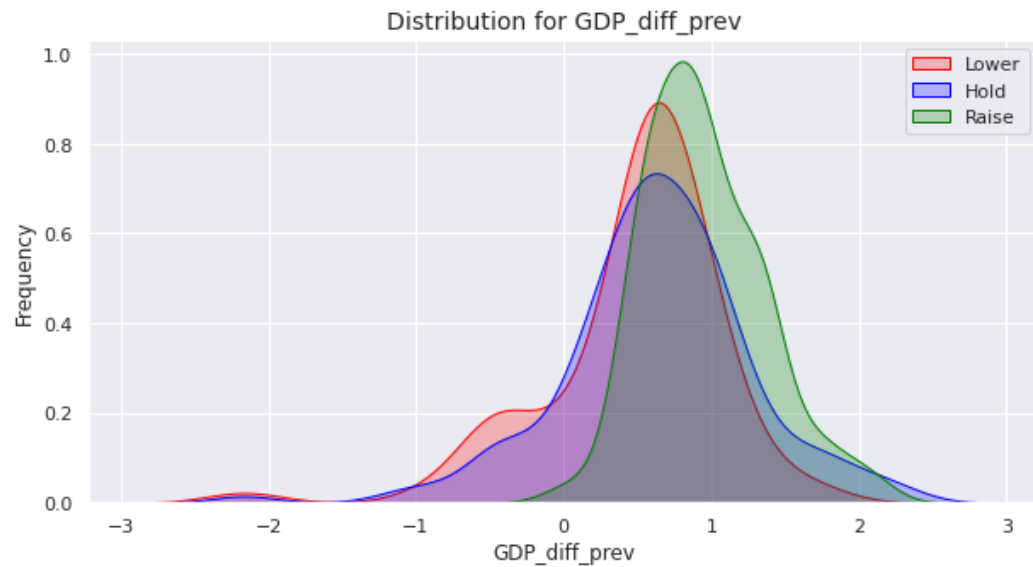
```
def plot_distribution(df, columns):
    for col in columns:
        fig, ax = plt.subplots(figsize=(10, 5))
        g = sns.kdeplot(df[col][(df["RateDecision"] == -1) & (df[col].notnull())], color="Red", shade=True)
        g = sns.kdeplot(df[col][(df["RateDecision"] == 0) & (df[col].notnull())], ax=g, color="Blue", shade=True)
        g = sns.kdeplot(df[col][(df["RateDecision"] == 1) & (df[col].notnull())], ax=g, color="Green", shade=True)
        g.set_xlabel(col)
        g.set_ylabel("Frequency")
        g.set_title("Distribution for " + col, fontsize=14)
        g = g.legend(["Lower", "Hold", "Raise"])

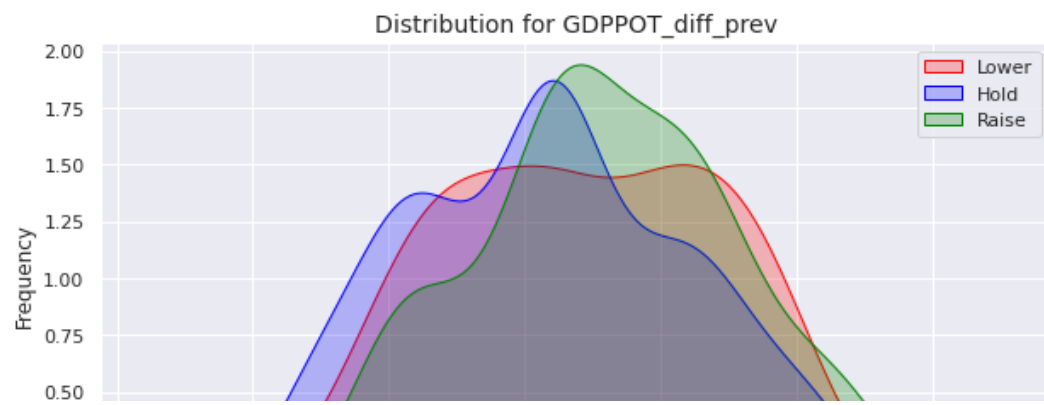
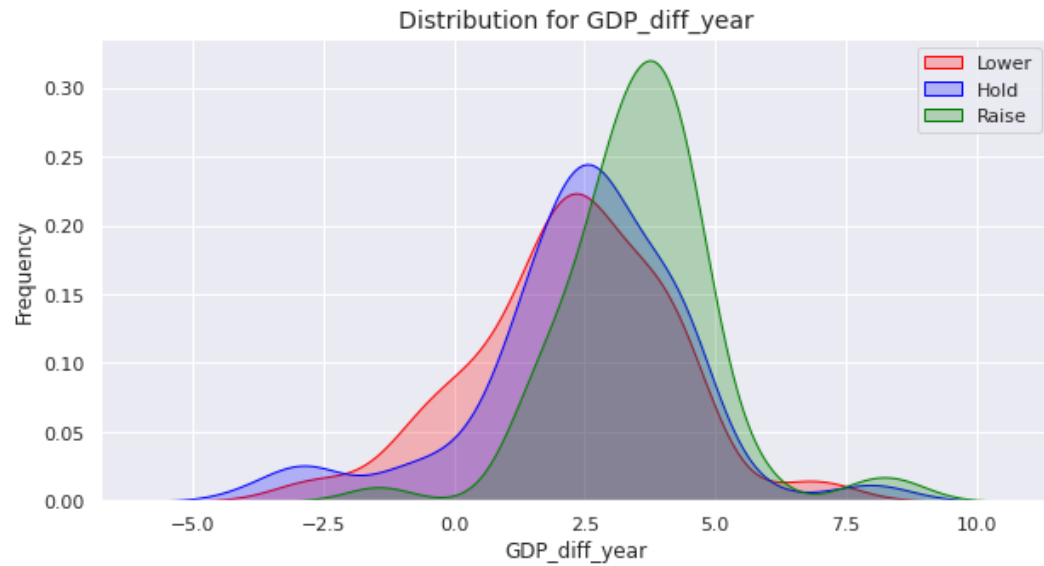
    g = sns.FacetGrid(df, col='RateDecision', height=3, aspect=1)
```

```
g.map(sns.distplot, col)
```

```
plot_distribution(nontext_data, ["GDP_diff_prev", "GDP_diff_year", "GDPPOT_diff_prev", "GDPPOT_diff_year"])
```

```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
```

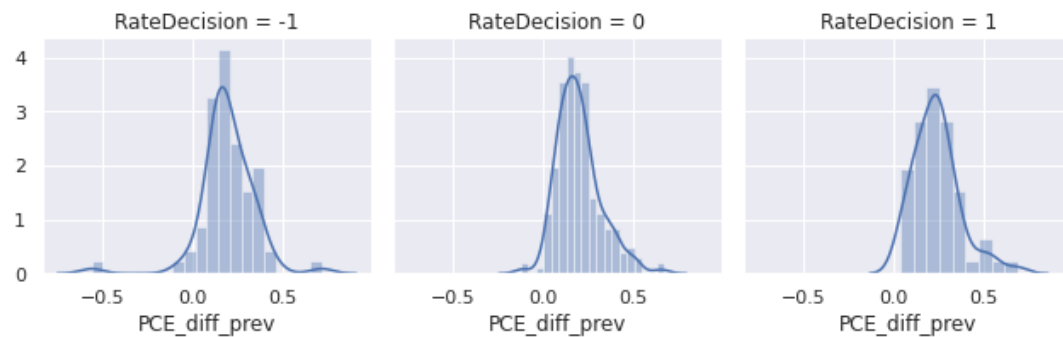
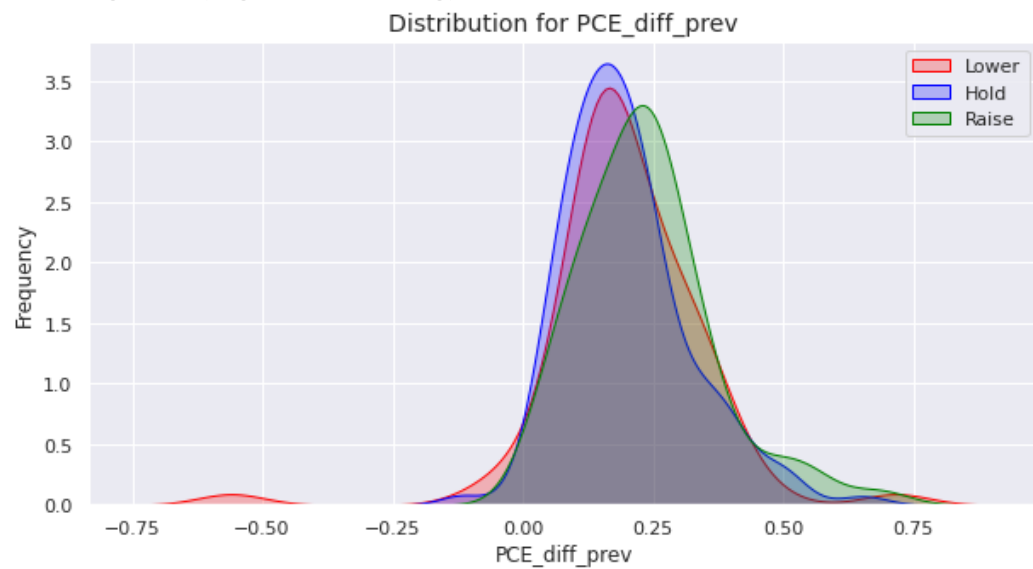


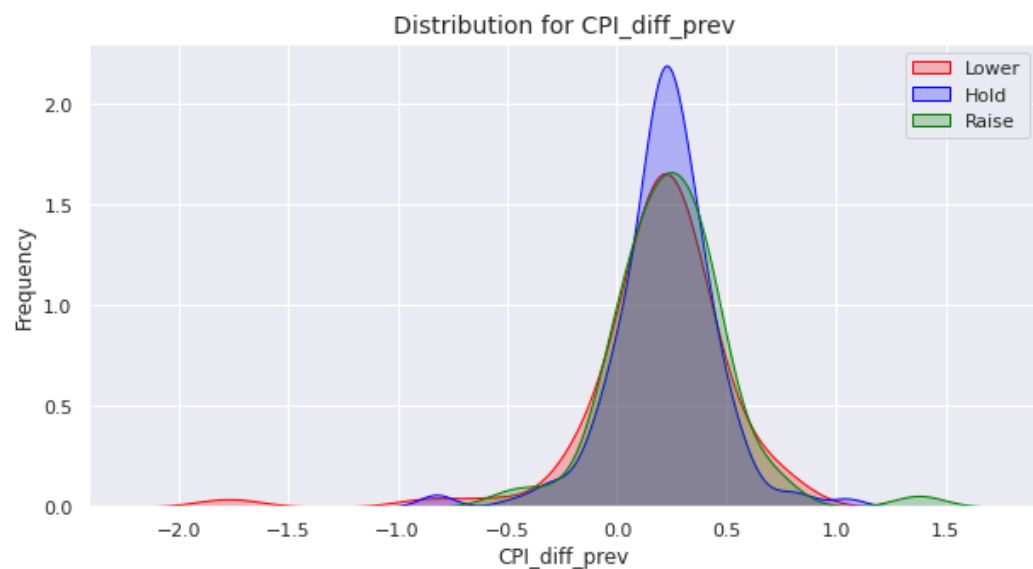
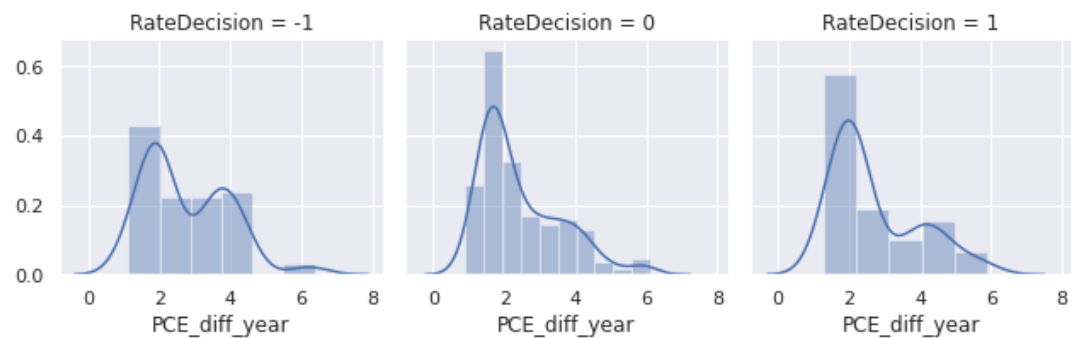
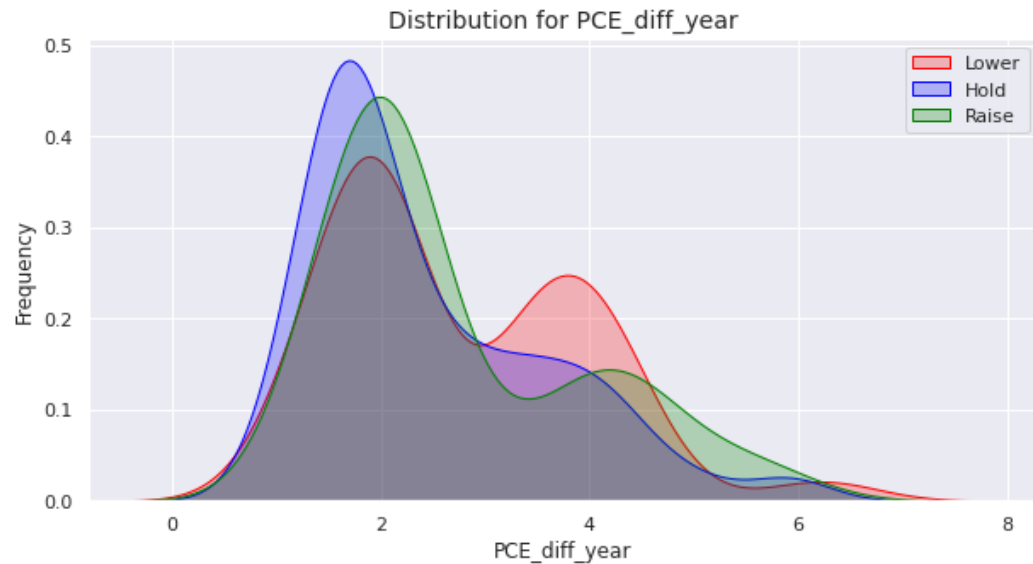


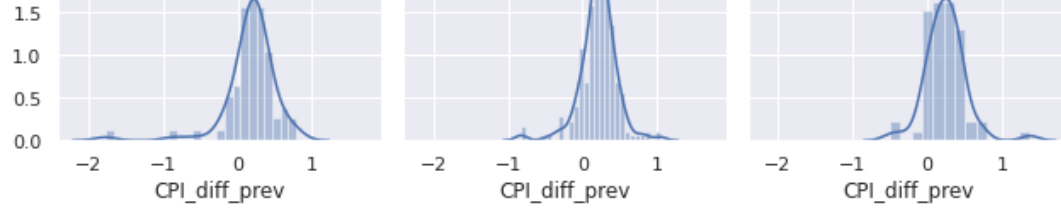
```
plot_distribution(nontext_data, ["PCE_diff_prev", "PCE_diff_year", "CPI_diff_prev", "CPI_diff_year"])
```



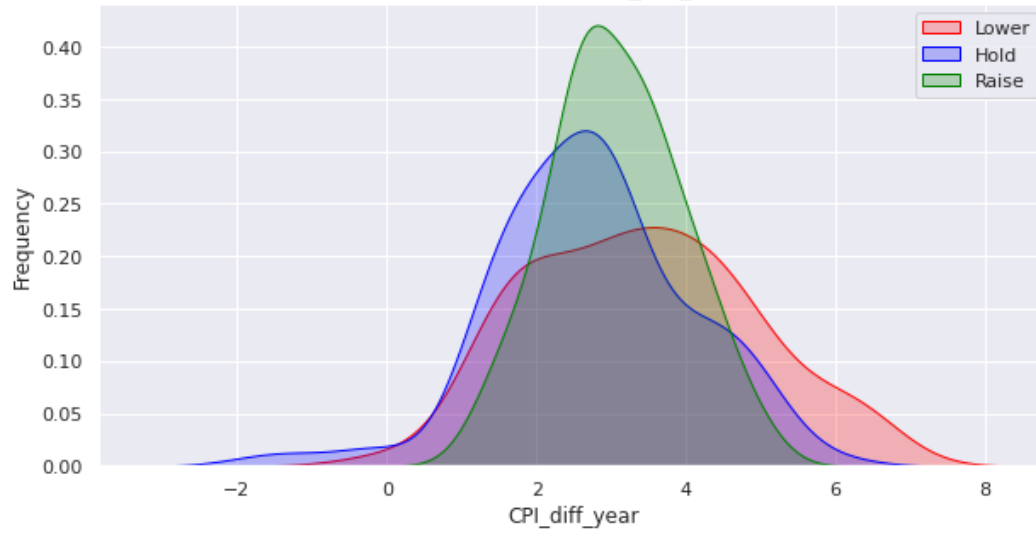
```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
```





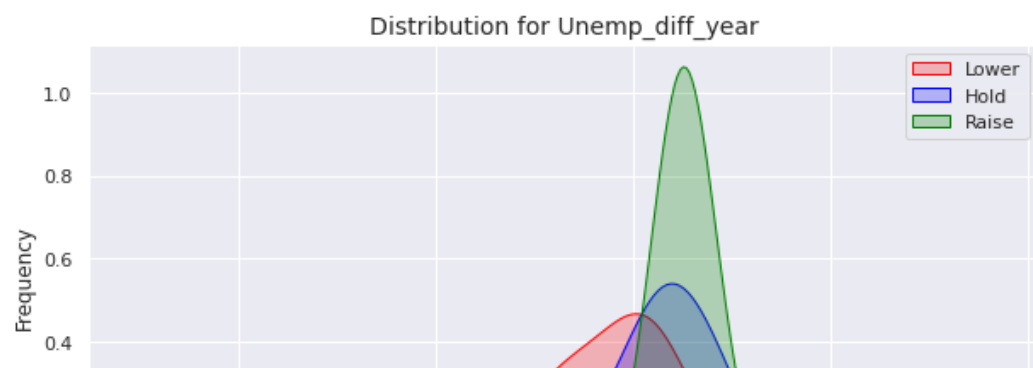
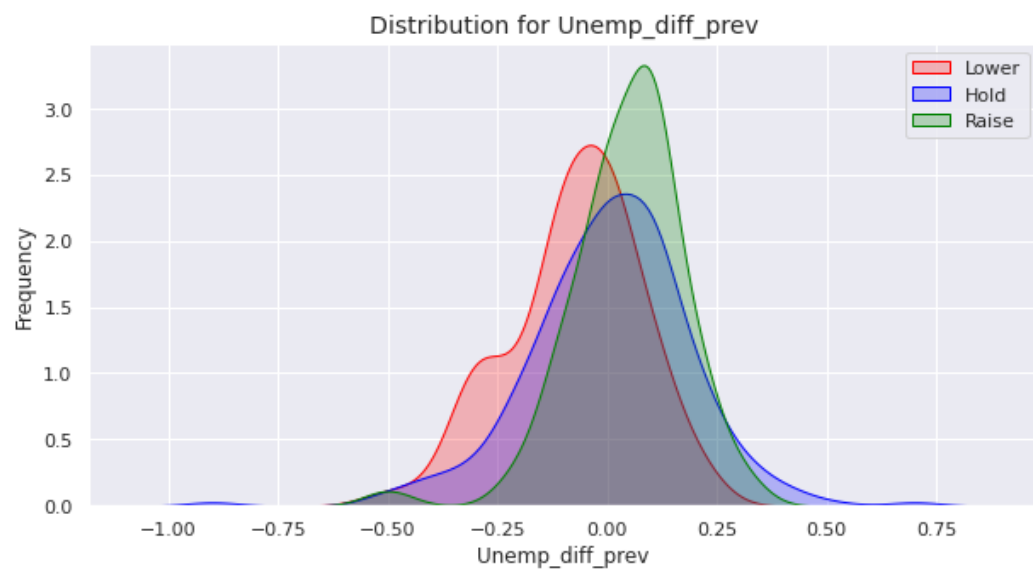
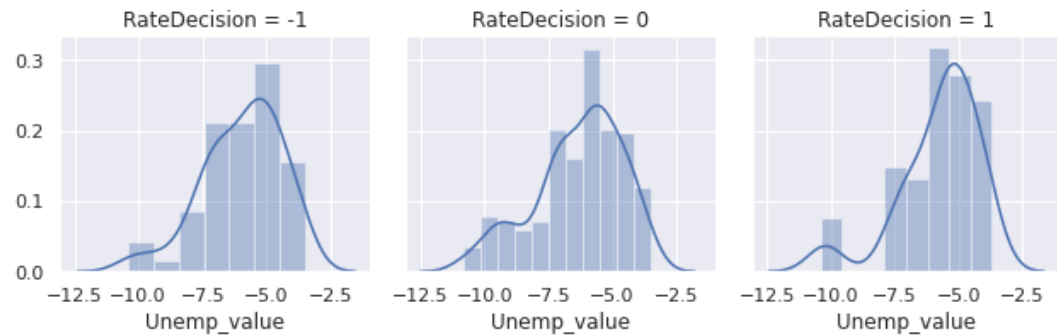


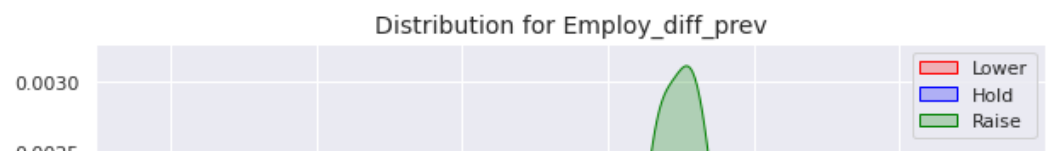
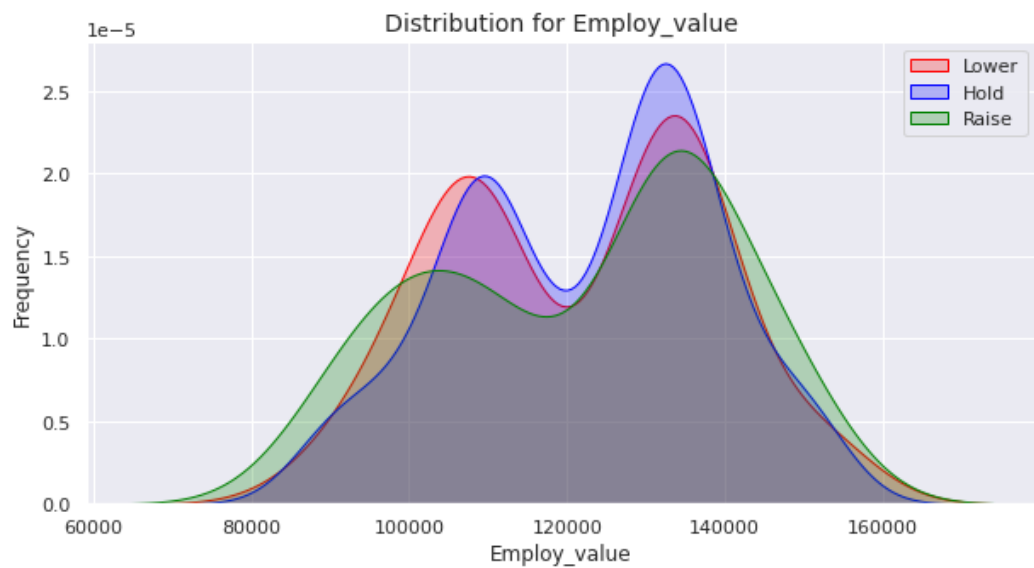
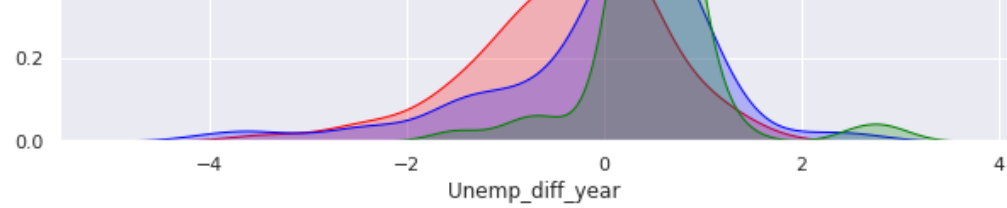
Distribution for CPI\_diff\_year



```
plot_distribution(nontext_data, ["Unemp_value", "Unemp_diff_prev", "Unemp_diff_year", "Employ_value", "Employ_diff_prev", "Employ_diff_year"])
```



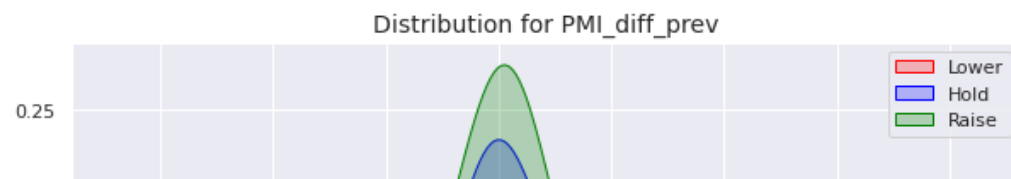
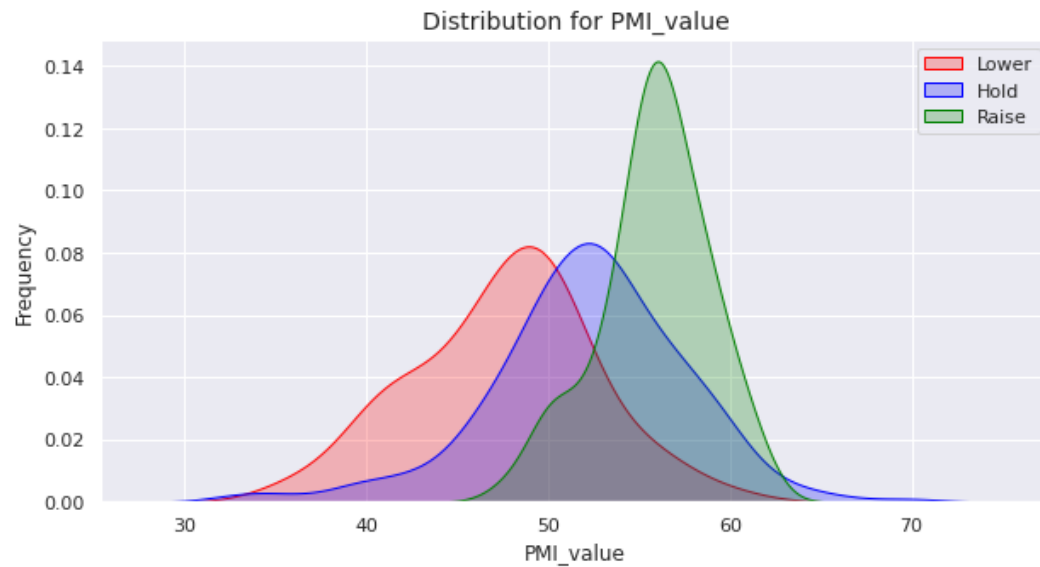




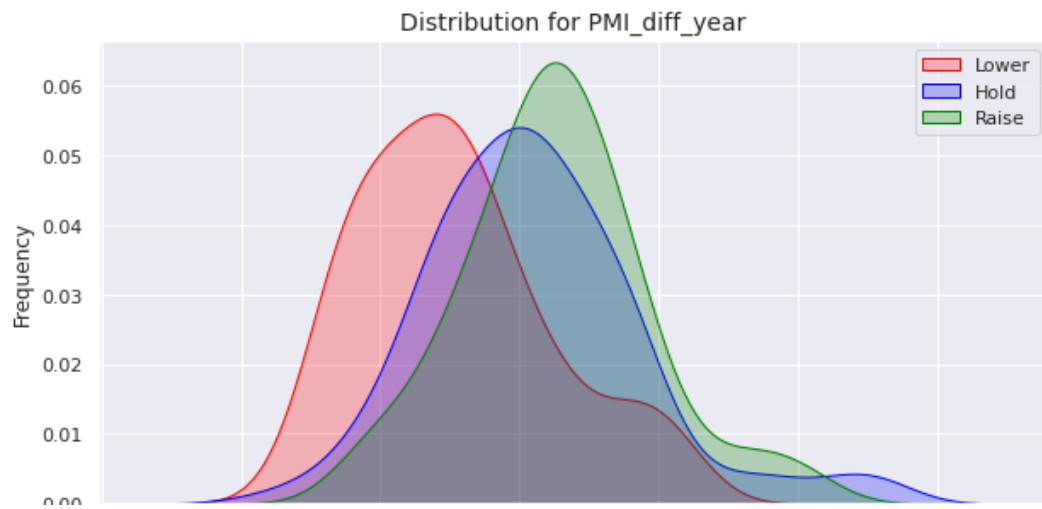
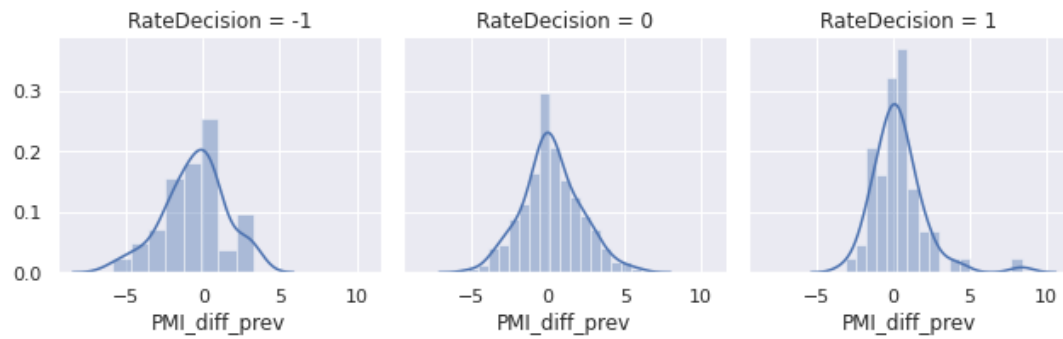
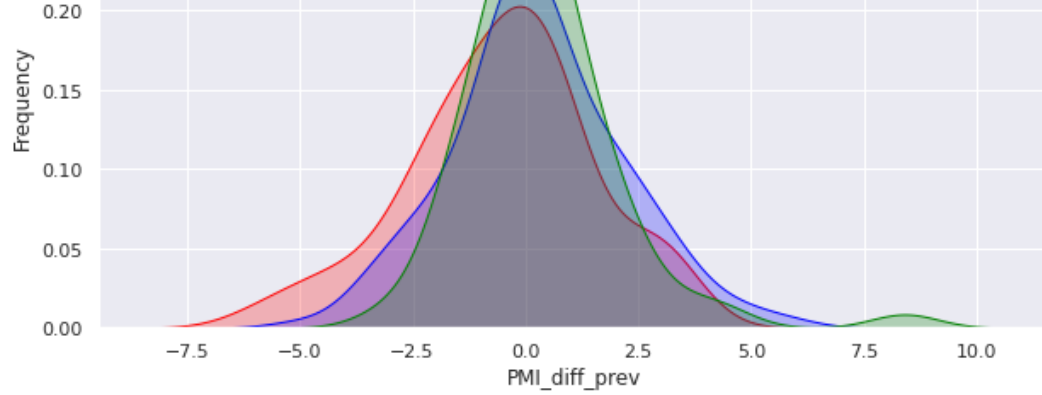


```
plot_distribution(nontext_data, ["PMI_value", "PMI_diff_prev", "PMI_diff_year"])
```

```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
```

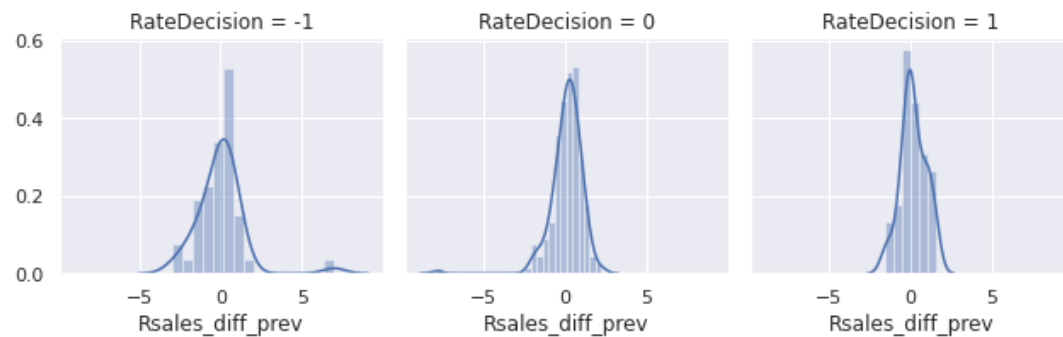
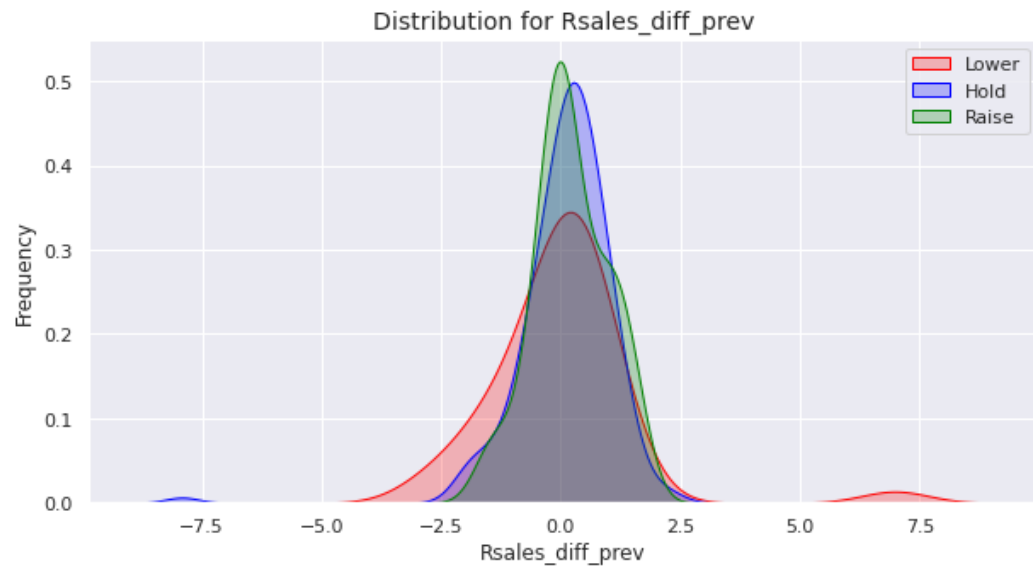




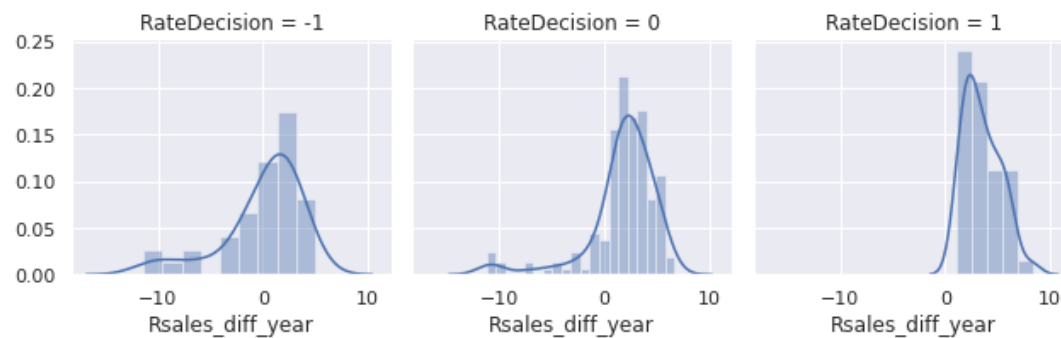
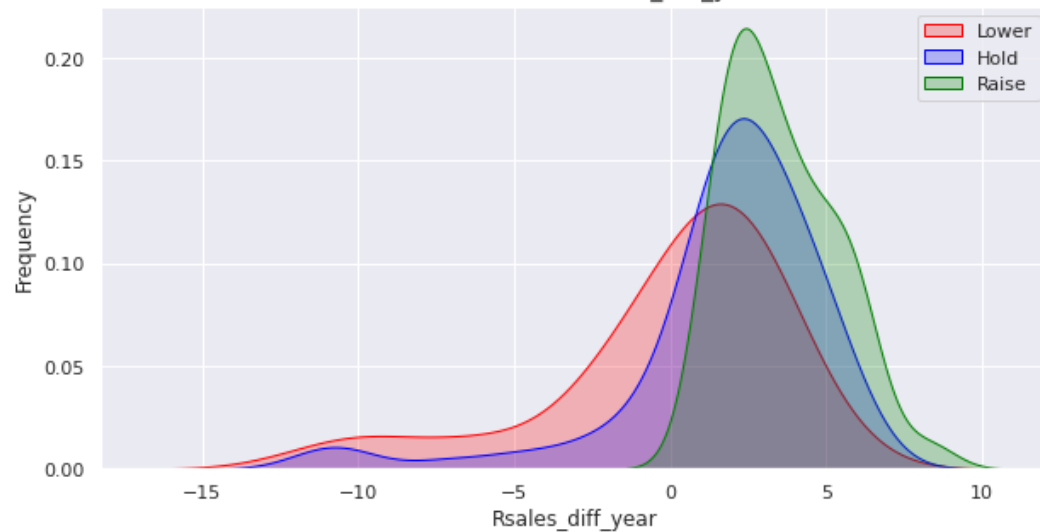


```
plot_distribution(nontext_data, ["Rsales_diff_prev", "Rsales_diff_year", "Hsales_diff_prev", "Hsales_diff_year"])
```

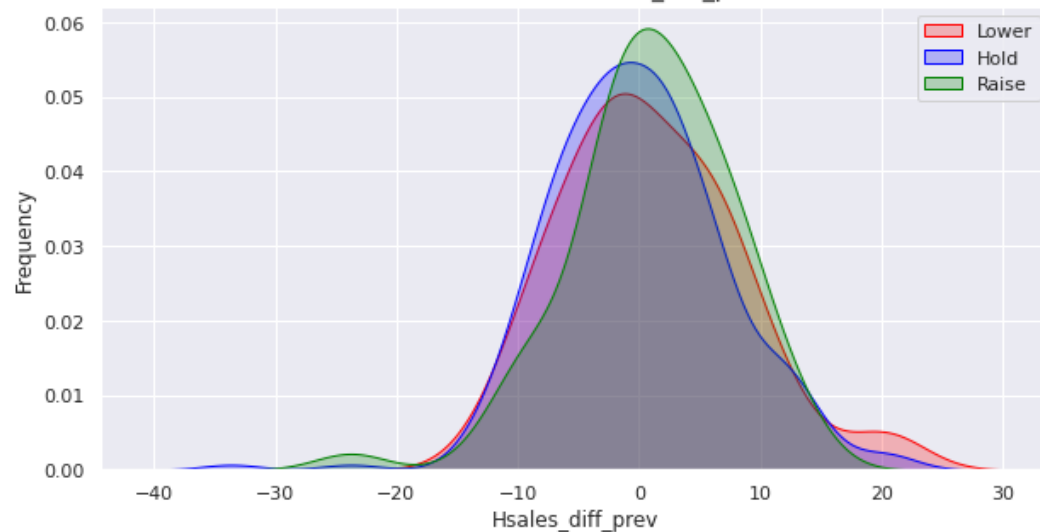
```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
```

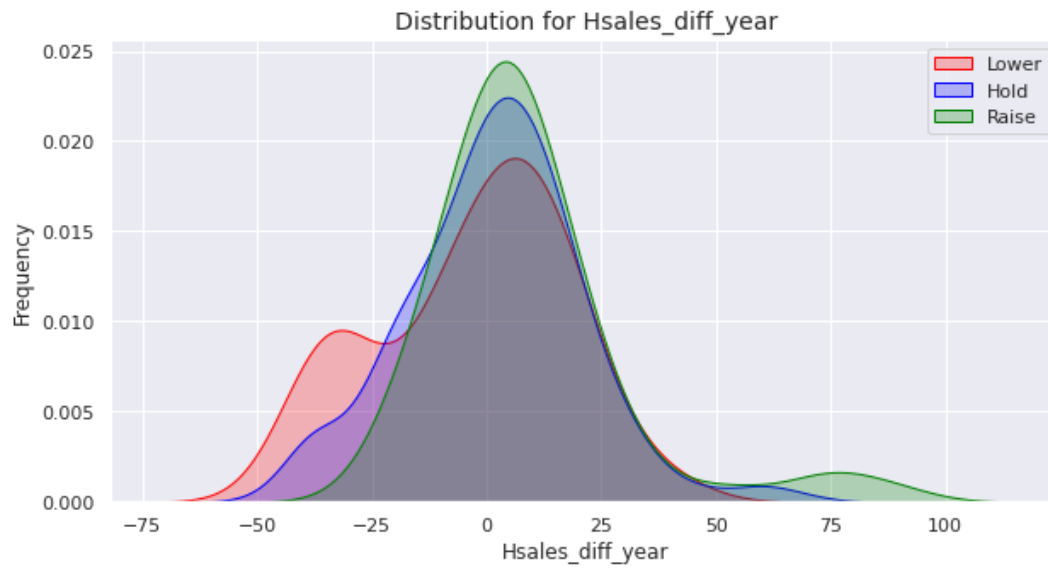
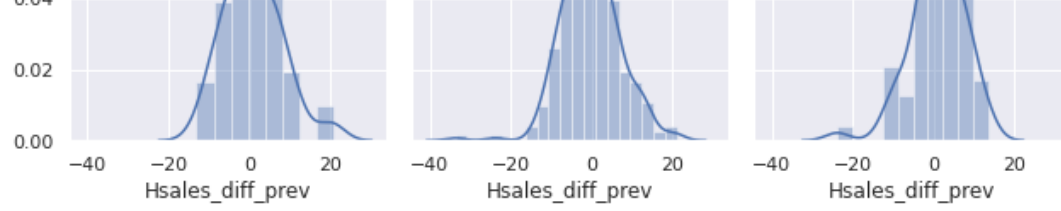


Distribution for Rsales\_diff\_year



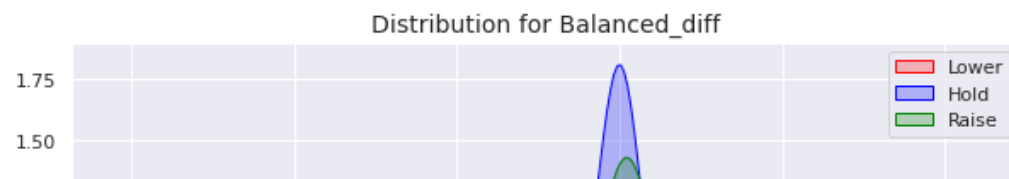
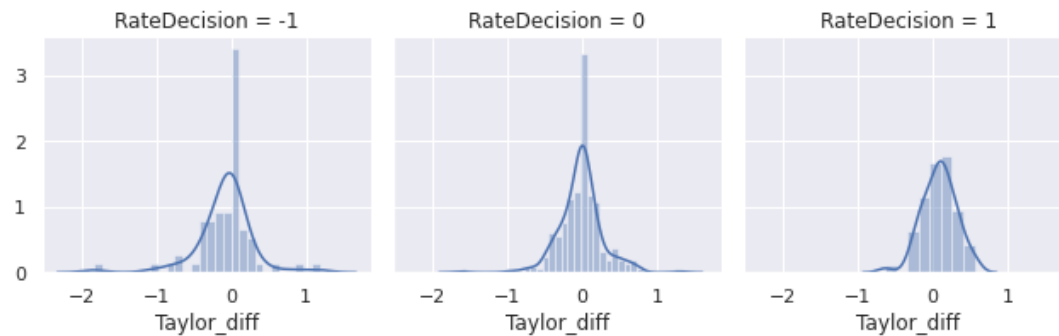
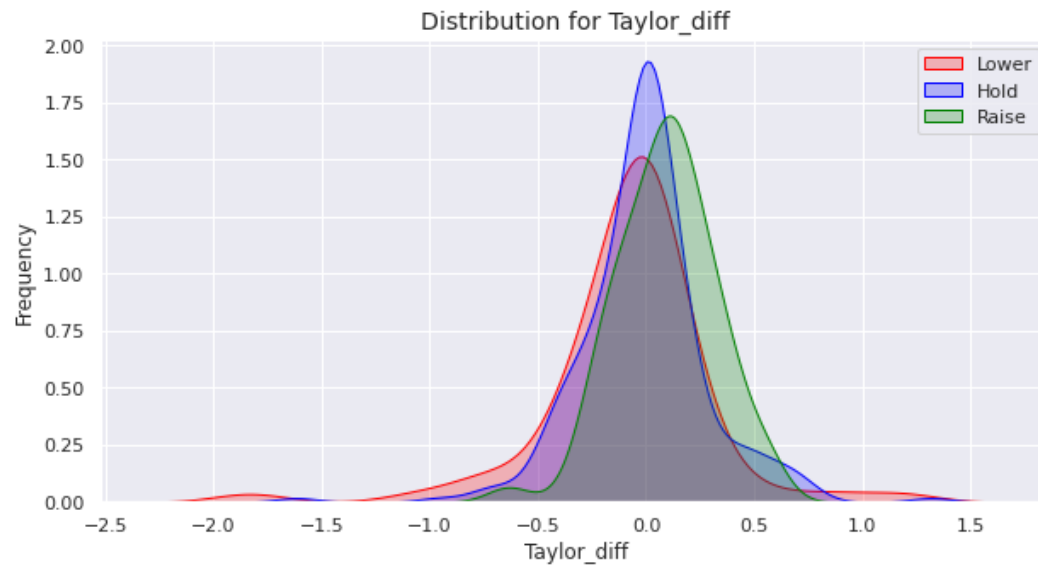
Distribution for Hsales\_diff\_prev

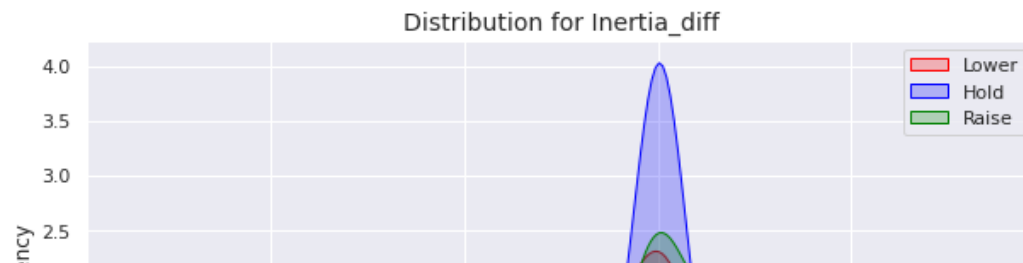
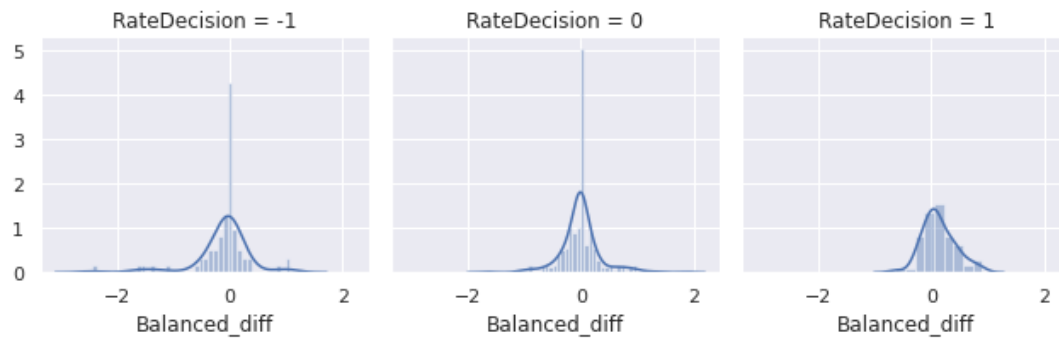
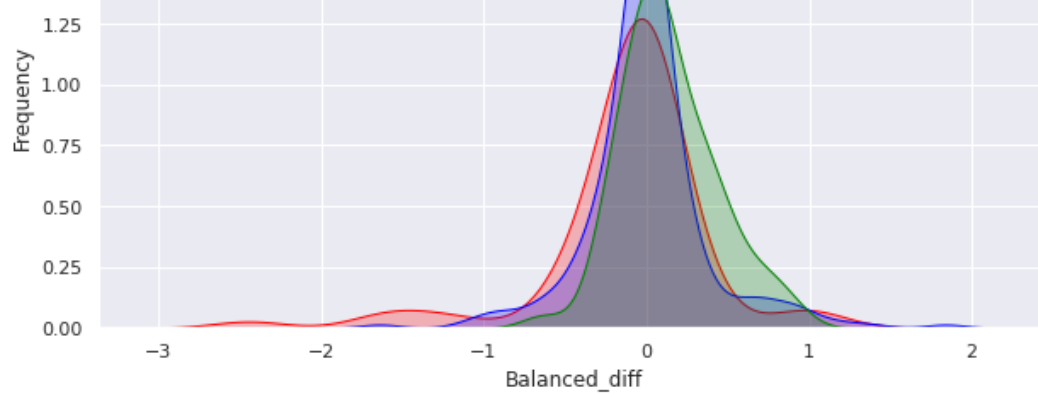




```
plot_distribution(nontext_data, ["Taylor_diff", "Balanced_diff", "Inertia_diff"])
```

```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future
warnings.warn(msg, FutureWarning)
```





## ▼ Create Training Data Set

```

nontext_train_small = pd.concat([nontext_data[['RateDecision', 'prev_decision', 'GDP_diff_prev', 'PMI_value']],
                                nontext_ma2[['Employ_diff_prev', 'Rsales_diff_year']],
                                nontext_ma3[['Unemp_diff_prev', 'Inertia_diff']],
                                nontext_ma12[['Hsales_diff_year', 'Balanced_diff']]], axis=1)
nontext_train_small.rename(columns={'RateDecision': 'target'}, inplace=True)

nontext_train_small.isnull().sum()

```

```

target          0
prev_decision    1
GDP_diff_prev    0
PMI_value        0
Employ_diff_prev  0
Rsales_diff_year 131
Unemp_diff_prev  0

```

```
Inertia_diff      4
Hsales_diff_year   0
Balanced_diff     15
dtype: int64
```

```
latest_columns = ['RateDecision',
                  'prev_decision',
                  'GDP_diff_prev',
                  'GDP_diff_year',
                  'GDPPOT_diff_prev',
                  'GDPPOT_diff_year',
                  'PCE_diff_prev',
                  'PCE_diff_year',
                  'CPI_diff_prev',
                  'CPI_diff_year',
                  'Unemp_value',
                  'Unemp_diff_prev',
                  'Unemp_diff_year',
                  'Employ_value',
                  'Employ_diff_prev',
                  'Employ_diff_year',
                  'PMI_value',
                  'PMI_diff_prev',
                  'PMI_diff_year',
                  'Rsales_diff_prev',
                  'Rsales_diff_year',
                  'Hsales_diff_prev',
                  'Hsales_diff_year',
                  'Taylor-Rate',
                  'Balanced-Rate',
                  'Inertia-Rate',
                  'Taylor_diff',
                  'Balanced_diff',
                  'Inertia_diff']
```

```
ma3_columns = [
    'GDP_diff_prev',
    'GDP_diff_year',
    'GDPPOT_diff_prev',
    'GDPPOT_diff_year',
    'PCE_diff_prev',
    'PCE_diff_year',
    'CPI_diff_prev',
    'CPI_diff_year',
    'Unemp_value',
    'Unemp_diff_prev',
    'Unemp_diff_year',
    'Employ_value',
    'Employ_diff_prev',
    'Employ_diff_year',
```

```

        'PMI_value',
        'PMI_diff_prev',
        'PMI_diff_year',
        'Rsales_diff_prev',
        'Rsales_diff_year',
        'Hsales_diff_prev',
        'Hsales_diff_year',
        'Taylor-Rate',
        'Balanced-Rate',
        'Inertia-Rate',
        'Taylor_diff',
        'Balanced_diff',
        'Inertia_diff'
    ]

```

```

nontext_train_large = pd.concat([nontext_data[latest_columns], nontext_ma3[ma3_columns].add_suffix('_ma3')], axis=1)
nontext_train_large.rename(columns={'RateDecision': 'target'}, inplace=True)
print(nontext_data[latest_columns].shape)
print(nontext_ma3[ma3_columns].shape)
print(nontext_train_large.shape)

```

```

(390, 29)
(390, 27)
(390, 56)

```

```

nontext_train_large.isnull().sum()

```

```

target                0
prev_decision          1
GDP_diff_prev         0
GDP_diff_year         0
GDPPOT_diff_prev      0
GDPPOT_diff_year      0
PCE_diff_prev         0
PCE_diff_year         0
CPI_diff_prev         0
CPI_diff_year         0
Unemp_value           0
Unemp_diff_prev       0
Unemp_diff_year       0
Employ_value          0
Employ_diff_prev      0
Employ_diff_year      0
PMI_value             0
PMI_diff_prev         0
PMI_diff_year         0
Rsales_diff_prev     117
Rsales_diff_year     129
Hsales_diff_prev      0
Hsales_diff_year      0
Taylor-Rate           0
Balanced-Rate         0
Inertia-Rate          0

```



```
Taylor_diff          1
Balanced_diff        1
Inertia_diff         1
GDP_diff_prev_ma3    0
GDP_diff_year_ma3    0
GDPPOT_diff_prev_ma3 0
GDPPOT_diff_year_ma3 0
PCE_diff_prev_ma3    0
PCE_diff_year_ma3    0
CPI_diff_prev_ma3    0
CPI_diff_year_ma3    0
Unemp_value_ma3      0
Unemp_diff_prev_ma3  0
Unemp_diff_year_ma3  0
Employ_value_ma3     0
Employ_diff_prev_ma3 0
Employ_diff_year_ma3 0
PMI_value_ma3        0
PMI_diff_prev_ma3    0
PMI_diff_year_ma3    0
Rsales_diff_prev_ma3 119
Rsales_diff_year_ma3 132
Hsales_diff_prev_ma3  0
Hsales_diff_year_ma3  0
Taylor-Rate_ma3      3
Balanced-Rate_ma3    3
Inertia-Rate_ma3     3
Taylor_diff_ma3      4
Balanced_diff_ma3    4
Inertia_diff_ma3     4
dtype: int64
```

## ▼ Missing Values

```
nontext_train_small['prev_decision'].fillna(0, inplace=True)
nontext_train_large['prev_decision'].fillna(0, inplace=True)
```

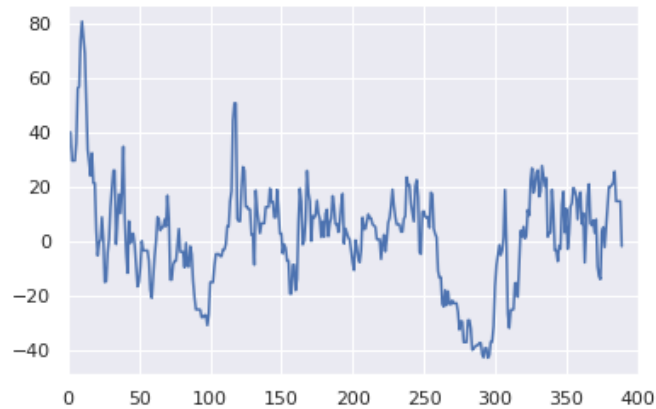
```
ax = sns.lineplot(data=nontext_train_small['Rsales_diff_year'].values)
ax.set_xlim(0, 400)
```

```
(0.0, 400.0)
```



```
ax = sns.lineplot(data=nontext_ma2['Hsales_diff_year'].values)
ax.set_xlim(0, 400)
```

```
(0.0, 400.0)
```



```
nontext_train_small['Rsales_diff_year'].fillna(nontext_train_small['Rsales_diff_year'].mean(), inplace=True)
nontext_train_large['Rsales_diff_prev'].fillna(nontext_train_large['Rsales_diff_prev'].mean(), inplace=True)
nontext_train_large['Rsales_diff_year'].fillna(nontext_train_large['Rsales_diff_year'].mean(), inplace=True)
```

```
nontext_train_small['Inertia_diff'].fillna(nontext_train_small['Inertia_diff'].mean(), inplace=True)
nontext_train_small['Balanced_diff'].fillna(nontext_train_small['Balanced_diff'].mean(), inplace=True)
nontext_train_large['Inertia_diff'].fillna(nontext_train_large['Inertia_diff'].mean(), inplace=True)
nontext_train_large['Balanced_diff'].fillna(nontext_train_large['Balanced_diff'].mean(), inplace=True)
nontext_train_large['Taylor_diff'].fillna(nontext_train_large['Taylor_diff'].mean(), inplace=True)
```

```
nontext_train_large['Rsales_diff_prev_ma3'].fillna(nontext_train_large['Rsales_diff_prev_ma3'].mean(), inplace=True)
nontext_train_large['Rsales_diff_year_ma3'].fillna(nontext_train_large['Rsales_diff_year_ma3'].mean(), inplace=True)
```

```
nontext_train_large['Inertia_diff_ma3'].fillna(nontext_train_large['Inertia_diff_ma3'].mean(), inplace=True)
nontext_train_large['Balanced_diff_ma3'].fillna(nontext_train_large['Balanced_diff_ma3'].mean(), inplace=True)
nontext_train_large['Taylor_diff_ma3'].fillna(nontext_train_large['Taylor_diff_ma3'].mean(), inplace=True)
nontext_train_large['Inertia-Rate_ma3'].fillna(nontext_train_large['Inertia-Rate_ma3'].mean(), inplace=True)
nontext_train_large['Balanced-Rate_ma3'].fillna(nontext_train_large['Balanced-Rate_ma3'].mean(), inplace=True)
nontext_train_large['Taylor-Rate_ma3'].fillna(nontext_train_large['Taylor-Rate_ma3'].mean(), inplace=True)
```

```
nontext_train_small.isnull().sum()
```

```
target          0
prev_decision    0
```

```
GDP_diff_prev      0
PMI_value          0
Employ_diff_prev   0
Rsales_diff_year   0
Unemp_diff_prev    0
Inertia_diff       0
Hsales_diff_year   0
Balanced_diff      0
dtype: int64
```

```
nontext_train_large.isnull().sum()
```

```
target            0
prev_decision     0
GDP_diff_prev     0
GDP_diff_year     0
GDPPOT_diff_prev  0
GDPPOT_diff_year  0
PCE_diff_prev     0
PCE_diff_year     0
CPI_diff_prev     0
CPI_diff_year     0
Unemp_value       0
Unemp_diff_prev   0
Unemp_diff_year   0
Employ_value      0
Employ_diff_prev  0
Employ_diff_year  0
PMI_value         0
PMI_diff_prev     0
PMI_diff_year     0
Rsales_diff_prev  0
Rsales_diff_year  0
Hsales_diff_prev  0
Hsales_diff_year  0
Taylor-Rate       0
Balanced-Rate     0
Inertia-Rate      0
Taylor_diff       0
Balanced_diff     0
Inertia_diff      0
GDP_diff_prev_ma3 0
GDP_diff_year_ma3 0
GDPPOT_diff_prev_ma3 0
GDPPOT_diff_year_ma3 0
PCE_diff_prev_ma3 0
PCE_diff_year_ma3 0
CPI_diff_prev_ma3 0
CPI_diff_year_ma3 0
Unemp_value_ma3   0
Unemp_diff_prev_ma3 0
Unemp_diff_year_ma3 0
Employ_value_ma3  0
Employ_diff_prev_ma3 0
```

```

Employ_diff_year_ma3      0
PMI_value_ma3             0
PMI_diff_prev_ma3        0
PMI_diff_year_ma3        0
Rsales_diff_prev_ma3     0
Rsales_diff_year_ma3     0
Hsales_diff_prev_ma3     0
Hsales_diff_year_ma3     0
Taylor-Rate_ma3          0
Balanced-Rate_ma3        0
Inertia-Rate_ma3         0
Taylor_diff_ma3          0
Balanced_diff_ma3        0
Inertia_diff_ma3        0
dtype: int64

```

## One-hot encoding

```

# nontext_train['Lower'] = nontext_train['RateDecision'].apply(lambda x: 1 if x == -1 else 0)
# nontext_train['Hold'] = nontext_train['RateDecision'].apply(lambda x: 1 if x == 0 else 0)
# nontext_train['Raise'] = nontext_train['RateDecision'].apply(lambda x: 1 if x == 1 else 0)
# nontext_train

```

## Save Data

```

if IN_COLAB:
    def save_data(df, file_name, dir_name=preprocessed_dir, index_csv=True):
        if not os.path.exists(dir_name):
            os.mkdir(dir_name)
        # Save results to a picke file
        file = open(dir_name + file_name + '.pickle', 'wb')
        pickle.dump(df, file)
        file.close()
        print('Successfully saved {}.pickle. in {}'.format(file_name, dir_name + file_name + '.pickle'))
        # Save results to a csv file
        df.to_csv(dir_name + file_name + '.csv', index=index_csv)
        print('Successfully saved {}.csv. in {}'.format(file_name, dir_name + file_name + '.csv'))

else:
    def save_data(df, file_name, dir_name=preprocessed_dir, index_csv=True):
        # Save results to a .picke file
        file = open(dir_name + file_name + '.pickle', 'wb')
        pickle.dump(df, file)
        file.close()
        print('Successfully saved {}.pickle. in {}'.format(file_name, dir_name + file_name + '.pickle'))
        # Save results to a .csv file
        df.to_csv(dir_name + file_name + '.csv', index=index_csv)

```

```
print('Successfully saved {}.csv. in {}'.format(file_name, dir_name + file_name + '.csv'))
```

```
# Save non-text data
```

```
save_data(nontext_train_small, 'nontext_train_small')
```

```
save_data(nontext_train_large, 'nontext_train_large')
```

```
Successfully saved nontext_train_small.pickle. in /content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_train_small.pickle
```

```
Successfully saved nontext_train_small.csv. in /content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_train_small.csv
```

```
Successfully saved nontext_train_large.pickle. in /content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_train_large.pickle
```

```
Successfully saved nontext_train_large.csv. in /content/drive/My Drive/Colab Notebooks/proj2/src/data/preprocessed/nontext_train_large.csv
```