

# Обработка пропусков в данных, кодирование категориальных признаков, масштабирование данных.

Для данной лабораторной работы я использую набор данных индекса свободы человека по странам за 2018 год:

## Контекст

Центральная цель Индекса свободы человека заключается в том, чтобы нарисовать широкую, но достаточно точную картину масштабов общей свободы в мире. Более широкая цель состоит в том, чтобы более тщательно изучить, что мы подразумеваем под свободой, и лучше понять ее связь с любым количеством других социально-экономических явлений.

## Содержание

Индекс свободы человека измеряет экономические свободы, такие как свобода торговли или использования здоровых денег, и отражает степень, в которой люди свободны свободно пользоваться основными свободами, часто называемыми гражданскими свободами - свободой слова, религии, ассоциаций и собраний - в странах, участвующих в обзоре. Кроме того, он включает в себя показатели верховенства права, преступности и насилия, свободы передвижения и правовой дискриминации в отношении однополых отношений. Мы также включаем девять переменных, относящихся к свободам, специфичным для женщин, которые содержатся в различных категориях индекса.

Набор данных имеет следующие столбцы:

- year - год
- ISO\_code - буквенное обозначение страны
- countries - страны
- region - регион
- pf\_rol\_procedural - процедурное правосудие
- pf\_rol\_civil - гражданская справедливость
- pf\_rol\_criminal - уголовное правосудие
- pf\_rol - верховенство права
- pf\_ss\_homicide - убийству
- pf\_ss\_disappearances\_disap - исчезновение
- ...

```
In [2]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
%matplotlib inline
sns.set(style="ticks")
```

## Первичный анализ данных

```
In [17]: # читаем данные набора
data = pd.read_csv('JupyterNotebooks/data/hfi_cc_2018.csv', sep=",")
```

```
In [18]: data.shape
```

```
Out[18]: (1458, 123)
```

```
In [19]: # типы данных столбцов
data.dtypes
```

```
Out[19]: year                int64
ISO_code                   object
countries                  object
region                    object
pf_rol_procedural         float64
...
ef_score                  float64
ef_rank                   float64
hf_score                  float64
hf_rank                   float64
hf_quartile               float64
Length: 123, dtype: object
```

```
In [21]: # суммы пропусков по столбцам
data.isnull().sum()
```

```
Out[21]: year                0
ISO_code                   0
countries                  0
region                    0
pf_rol_procedural         578
...
ef_score                  80
ef_rank                   80
hf_score                  80
hf_rank                   80
hf_quartile               80
Length: 123, dtype: int64
```

```
In [29]: data.head()
```

```
Out[29]:
```

	year	ISO_code	countries	region	pf_rol_procedural	pf_rol_civil	pf_rol_criminal	
0	2016	ALB	Albania	Eastern Europe	6.661503	4.547244	4.666508	5.2
1	2016	DZA	Algeria	Middle East & North Africa	NaN	NaN	NaN	3.8
2	2016	AGO	Angola	Sub-Saharan Africa	NaN	NaN	NaN	3.4
3	2016	ARG	Argentina	Latin America & the Caribbean	7.098483	5.791960	4.343930	5.7
4	2016	ARM	Armenia	Caucasus	NaN	NaN	NaN	5.0

5 rows × 123 columns

```
In [30]: total_count = data.shape[0]
```

## Обработка пропусков значений в наборе данных

Удалим те строки данных, в которых встречаются пустые столбцы:

```
In [22]: # Удаление строк, содержащих пустые значения
data_new_2 = data.dropna(axis=0, how='any')
(data.shape, data_new_2.shape)
```

```
Out[22]: ((1458, 123), (0, 123))
```

Итоговая сетка данных оказывается пустой, так как в каждой строке так или иначе встречаются пропуски данных. В таком случае попробуем заполнить пустые значения например нулями.

```
In [26]: # Заполнение всех пропущенных значений нулями
data_new_3 = data.fillna(0)
data_new_3.head()
```

```
Out[26]:
```

	year	ISO_code	countries	region	pf_rol_procedural	pf_rol_civil	pf_rol_criminal	
0	2016	ALB	Albania	Eastern Europe	6.661503	4.547244	4.666508	5.2
1	2016	DZA	Algeria	Middle East & North Africa	0.000000	0.000000	0.000000	3.8
2	2016	AGO	Angola	Sub-Saharan Africa	0.000000	0.000000	0.000000	3.4
3	2016	ARG	Argentina	Latin America & the Caribbean	7.098483	5.791960	4.343930	5.7
4	2016	ARM	Armenia	Caucasus & Central Asia	0.000000	0.000000	0.000000	5.0

5 rows × 123 columns

Это не очень корректно, так как и категориальные данные заполняются нулями.

## Импьютация

### Обработка числовых данных

```
In [31]: # Выберем числовые колонки с пропущенными значениями
# Цикл по колонкам датасета
num_cols = []
for col in data.columns:
```

```
# Количество пустых значений
temp_null_count = data[data[col].isnull()].shape[0]
dt = str(data[col].dtype)
if temp_null_count>0 and (dt=='float64' or dt=='int64'):
    num_cols.append(col)
    temp_perc = round((temp_null_count / total_count) * 100.0, 2)
    print('Колонка {}'.format(col).ljust(15) + 'Тип данных {}'.format(dt).ljust(15) + 'Количество пустых значений {}'.format(temp_perc).ljust(15))
```

Колонка pf\_rol\_procedural. Тип данных float64. Количество пустых значений 578, 39.64%.

Колонка pf\_rol\_civil. Тип данных float64. Количество пустых значений 578, 39.64%.

Колонка pf\_rol\_criminal. Тип данных float64. Количество пустых значений 578, 39.64%.

Колонка pf\_rol. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_ss\_homicide. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_ss\_disappearances\_disap. Тип данных float64. Количество пустых значений 89, 6.1%.

Колонка pf\_ss\_disappearances\_violent. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_ss\_disappearances\_organized. Тип данных float64. Количество пустых значений 179, 12.28%.

Колонка pf\_ss\_disappearances\_fatalities. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_ss\_disappearances\_injuries. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_ss\_disappearances. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_ss\_women\_fgm. Тип данных float64. Количество пустых значений 172, 11.8%.

Колонка pf\_ss\_women\_missing. Тип данных float64. Количество пустых значений 120, 8.23%.

Колонка pf\_ss\_women\_inheritance\_widows. Тип данных float64. Количество пустых значений 541, 37.11%.

Колонка pf\_ss\_women\_inheritance\_daughters. Тип данных float64. Количество пустых значений 541, 37.11%.

Колонка pf\_ss\_women\_inheritance. Тип данных float64. Количество пустых значений 119, 8.16%.

Колонка pf\_ss\_women. Тип данных float64. Количество пустых значений 100, 6.86%.

Колонка pf\_ss. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_movement\_domestic. Тип данных float64. Количество пустых значений 98, 6.72%.

Колонка pf\_movement\_foreign. Тип данных float64. Количество пустых значений 98, 6.72%.

Колонка pf\_movement\_women. Тип данных float64. Количество пустых значений 141, 9.67%.

Колонка pf\_movement. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_religion\_estop\_establish. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_religion\_estop\_operate. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_religion\_estop. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_religion\_harassment. Тип данных float64. Количество пустых значений 94, 6.45%.

Колонка pf\_religion\_restrictions. Тип данных float64. Количество пустых значений 94, 6.45%.

Колонка pf\_religion. Тип данных float64. Количество пустых значений 90, 6.17%.

Колонка pf\_association\_association. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_association\_assembly. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_association\_political\_establish. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_association\_political\_operate. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_association\_political. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_association\_prof\_establish. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_association\_prof\_operate. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_association\_prof. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_association\_sport\_establish. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_association\_sport\_operate. Тип данных float64. Количество пустых значений 731, 50.14%.

Колонка pf\_association\_sport. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_association. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_expression\_killed. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_expression\_jailed. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_expression\_influence. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_expression\_control. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_expression\_cable. Тип данных float64. Количество пустых значений 335, 22.98%.

Колонка pf\_expression\_newspapers. Тип данных float64. Количество пустых значений 335, 22.98%.

Колонка pf\_expression\_internet. Тип данных float64. Количество пустых значений 329, 22.57%.

Колонка pf\_expression. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_identity\_legal. Тип данных float64. Количество пустых значений 1253, 85.94%.

Колонка pf\_identity\_parental\_marriage. Тип данных float64. Количество пустых значений 535, 36.69%.

Колонка pf\_identity\_parental\_divorce. Тип данных float64. Количество пустых значений 535, 36.69%.

Колонка pf\_identity\_parental. Тип данных float64. Количество пустых значений 100, 6.86%.

Колонка pf\_identity\_sex\_male. Тип данных float64. Количество пустых значений 83, 5.69%.

Колонка pf\_identity\_sex\_female. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_identity\_sex. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_identity\_divorce. Тип данных float64. Количество пустых значений 873, 59.88%.

Колонка pf\_identity. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_score. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка pf\_rank. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_government\_consumption. Тип данных float64. Количество пустых значений 66, 4.53%.

Колонка ef\_government\_transfers. Тип данных float64. Количество пустых значений 160, 10.97%.

Колонка ef\_government\_enterprises. Тип данных float64. Количество пустых значений 104, 7.13%.

Колонка ef\_government\_tax\_income. Тип данных float64. Количество пустых значений 124, 8.5%.

Колонка ef\_government\_tax\_payroll. Тип данных float64. Количество пустых значений 193, 13.24%.

Колонка ef\_government\_tax. Тип данных float64. Количество пустых значений 124, 8.5%.

Колонка ef\_government. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_legal\_judicial. Тип данных float64. Количество пустых значений 167, 11.45%.

Колонка ef\_legal\_courts. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_legal\_protection. Тип данных float64. Количество пустых значений 169, 11.59%.

Колонка ef\_legal\_military. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_legal\_integrity. Тип данных float64. Количество пустых значений 277, 19.0%.

Колонка ef\_legal\_enforcement. Тип данных float64. Количество пустых значений 90, 6.17%.

Колонка ef\_legal\_restrictions. Тип данных float64. Количество пустых значений 100, 6.86%.

Колонка ef\_legal\_police. Тип данных float64. Количество пустых значений 169, 11.59%.

Колонка ef\_legal\_crime. Тип данных float64. Количество пустых значений 169, 11.59%.

Колонка ef\_legal\_gender. Тип данных float64. Количество пустых значений 24, 1.65%.

Колонка ef\_legal. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_money\_growth. Тип данных float64. Количество пустых значений 70, 4.8%.

Колонка ef\_money\_sd. Тип данных float64. Количество пустых значений 72, 4.94%.

Колонка ef\_money\_inflation. Тип данных float64. Количество пустых значений 72, 4.94%.

Колонка ef\_money\_currency. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_money. Тип данных float64. Количество пустых значений 82, 5.62%.

Колонка ef\_trade\_tariffs\_revenue. Тип данных float64. Количество пустых значений 169, 11.59%.

Колонка ef\_trade\_tariffs\_mean. Тип данных float64. Количество пустых значений 92, 6.31%.

Колонка ef\_trade\_tariffs\_sd. Тип данных float64. Количество пустых значений 91, 6.24%.

Колонка ef\_trade\_tariffs. Тип данных float64. Количество пустых значений 85, 5.83%.

Колонка ef\_trade\_regulatory\_nontariff. Тип данных float64. Количество пустых значений 170, 11.66%.

Колонка ef\_trade\_regulatory\_compliance. Тип данных float64. Количество пустых значений 90, 6.17%.

Колонка ef\_trade\_regulatory. Тип данных float64. Количество пустых значений 84, 5.76%.

Колонка ef\_trade\_black. Тип данных float64. Количество пустых значений 87, 5.97%.

Колонка ef\_trade\_movement\_foreign. Тип данных float64. Количество пустых значений 164, 11.25%.

Колонка ef\_trade\_movement\_capital. Тип данных float64. Количество пустых значений 89, 6.1%.

Колонка ef\_trade\_movement\_visit. Тип данных float64. Количество пустых значений 85, 5.83%.

Колонка ef\_trade\_movement. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_trade. Тип данных float64. Количество пустых значений 81, 5.56%.

Колонка ef\_regulation\_credit\_ownership. Тип данных float64. Количество пустых значений 172, 11.8%.

Колонка ef\_regulation\_credit\_private. Тип данных float64. Количество пустых значений 72, 4.94%.

Колонка ef\_regulation\_credit\_interest. Тип данных float64. Количество пустых значений 100, 6.86%.

Колонка ef\_regulation\_credit. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_regulation\_labor\_minwage. Тип данных float64. Количество пустых значений 91, 6.24%.

Колонка ef\_regulation\_labor\_firing. Тип данных float64. Количество пустых значений 171, 11.73%.

Колонка ef\_regulation\_labor\_bargain. Тип данных float64. Количество пустых значений 170, 11.66%.

Колонка ef\_regulation\_labor\_hours. Тип данных float64. Количество пустых значений 88, 6.04%.

Колонка ef\_regulation\_labor\_dismissal. Тип данных float64. Количество пустых значений 110, 7.54%.

Колонка ef\_regulation\_labor\_conscription. Тип данных float64. Количество пустых значений 81, 5.56%.

Колонка ef\_regulation\_labor. Тип данных float64. Количество пустых значений 84, 5.76%.

Колонка ef\_regulation\_business\_adm. Тип данных float64. Количество пустых значений 169, 11.59%.

Колонка ef\_regulation\_business\_bureaucracy. Тип данных float64. Количество пустых значений 102, 7.0%.

Колонка ef\_regulation\_business\_start. Тип данных float64. Количество пустых значений 90, 6.17%.

Колонка ef\_regulation\_business\_bribes. Тип данных float64. Количество пустых значений 175, 12.0%.

Колонка ef\_regulation\_business\_licensing. Тип данных float64. Количество пустых значений 101, 6.93%.

Колонка ef\_regulation\_business\_compliance. Тип данных float64. Количество пустых значений 90, 6.17%.

Колонка ef\_regulation\_business. Тип данных float64. Количество пустых значений 84, 5.76%.

Колонка ef\_regulation. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_score. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка ef\_rank. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка hf\_score. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка hf\_rank. Тип данных float64. Количество пустых значений 80, 5.49%.

Колонка hf\_quartile. Тип данных float64. Количество пустых значений 80, 5.49%.

```
In [32]: # Фильтр по колонкам с пропущенными значениями
data_num = data[num_cols]
data_num
```

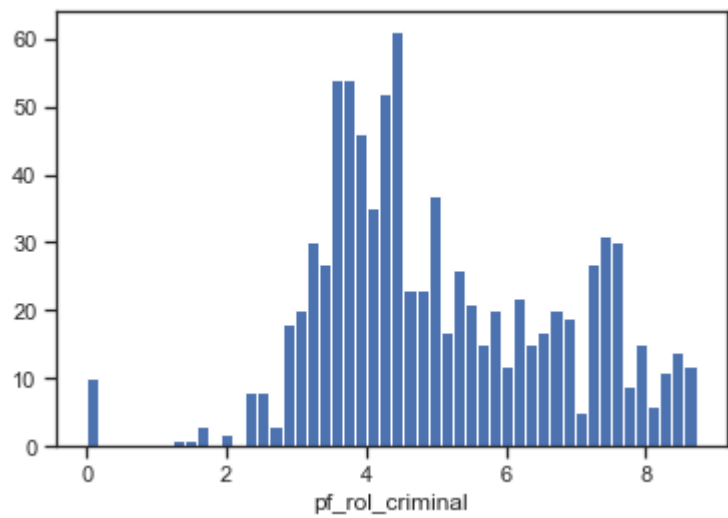
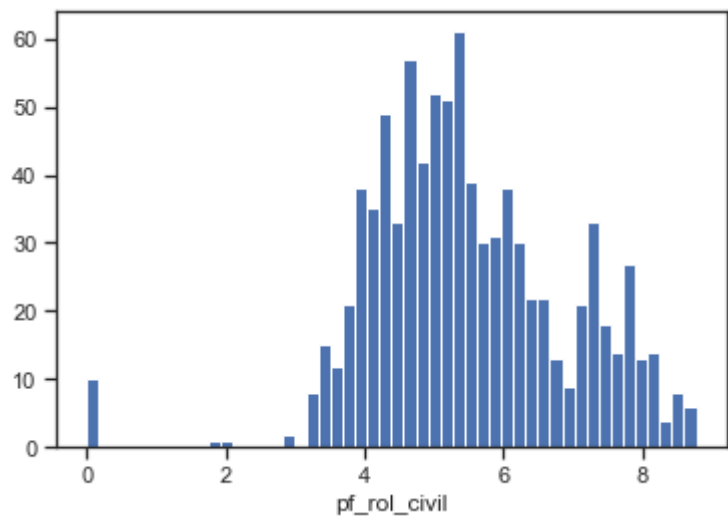
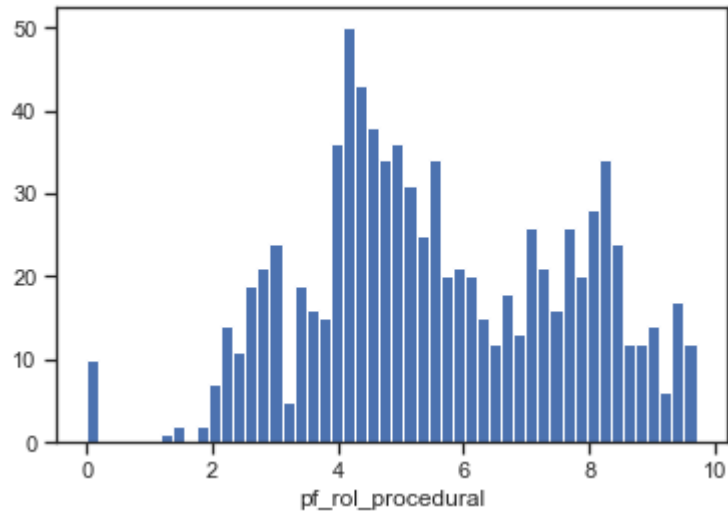
```
Out[32]:
```

	pf_rol_procedural	pf_rol_civil	pf_rol_criminal	pf_rol	pf_ss_homicide	pf_ss_disapp
0	6.661503	4.547244	4.666508	5.291752	8.920429	
1	NaN	NaN	NaN	3.819566	9.456254	
2	NaN	NaN	NaN	3.451814	8.060260	
3	7.098483	5.791960	4.343930	5.744791	7.622974	
4	NaN	NaN	NaN	5.003205	8.808750	
...	...	...	...	...	...	...
1453	3.000000	3.781688	2.369239	3.100000	0.000000	
1454	6.666667	4.349101	5.694847	5.600000	9.496239	
1455	NaN	NaN	NaN	NaN	NaN	
1456	4.800000	4.578003	3.688652	4.400000	7.878084	
1457	2.700000	3.991582	4.327660	3.700000	7.981019	

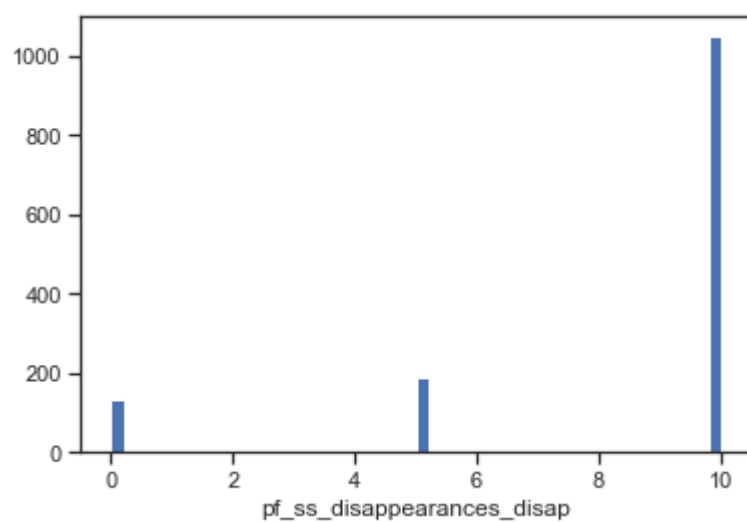
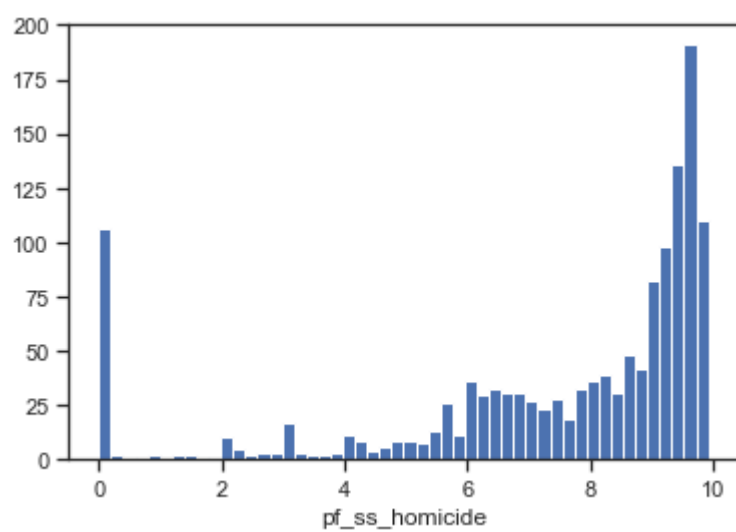
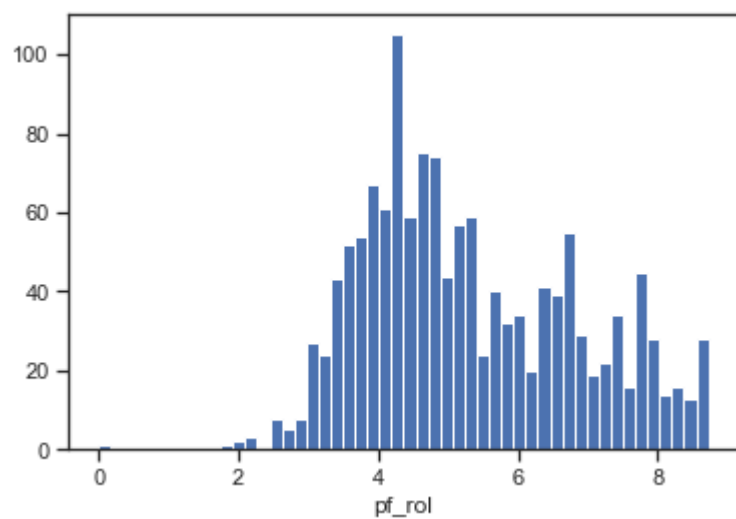
1458 rows × 119 columns

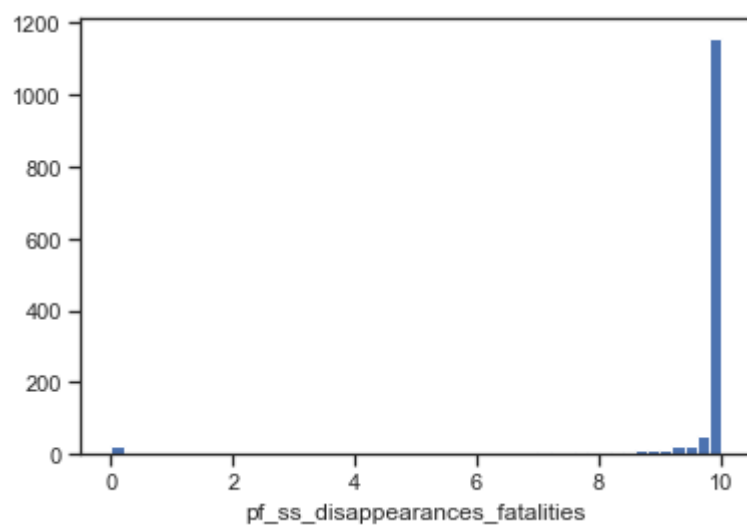
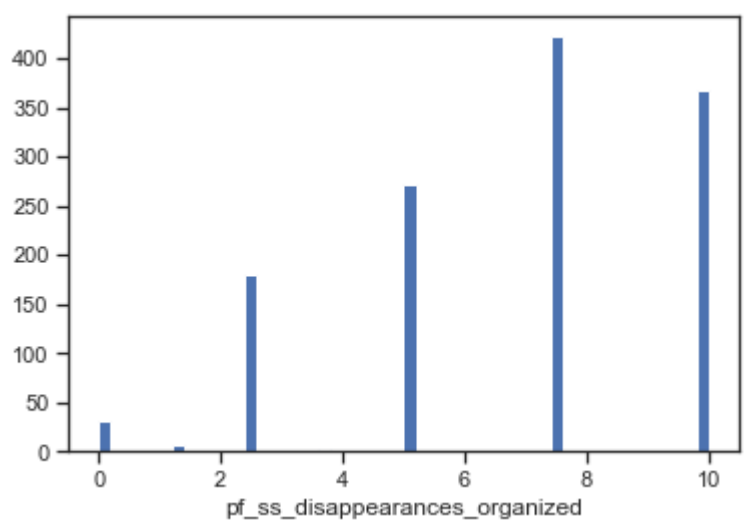
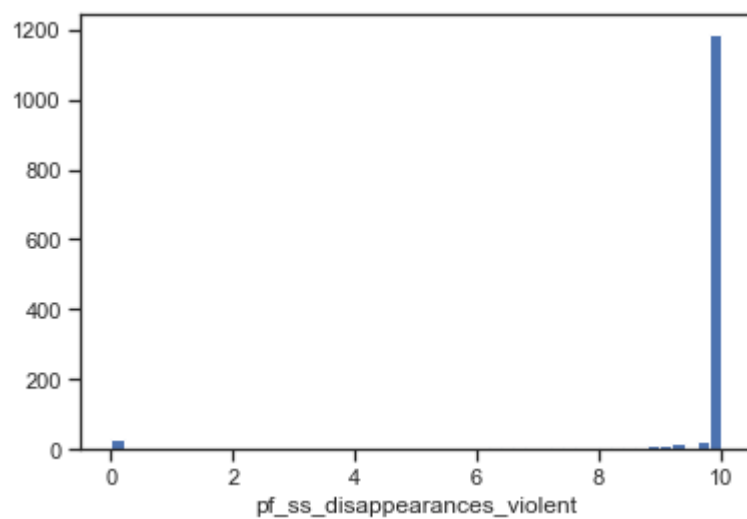
119 колонок из 123 являются численными и имеют пропуски данных.

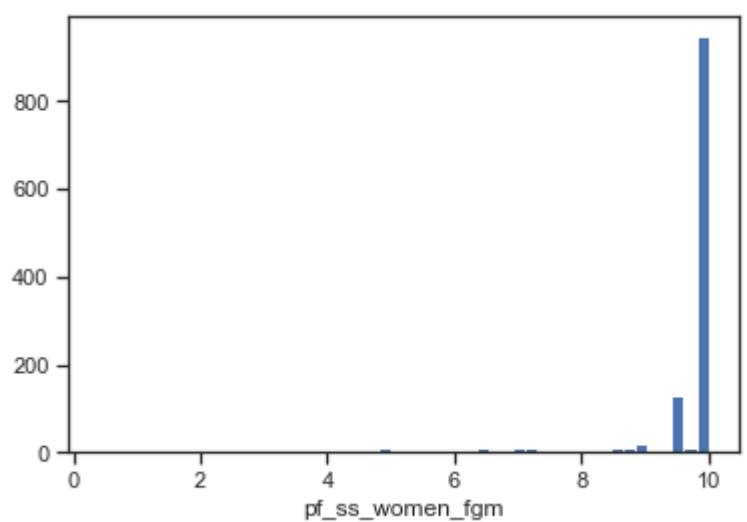
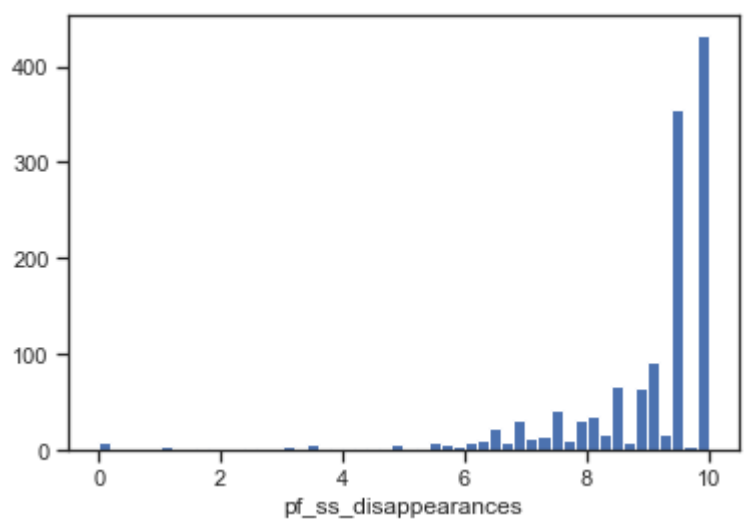
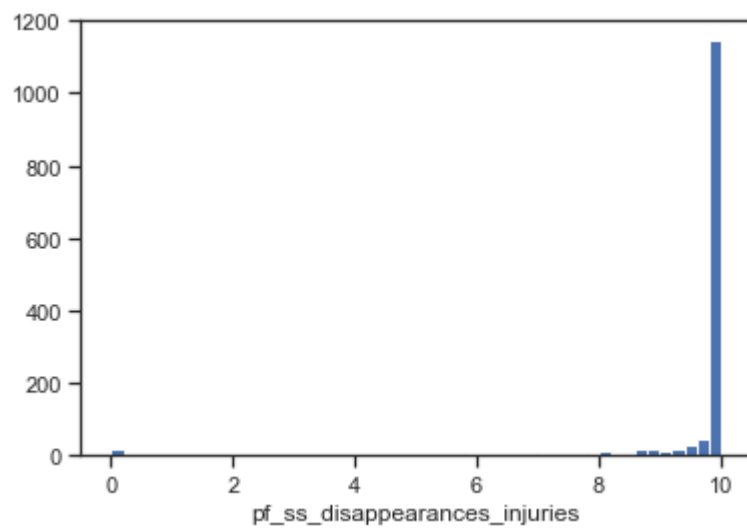
```
In [35]: # Гистограмма по признакам
for col in data_num:
    plt.hist(data[col], 50)
    plt.xlabel(col)
    plt.show()
```

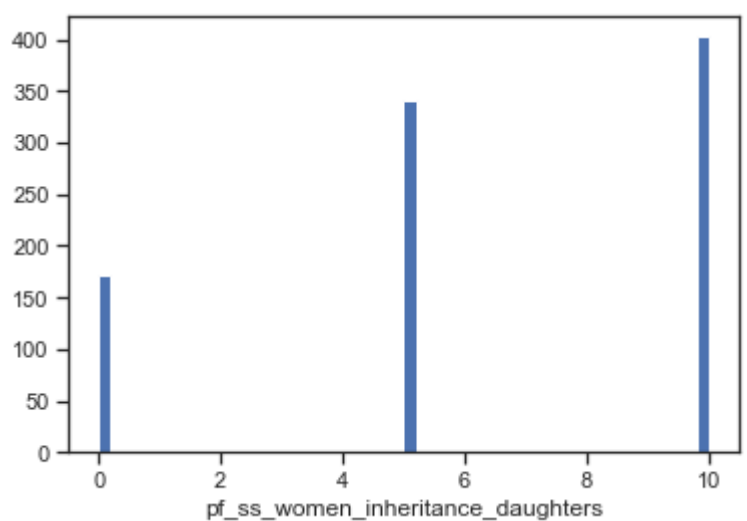
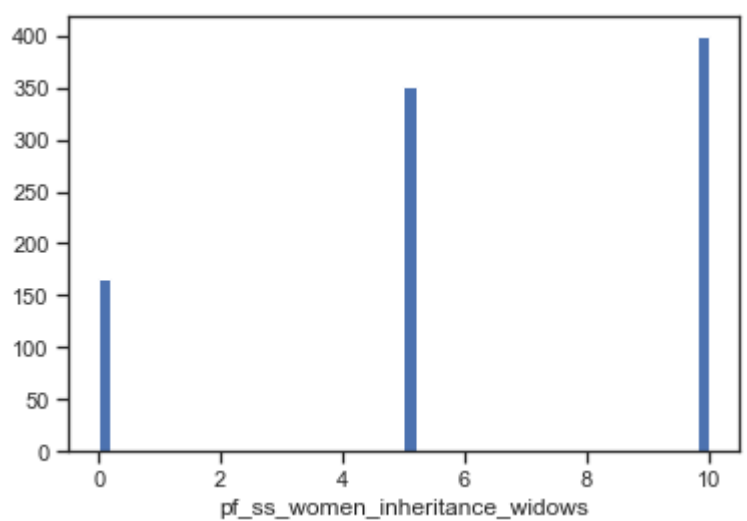
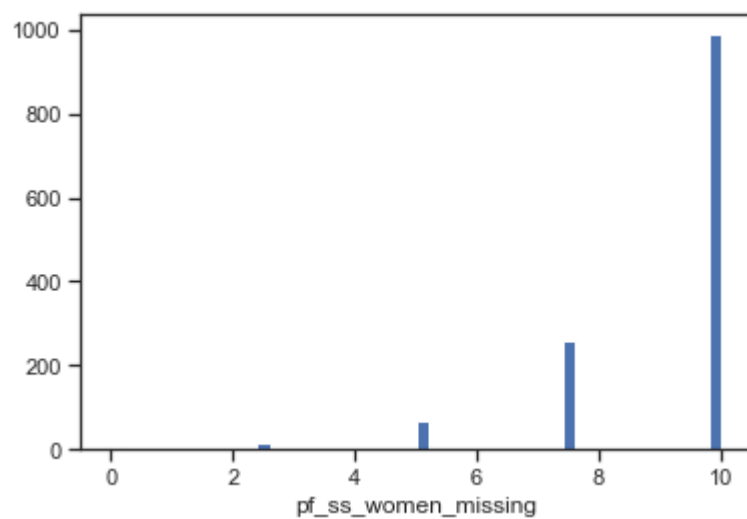


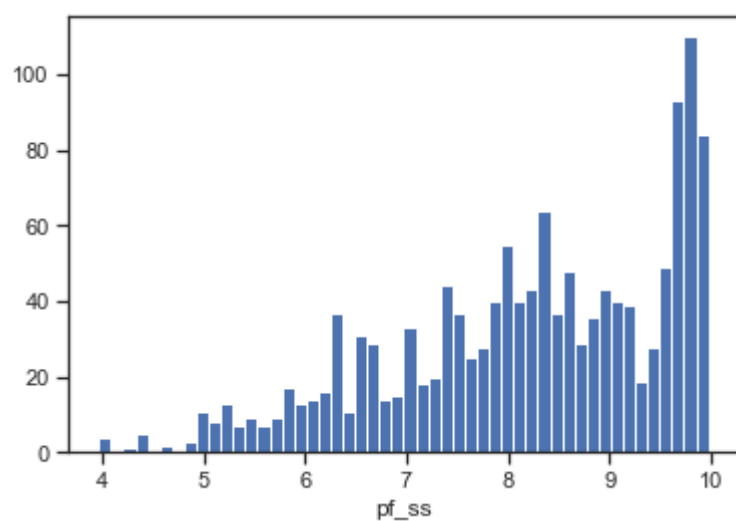
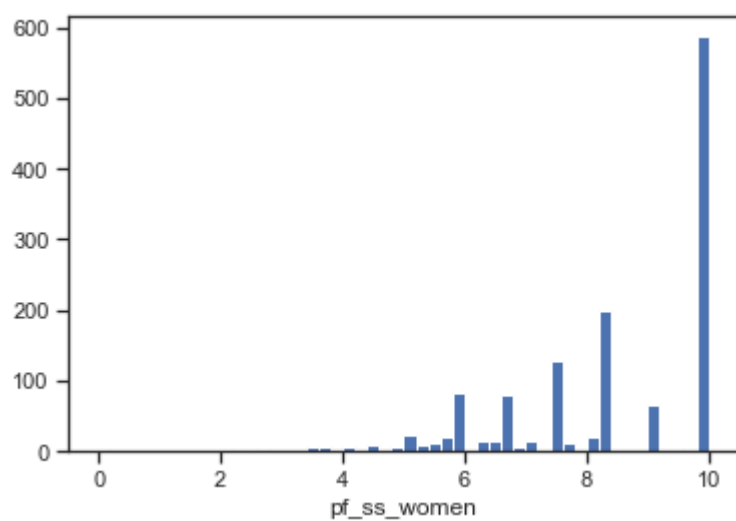
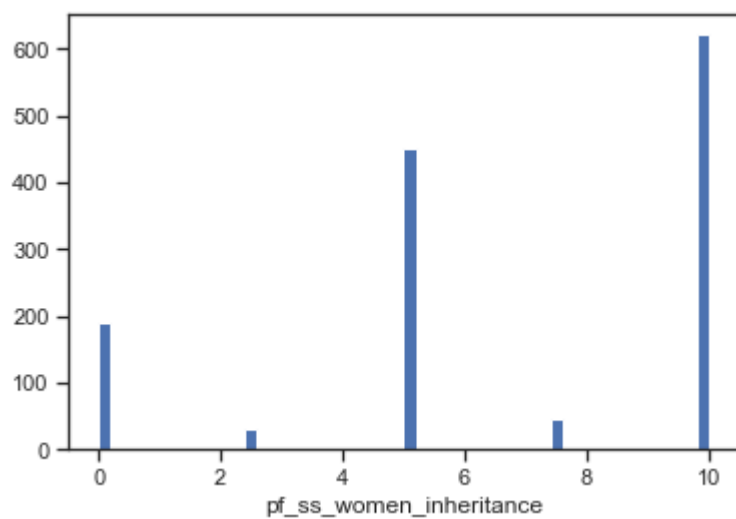


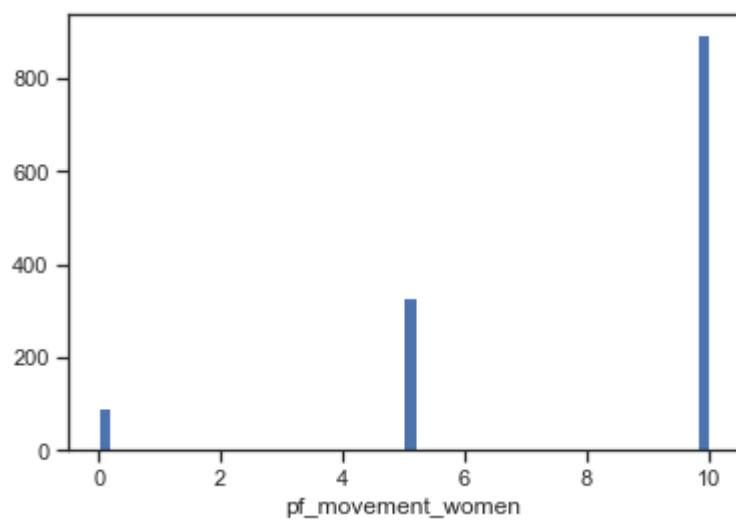
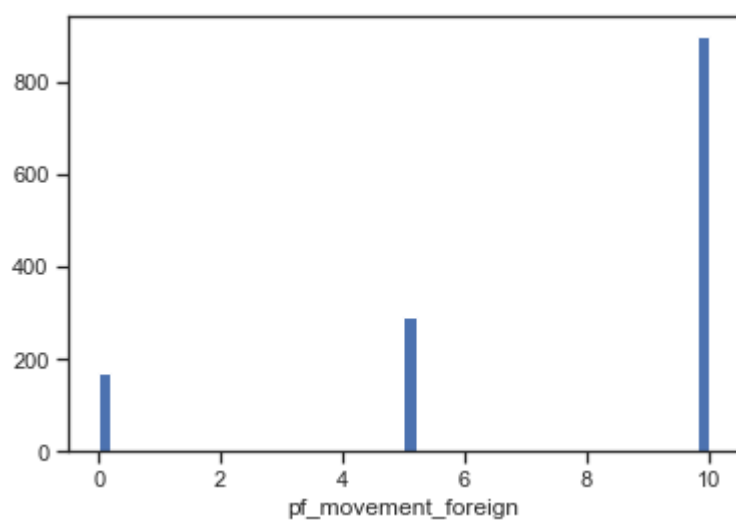
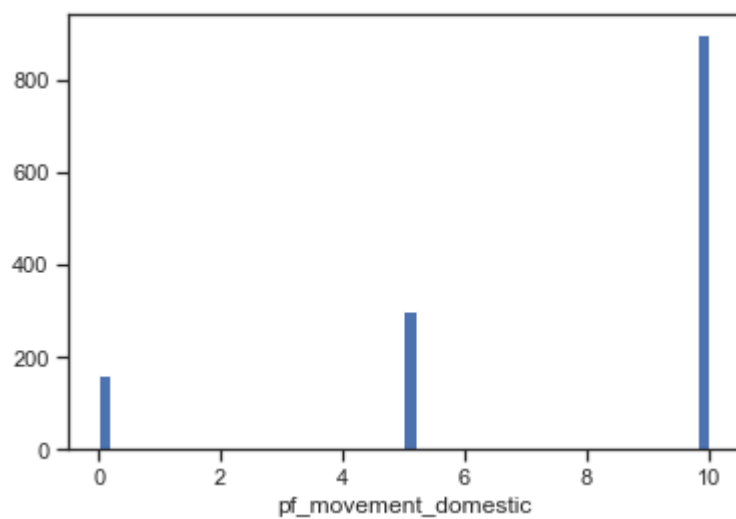


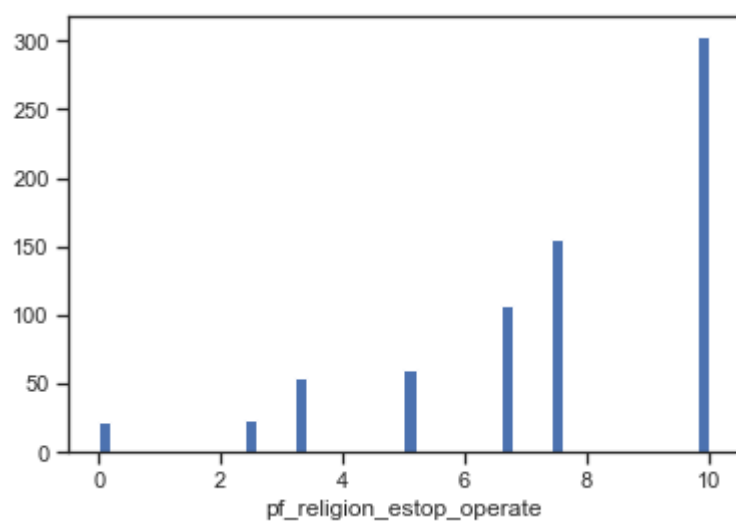
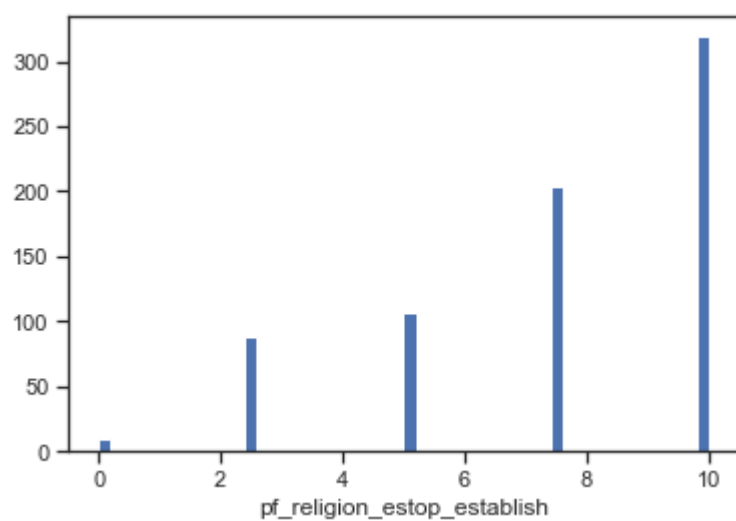
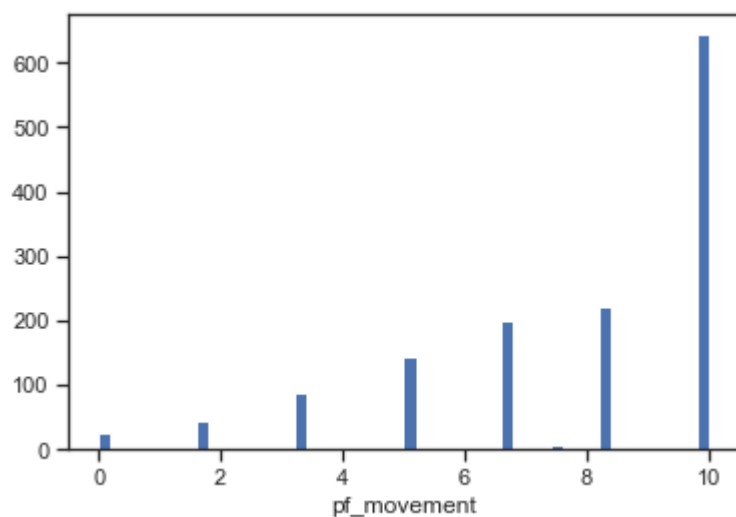


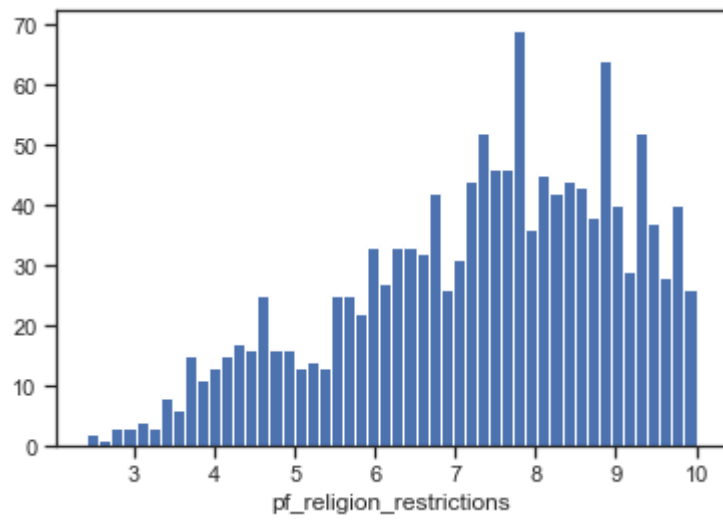
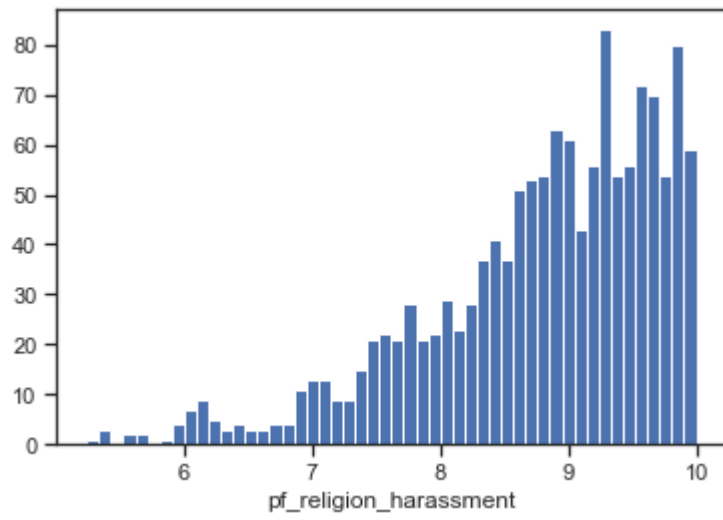
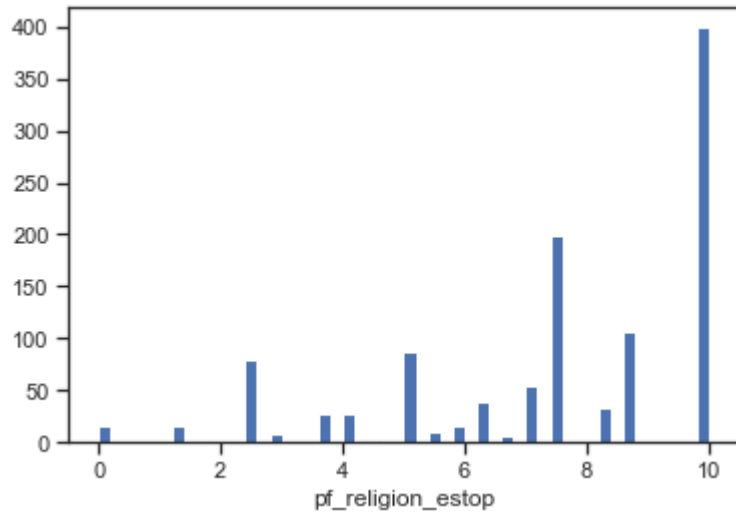




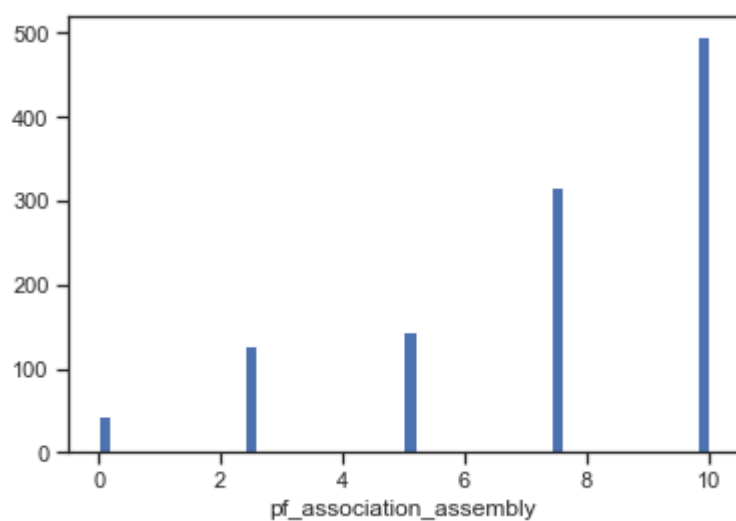
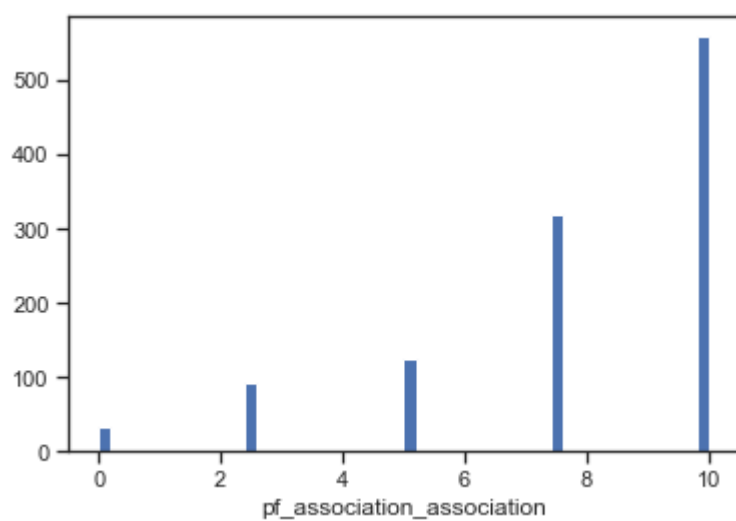
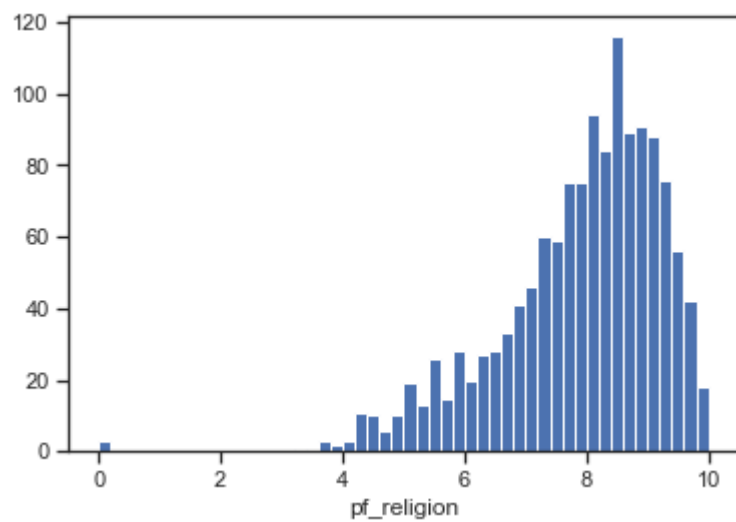


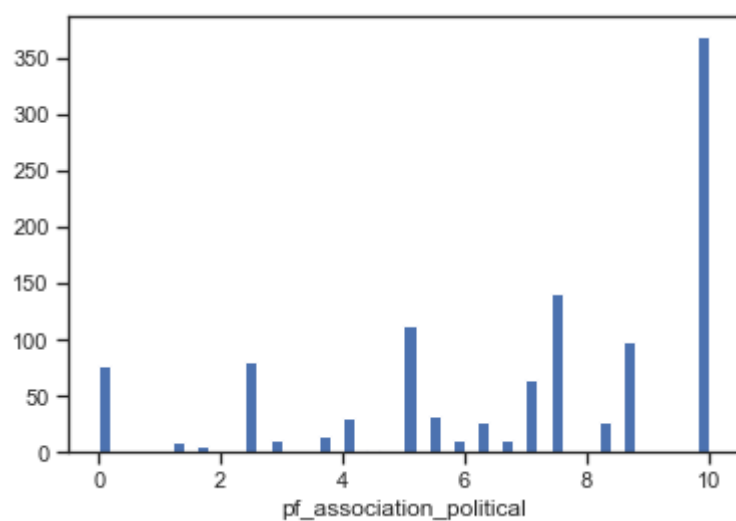
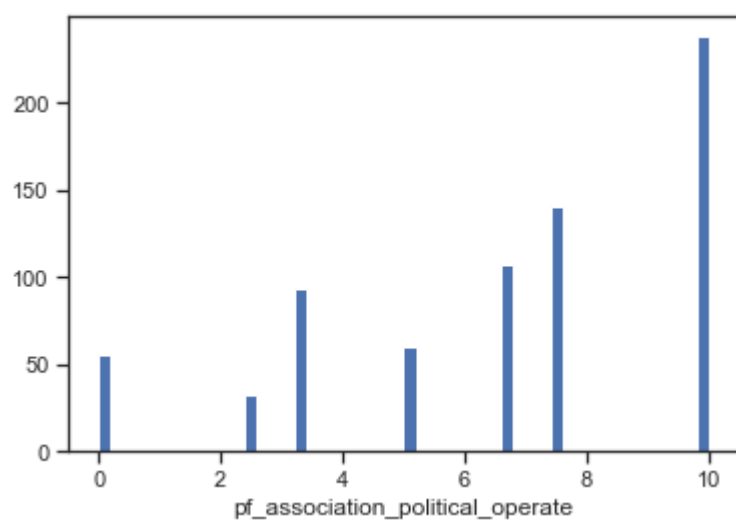
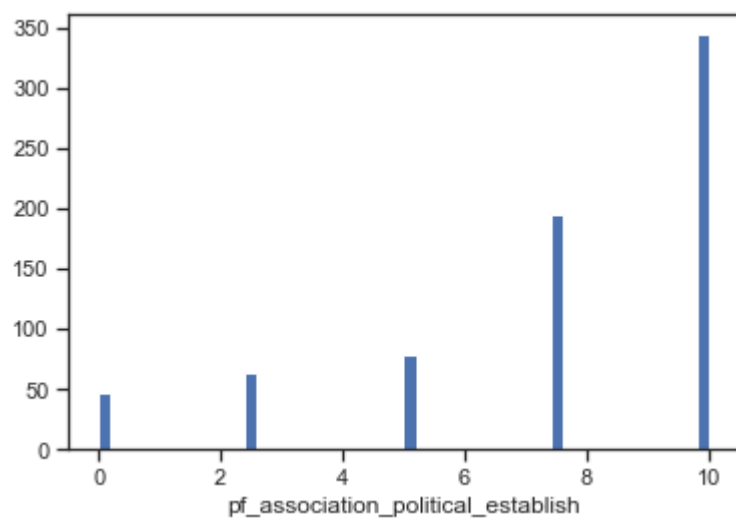


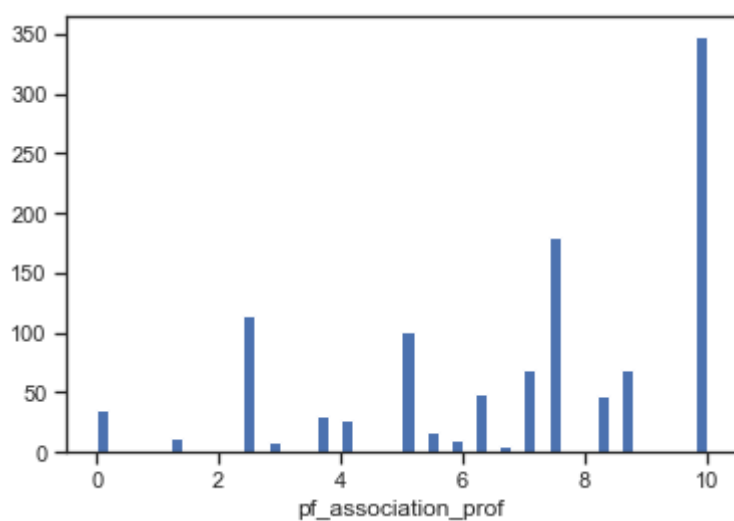
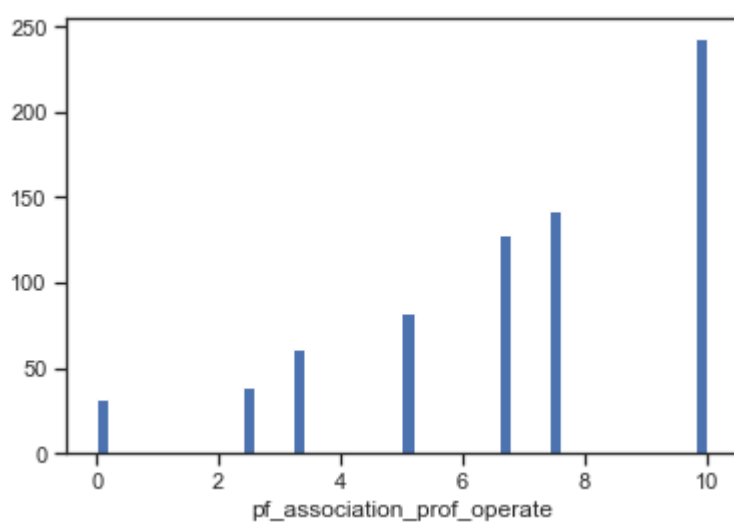
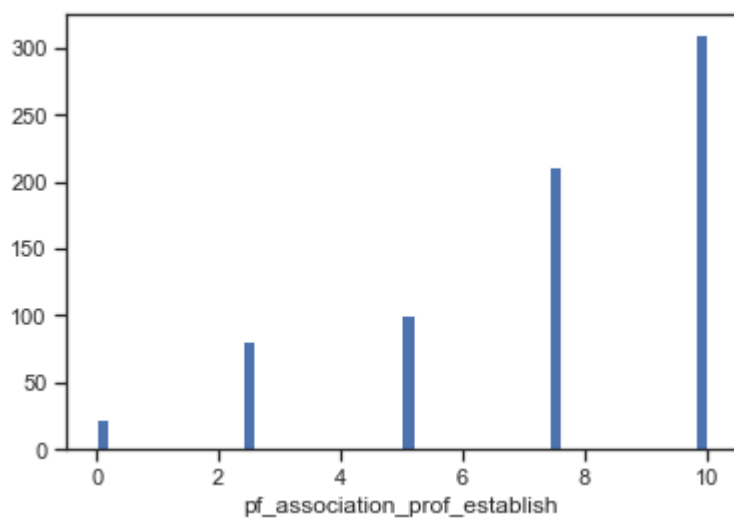


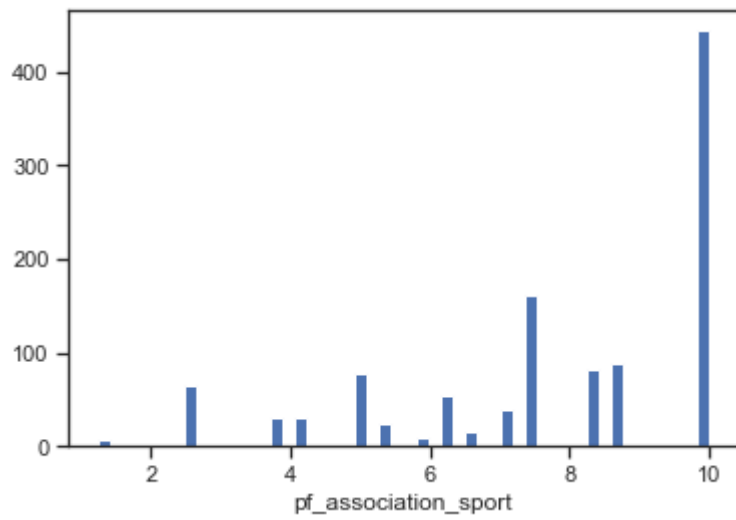
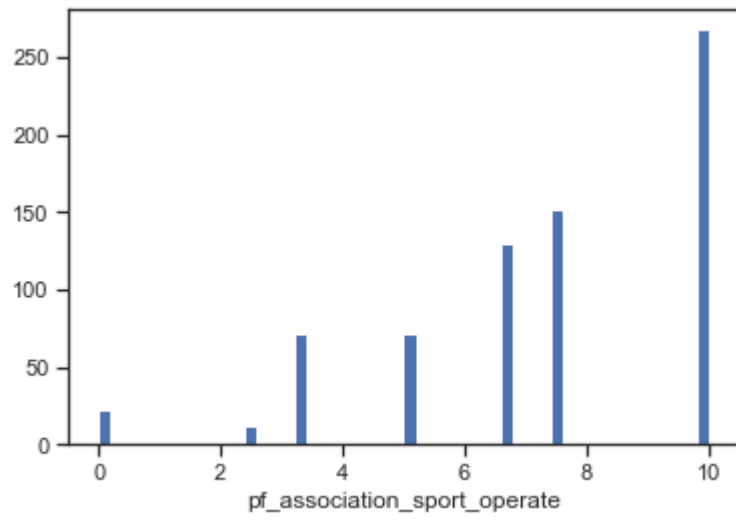
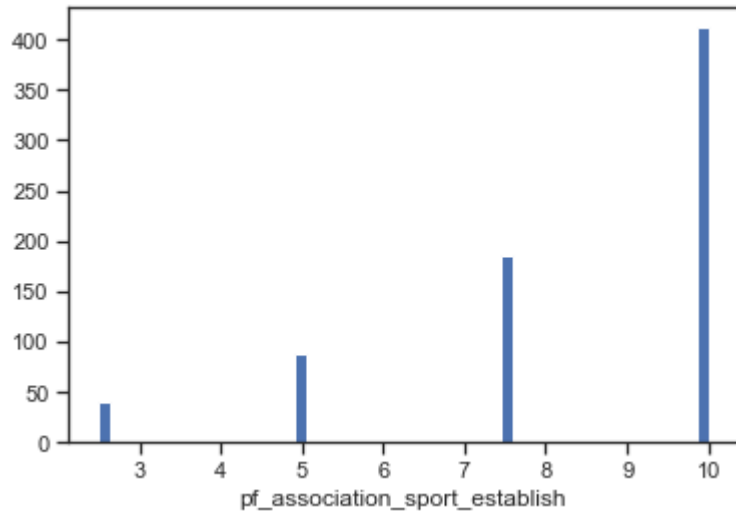


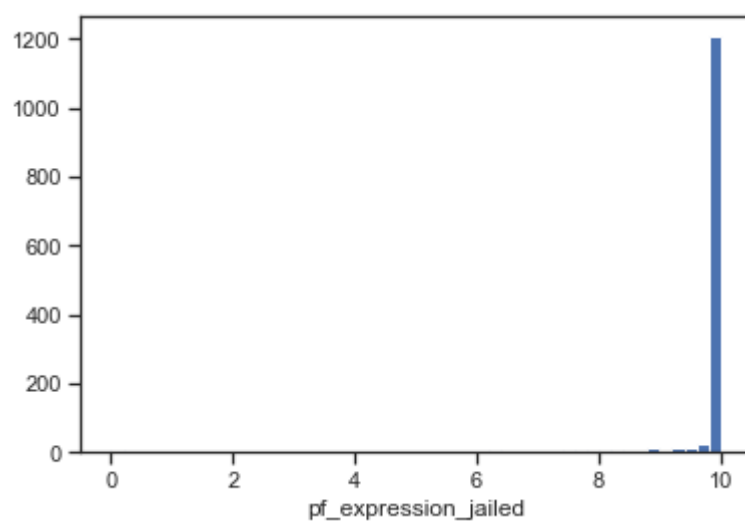
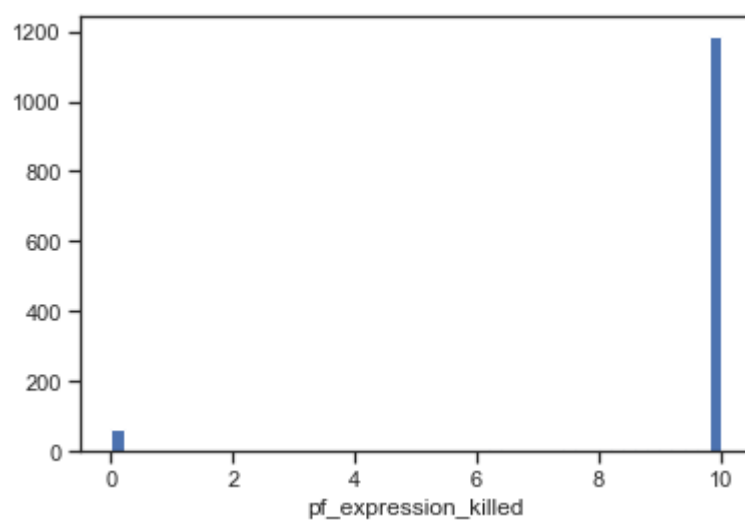
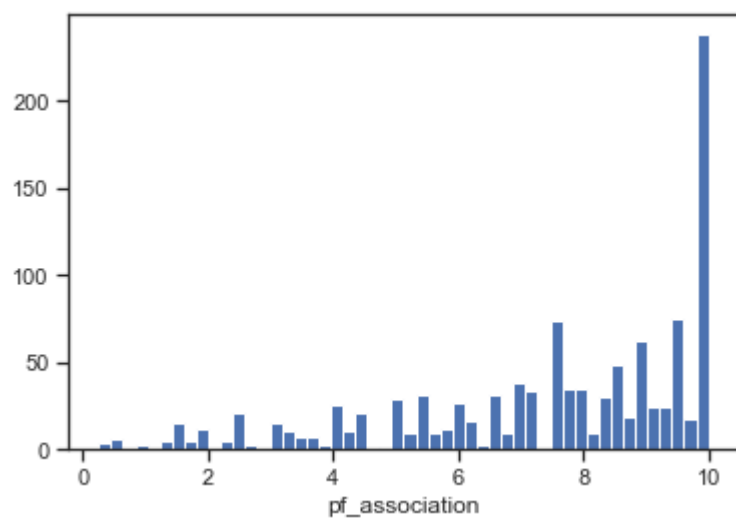


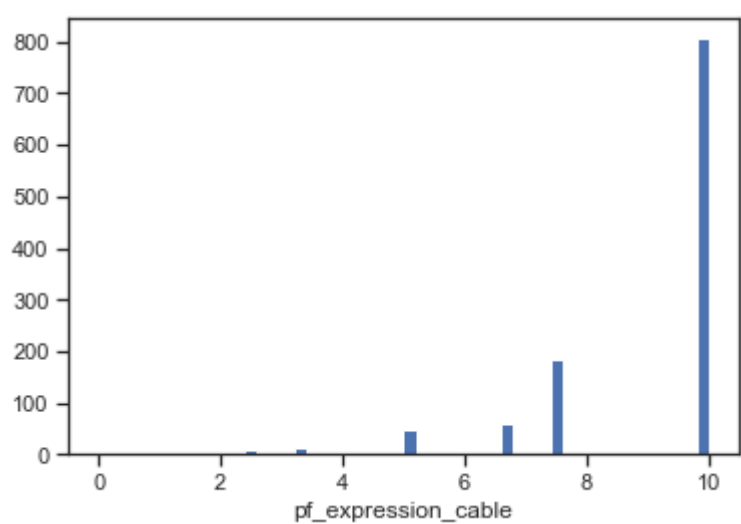
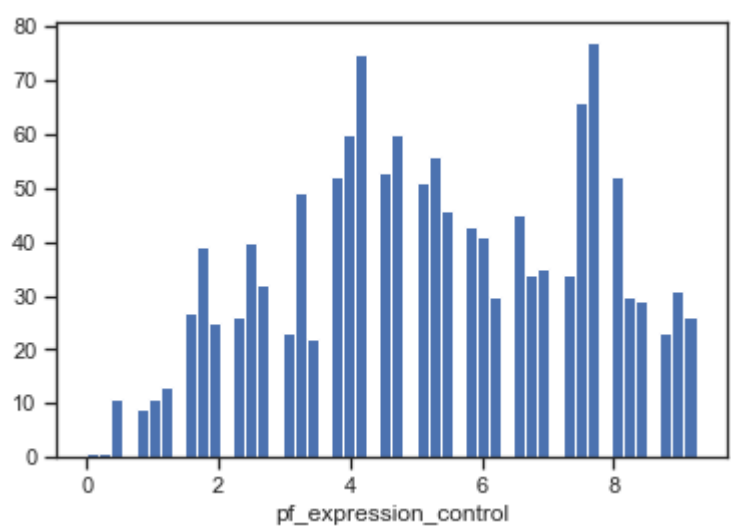
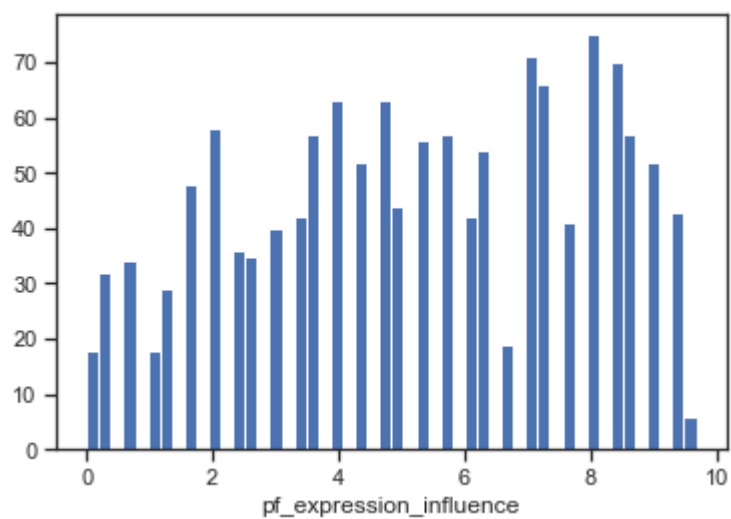


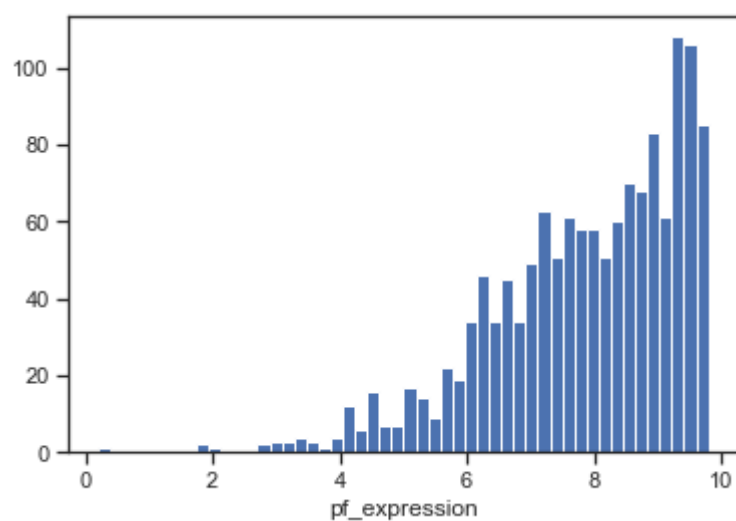
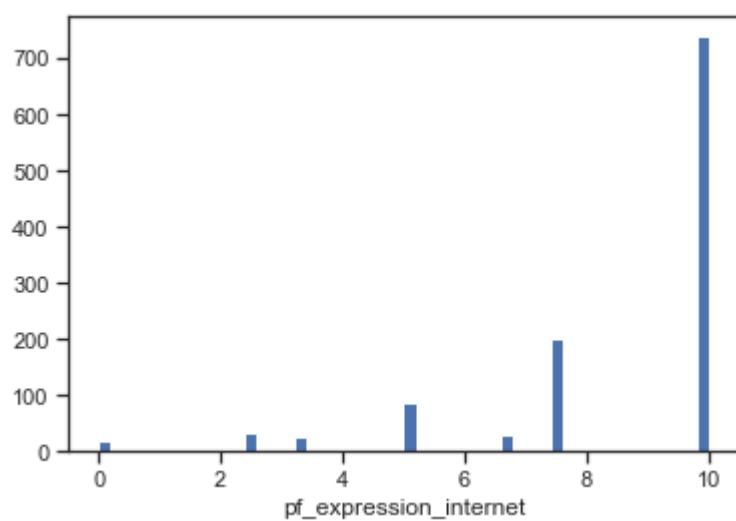
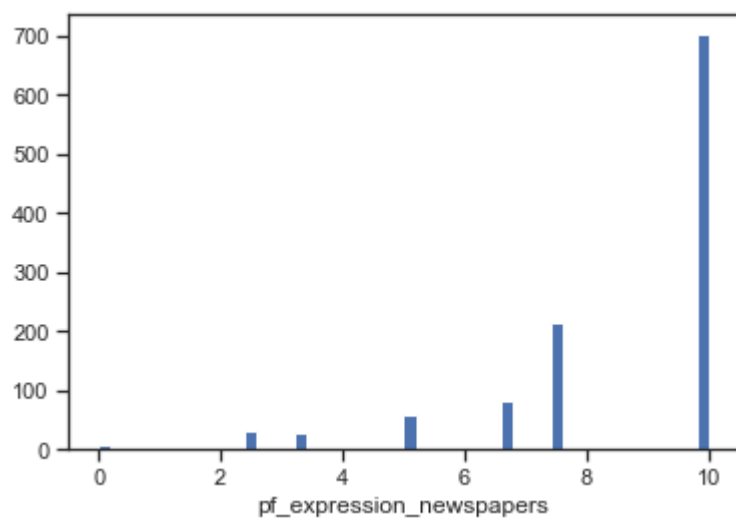


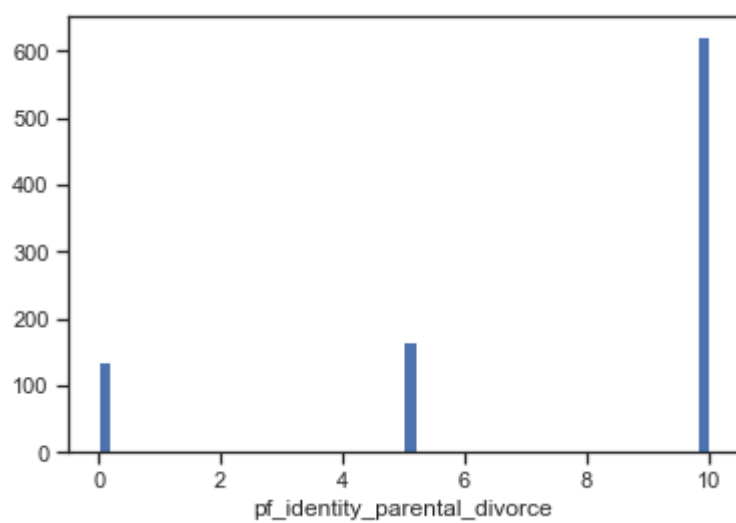
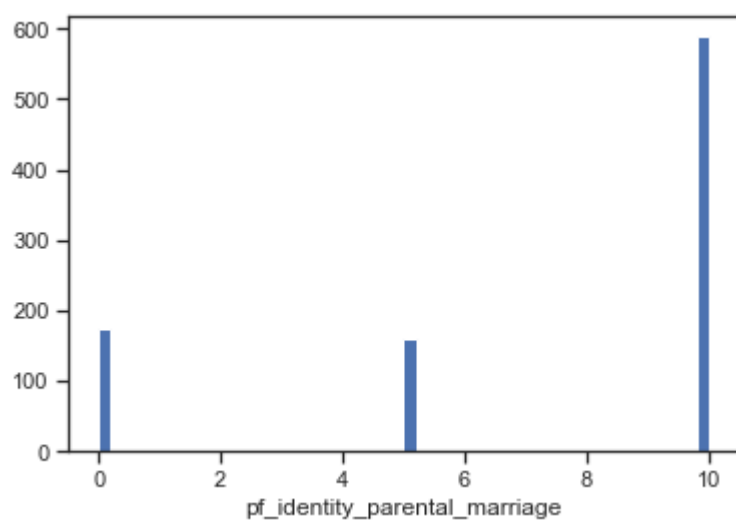
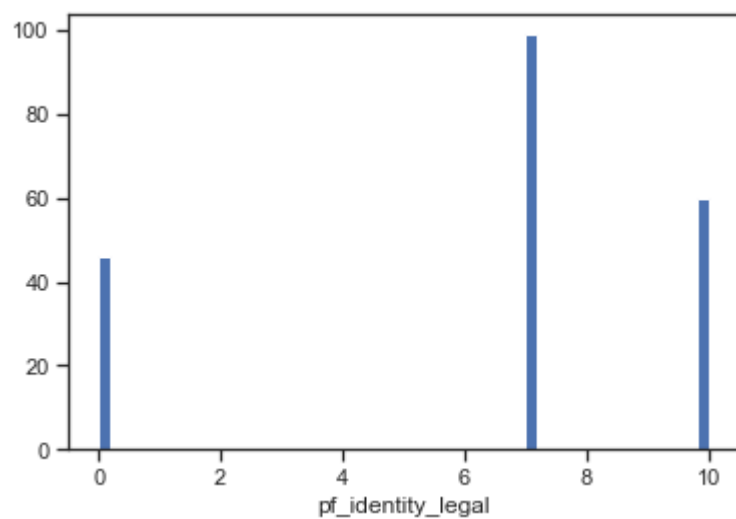




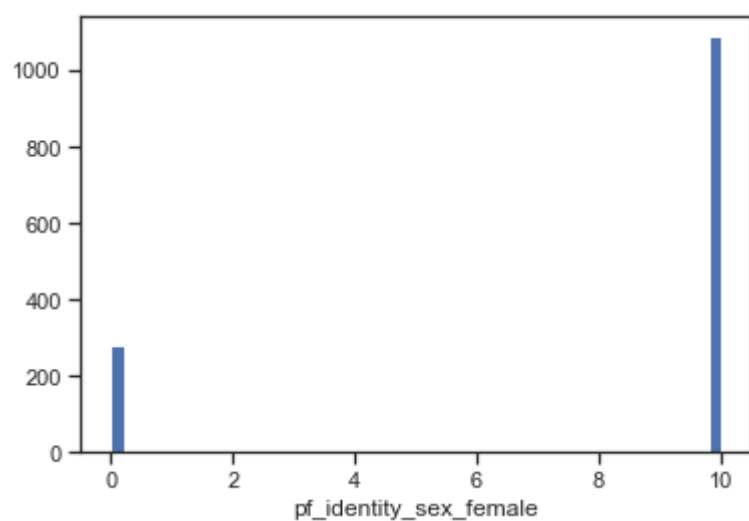
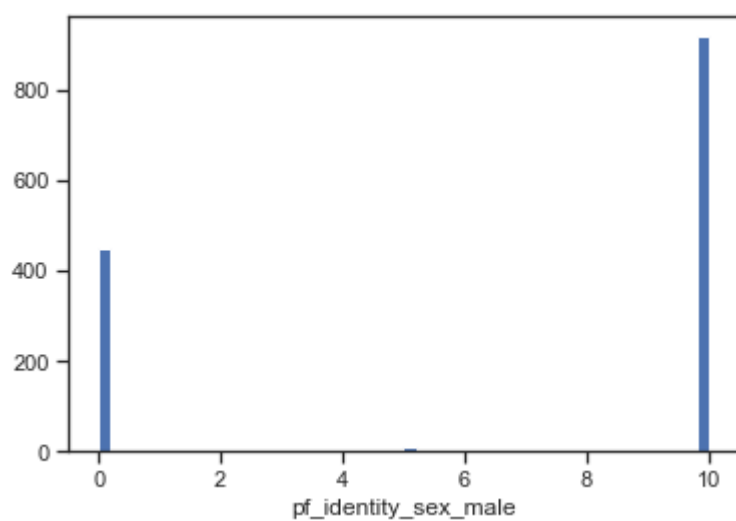
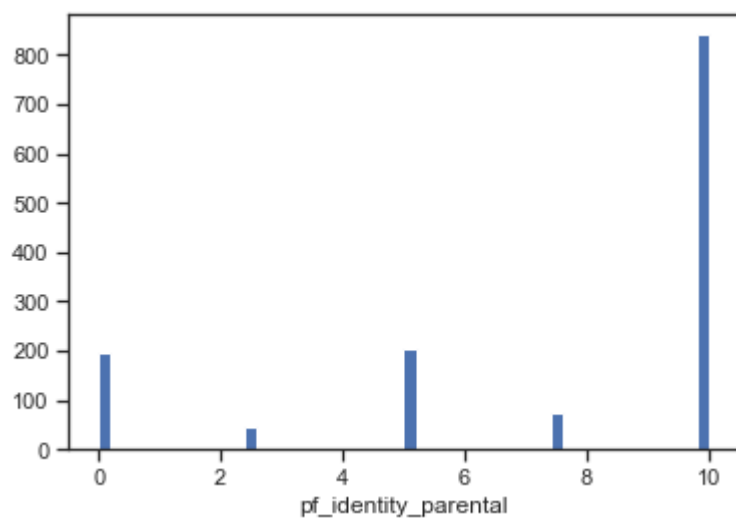


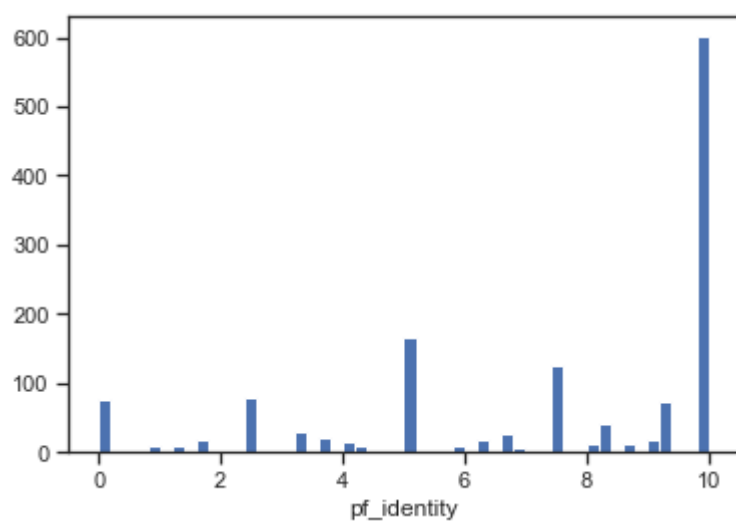
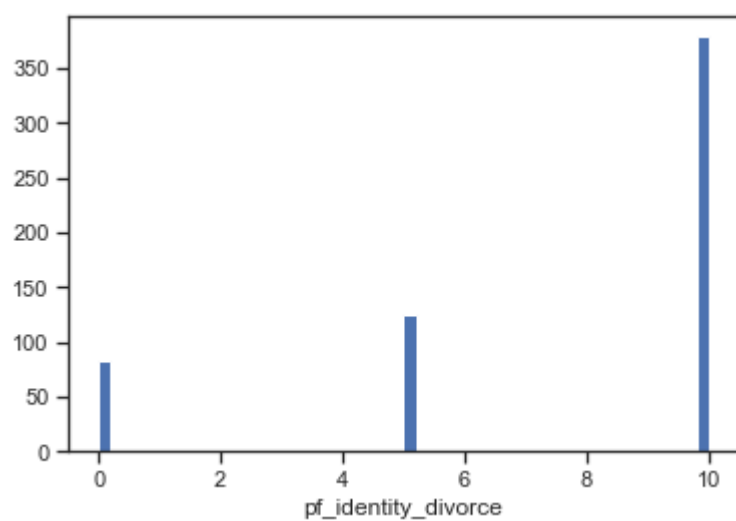
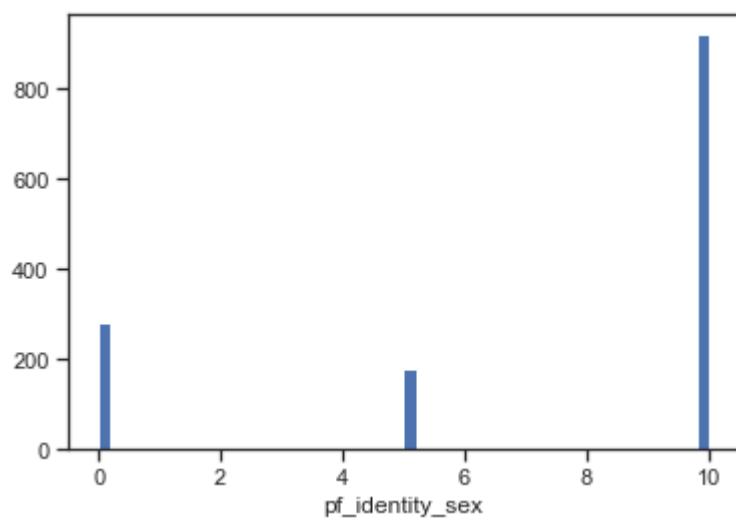


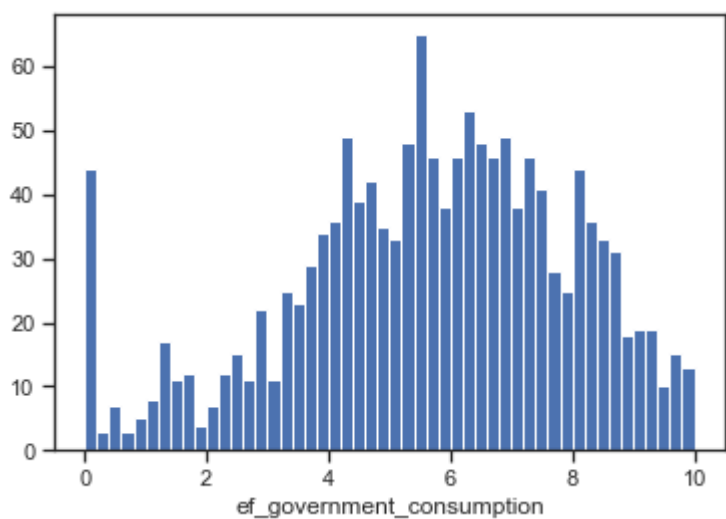
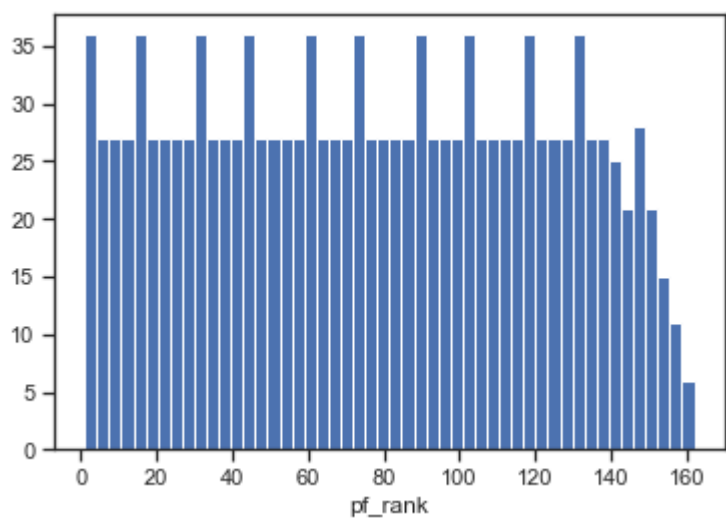
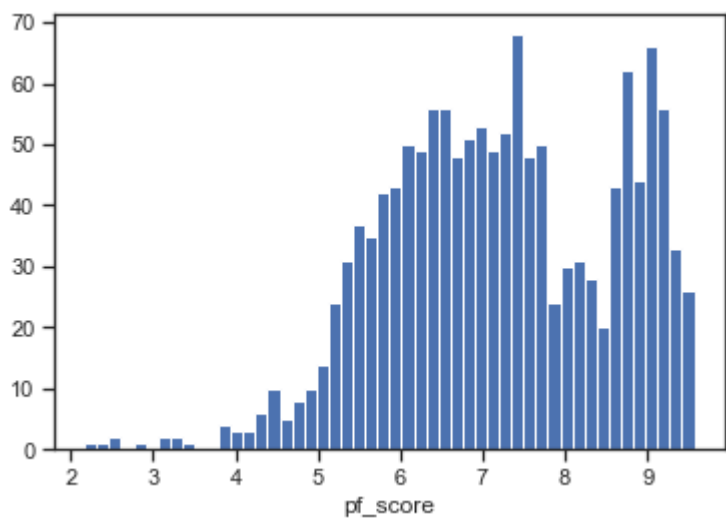


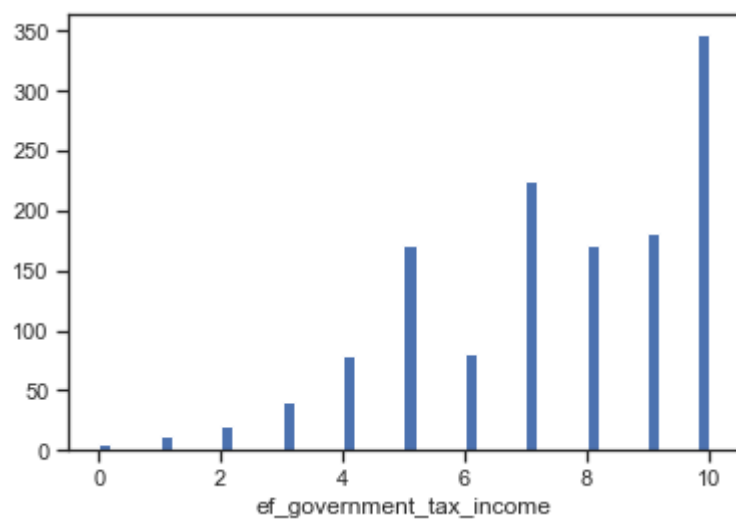
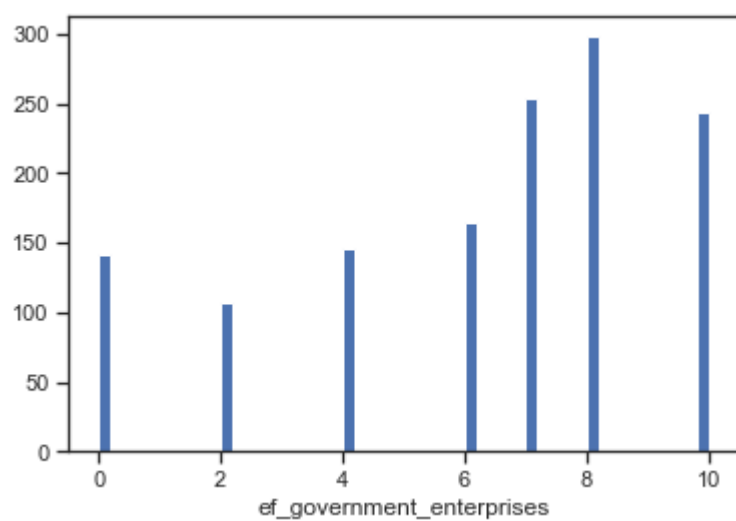
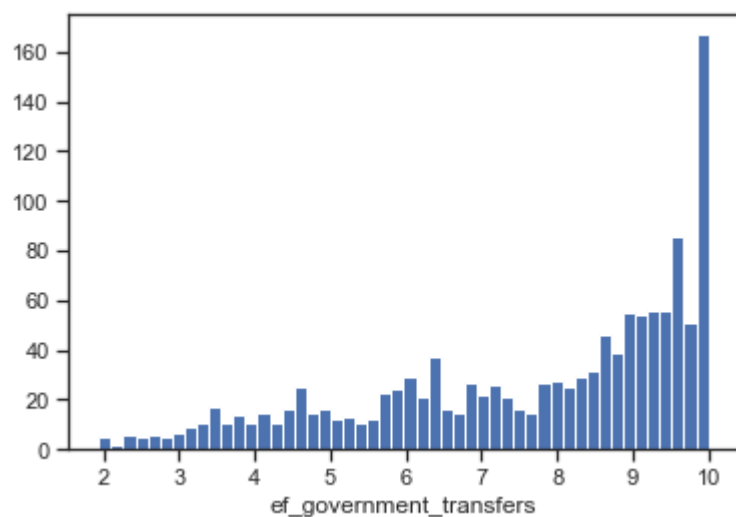


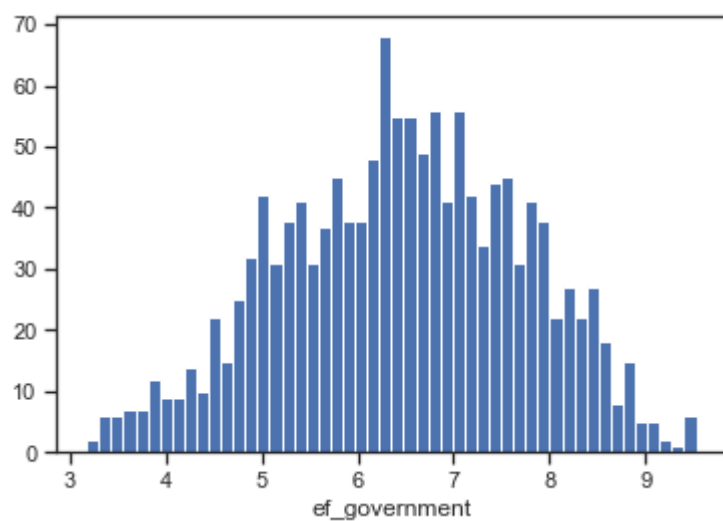
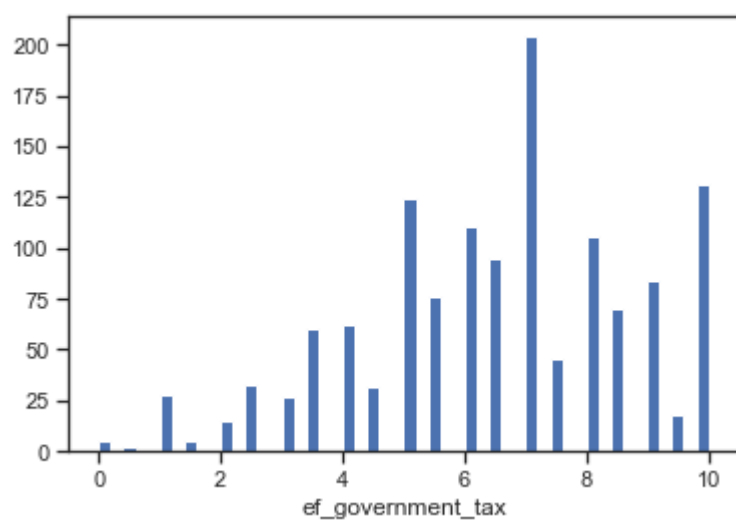
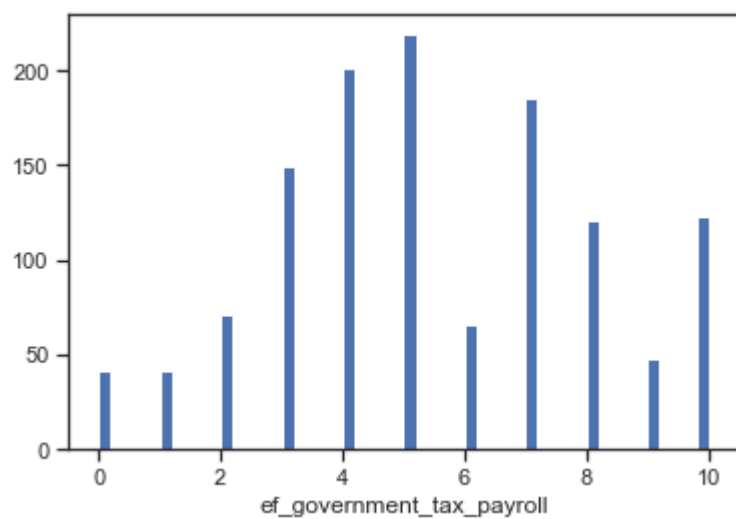


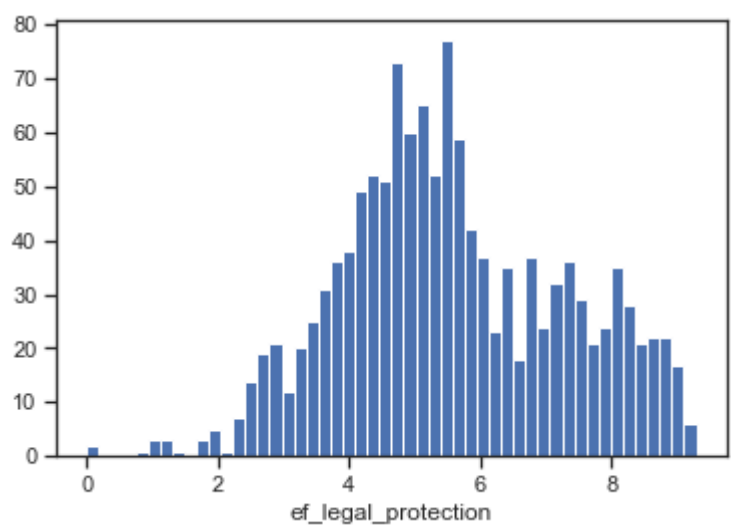
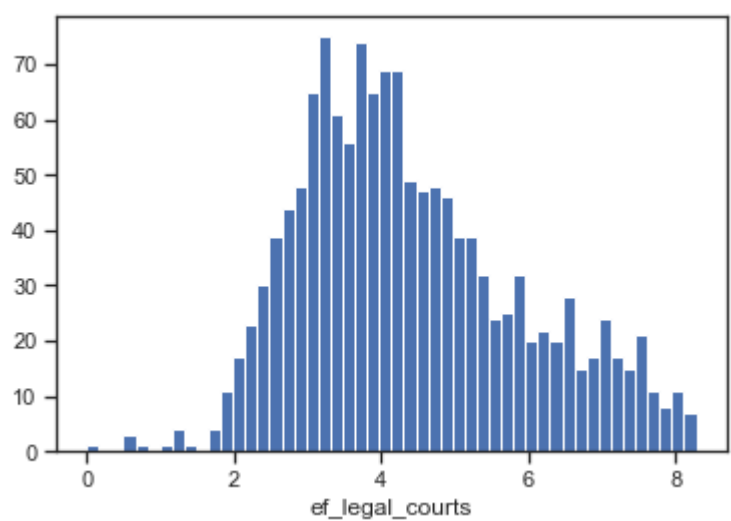
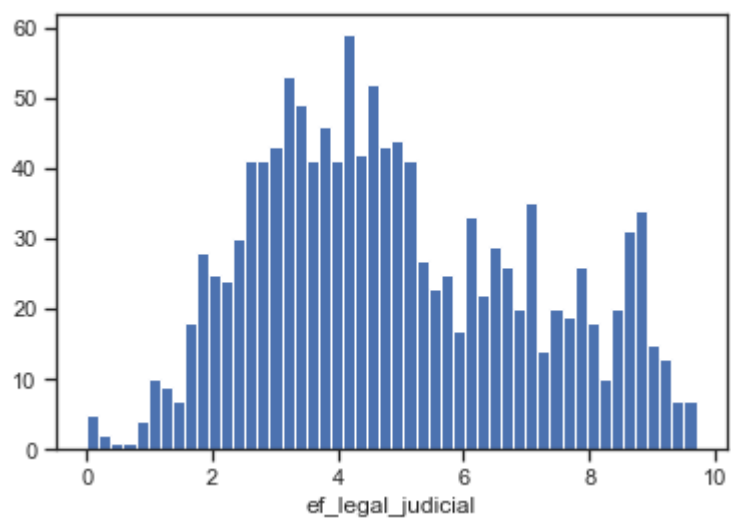


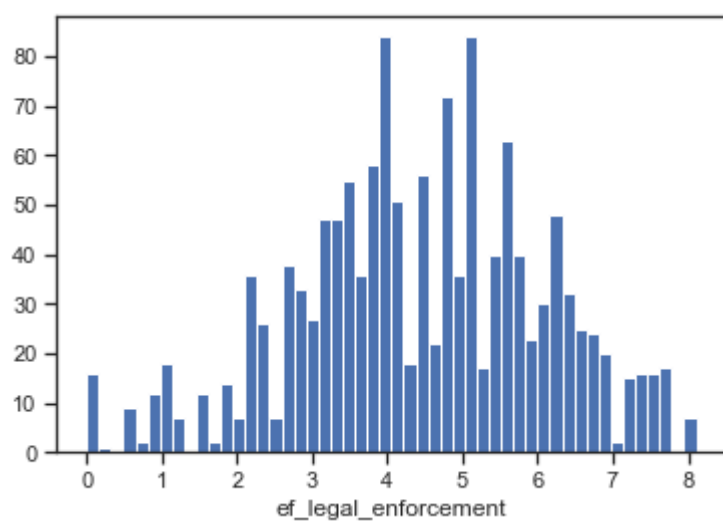
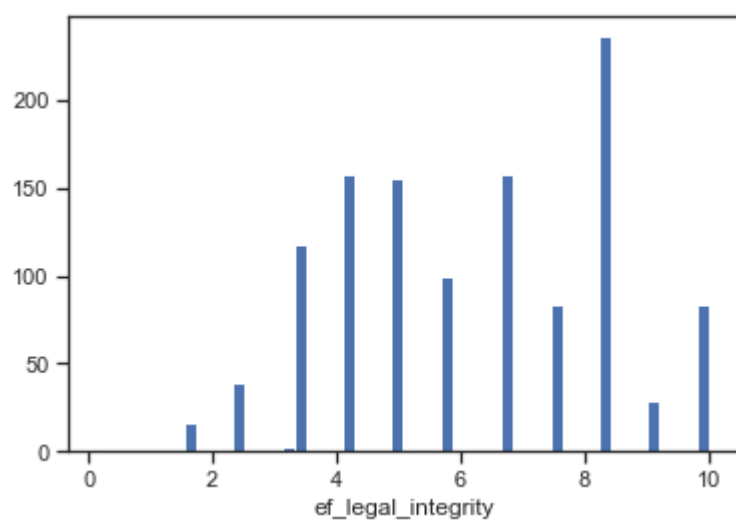
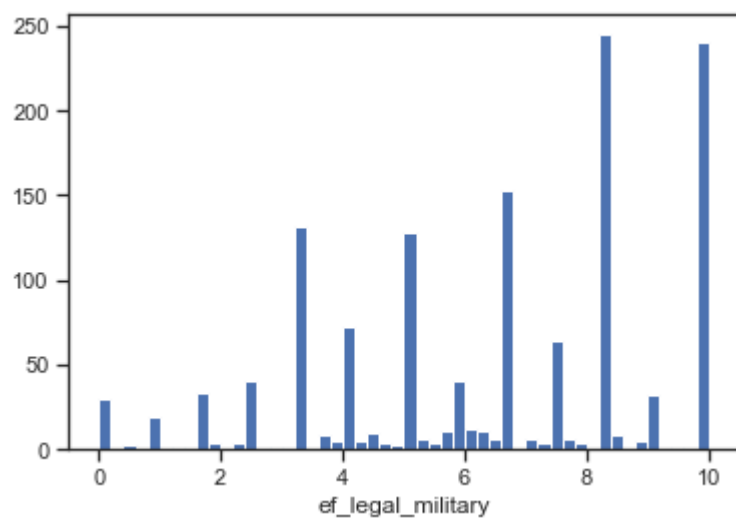


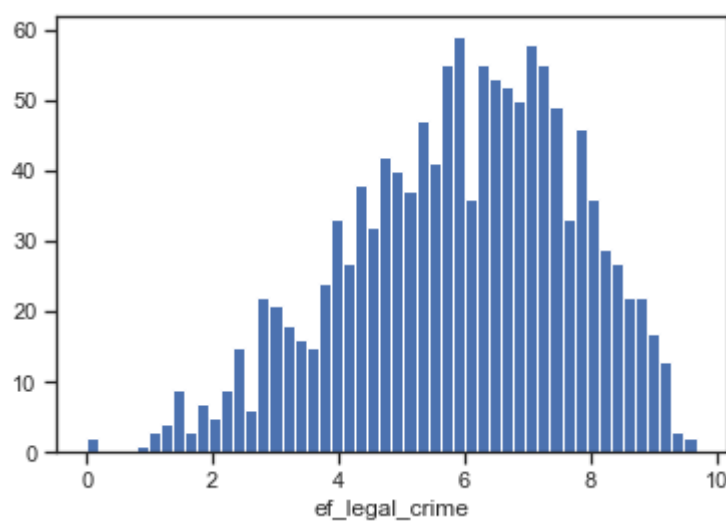
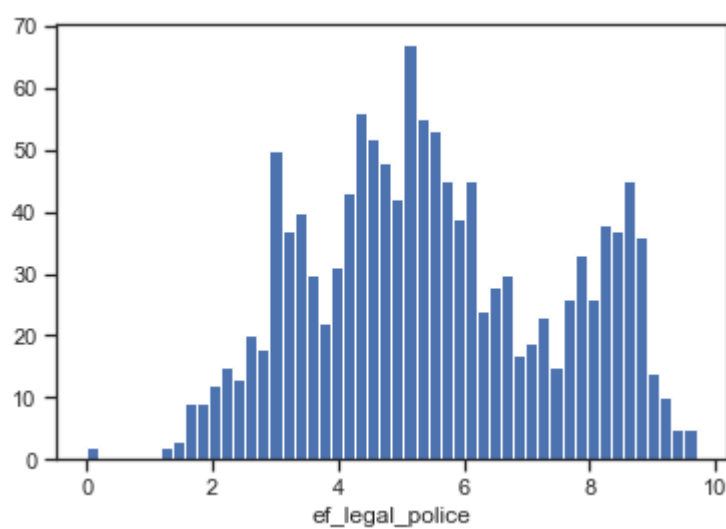
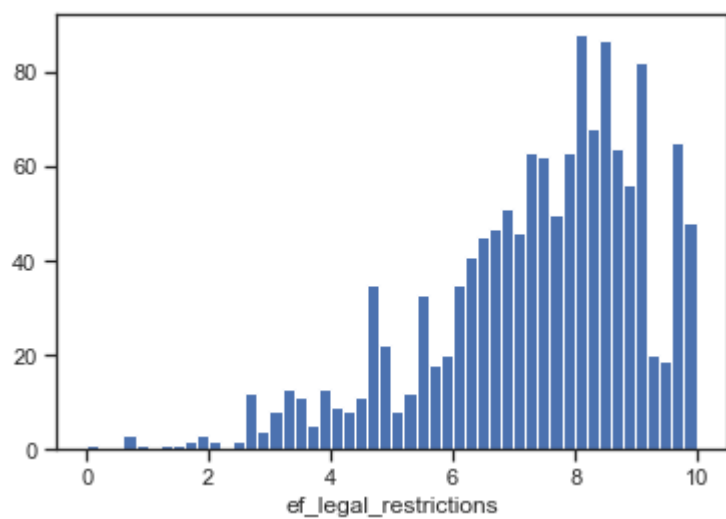




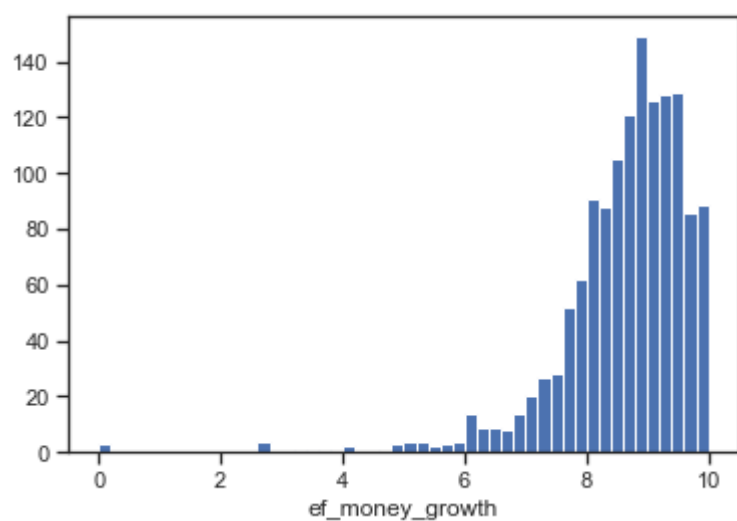
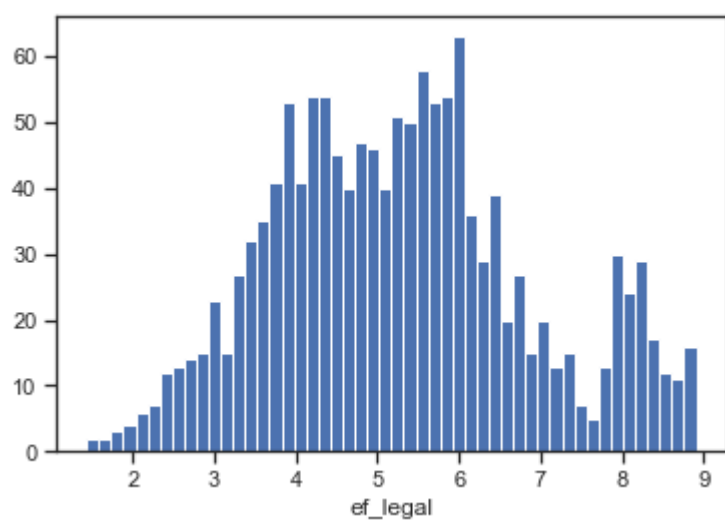
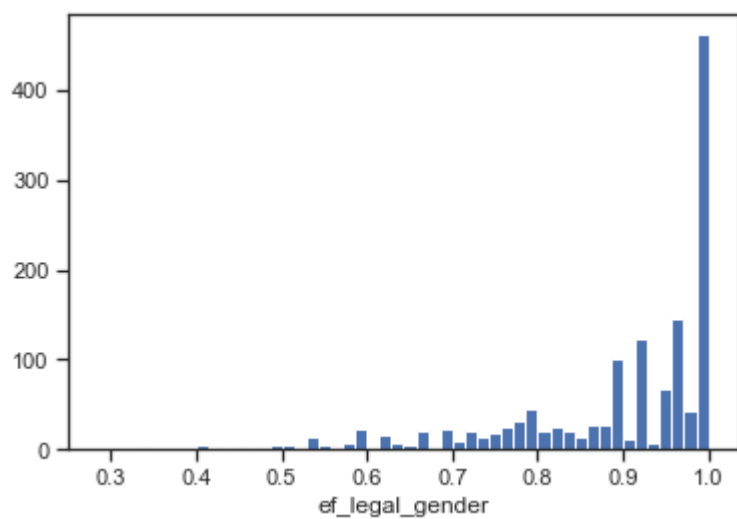


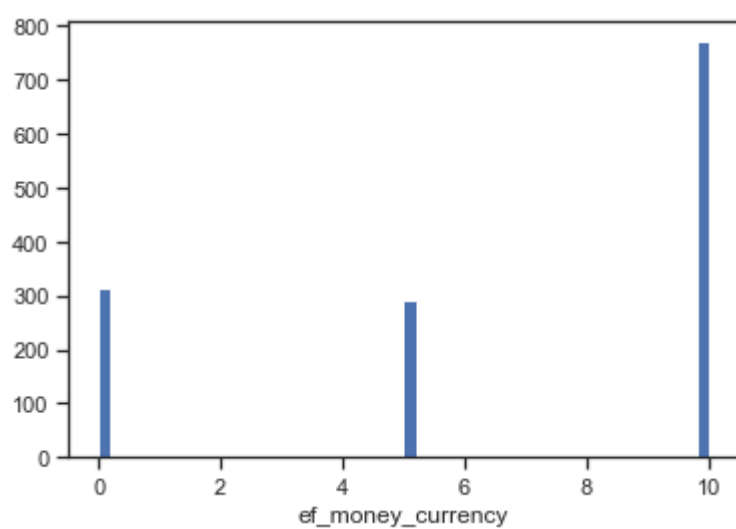
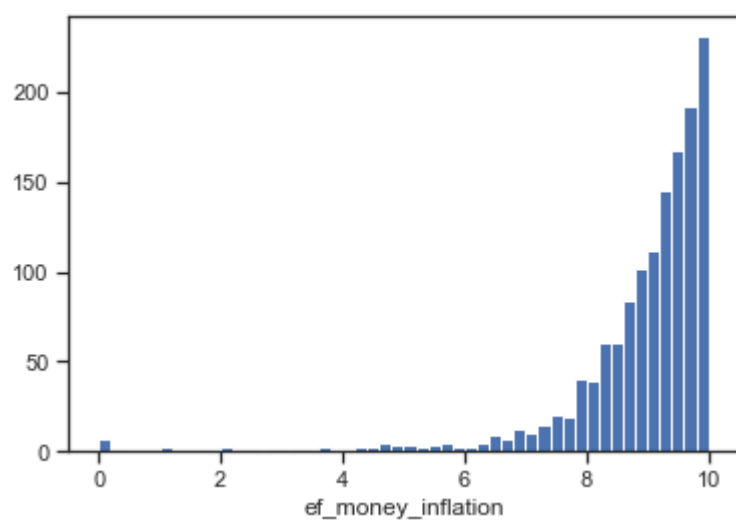
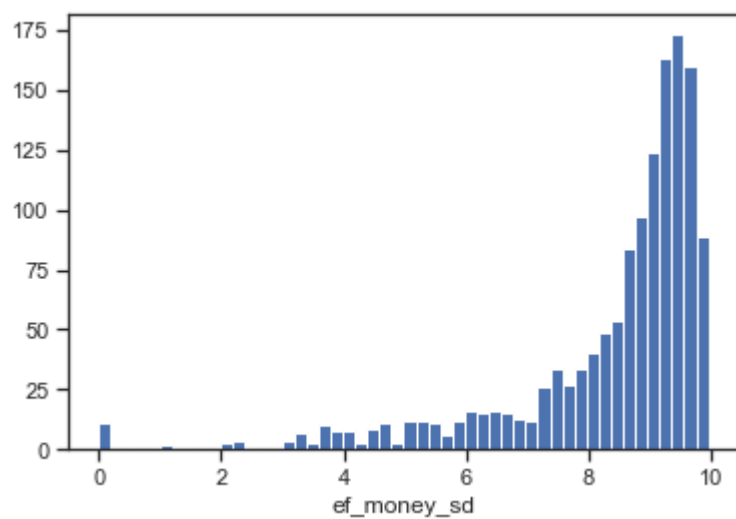


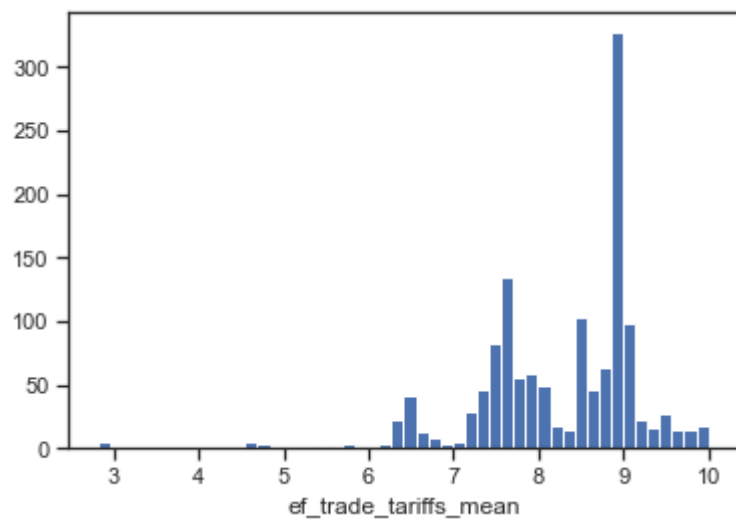
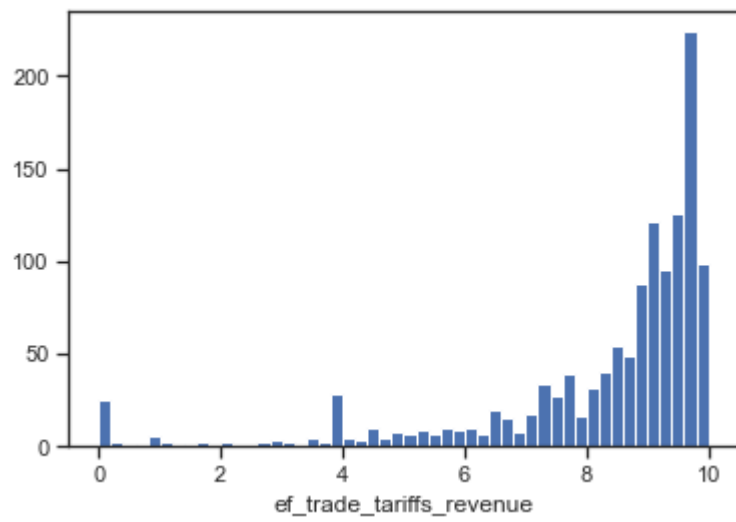
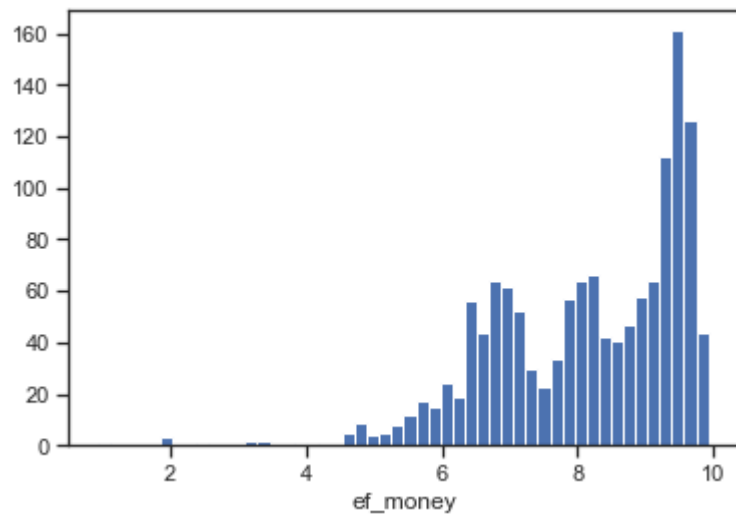


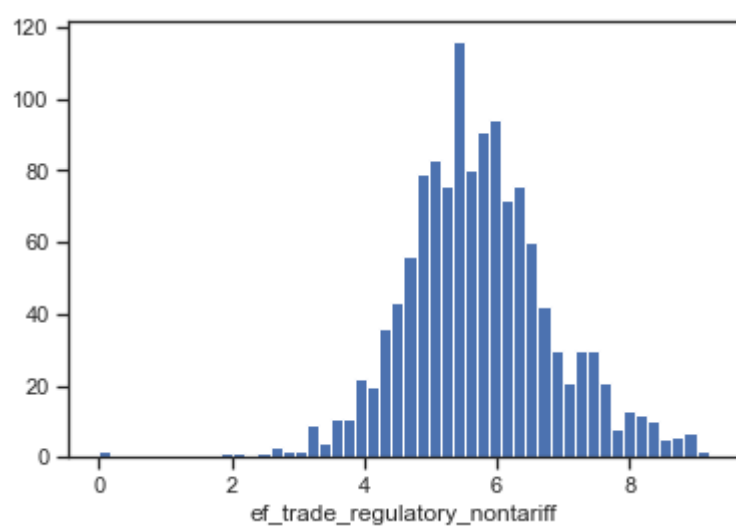
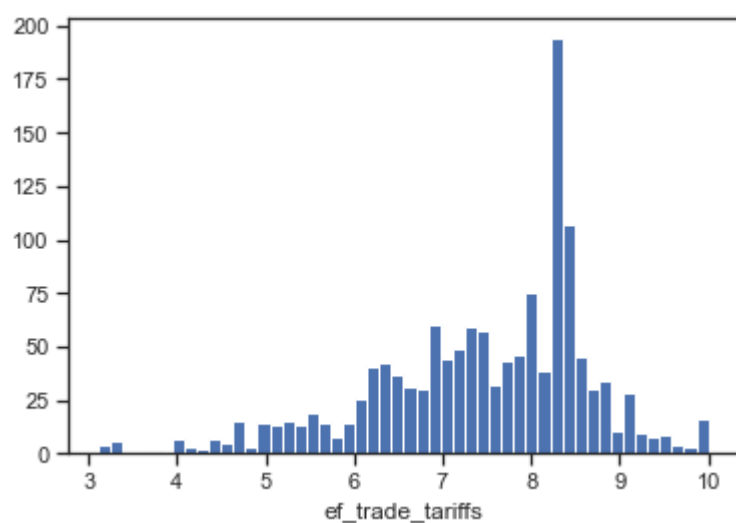
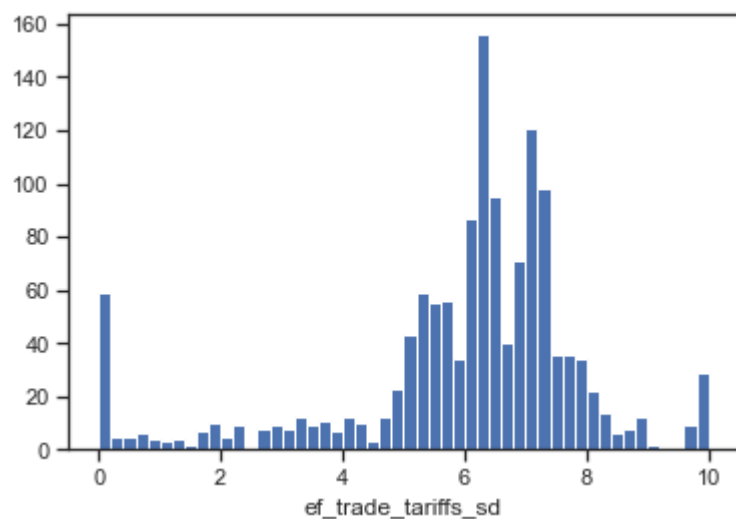


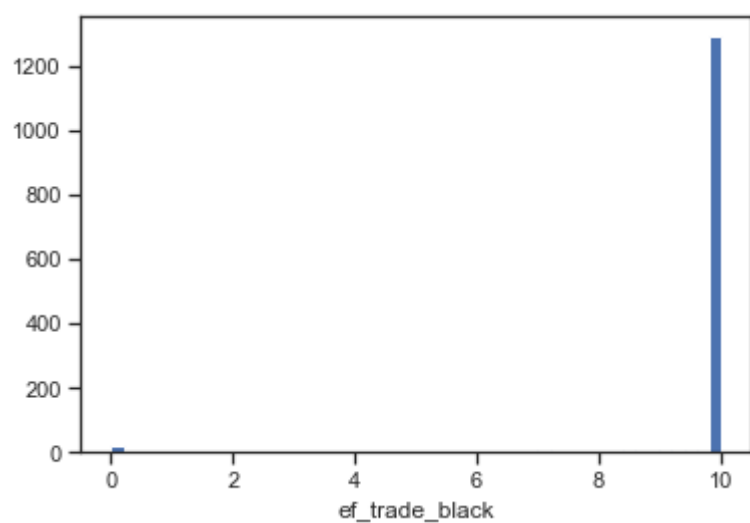
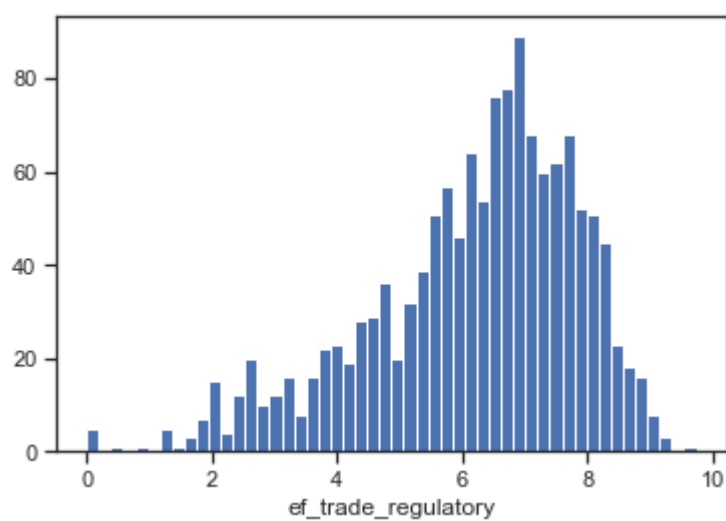
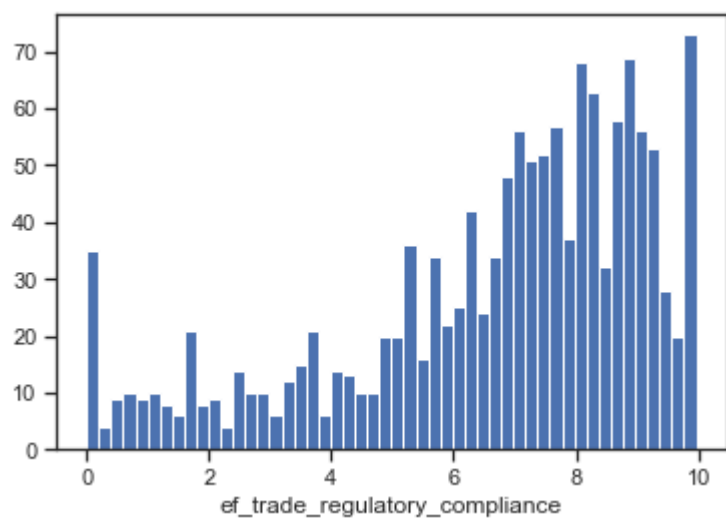


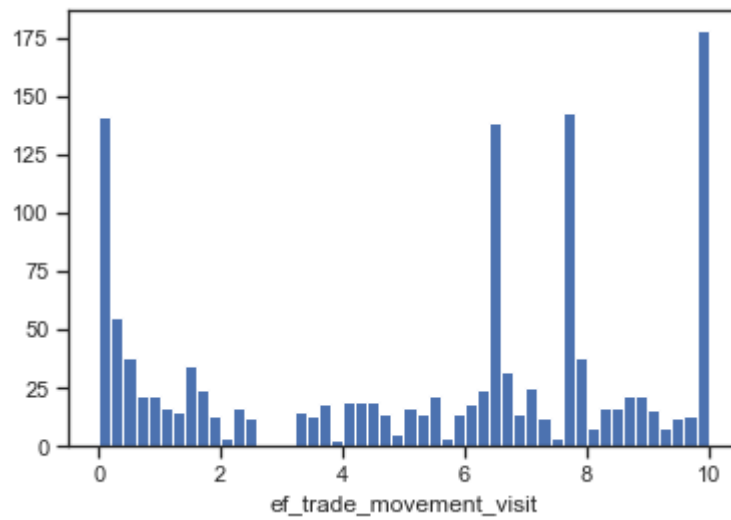
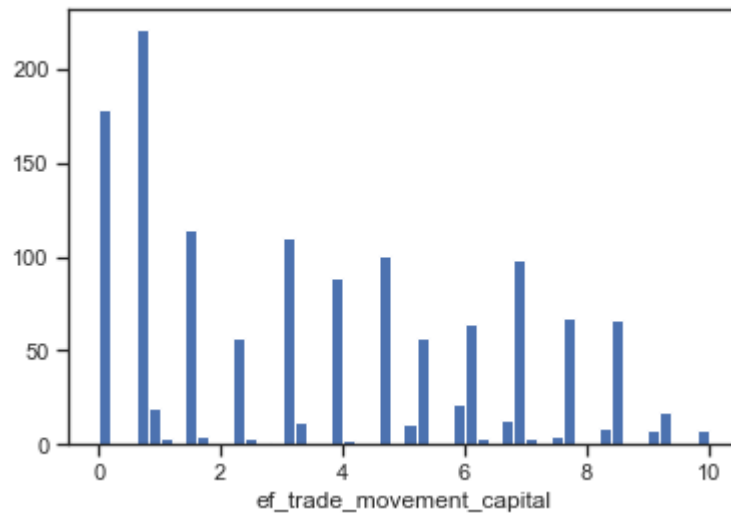
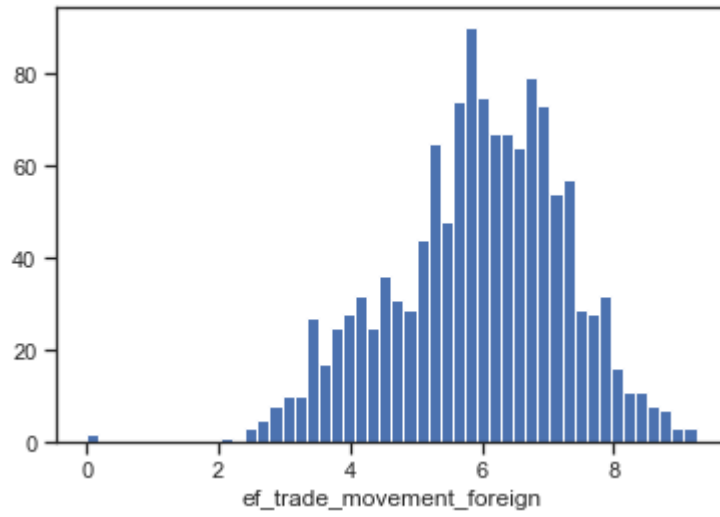


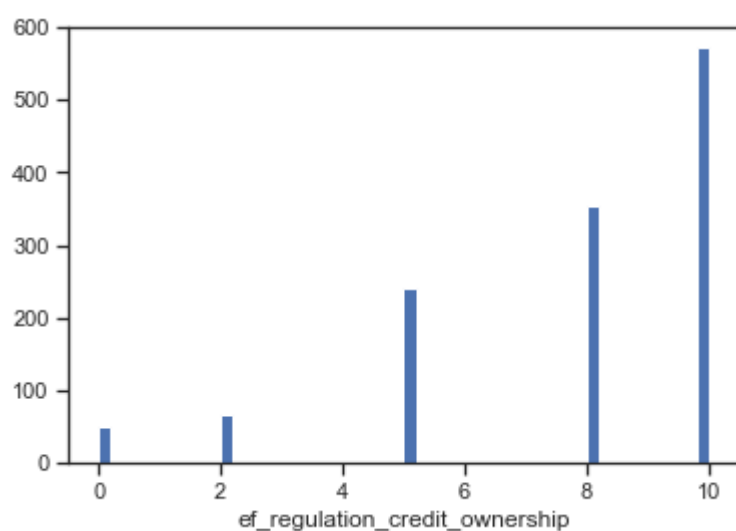
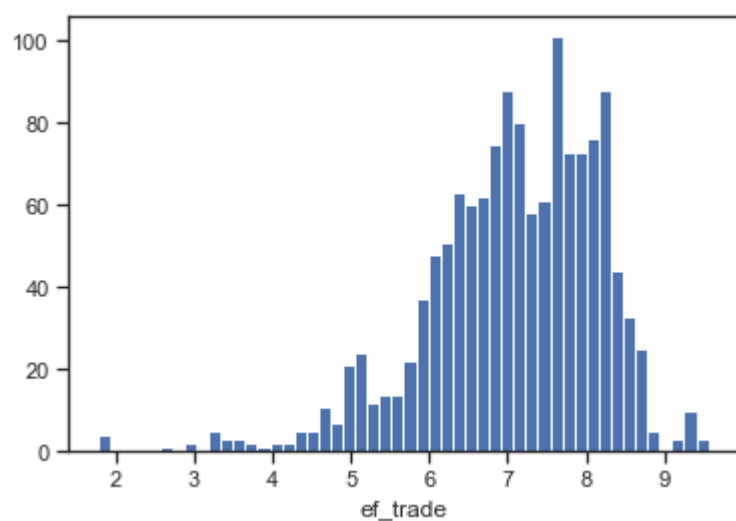
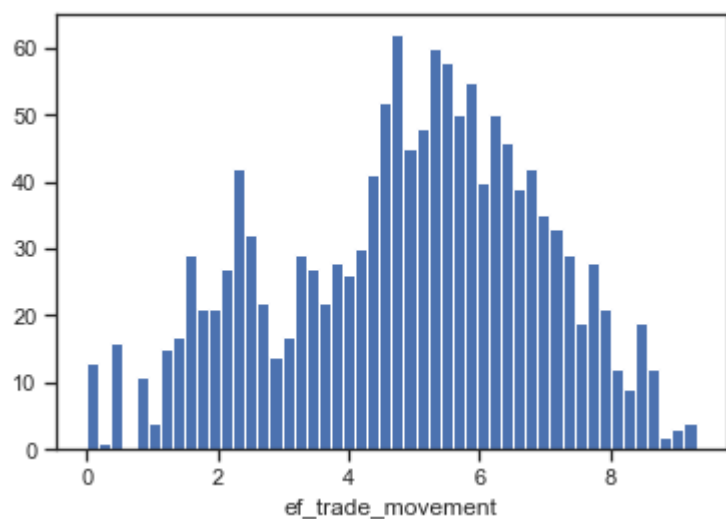


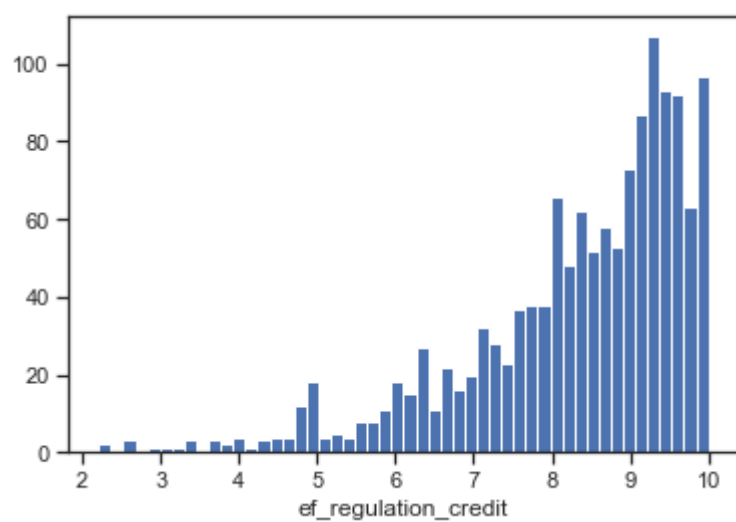
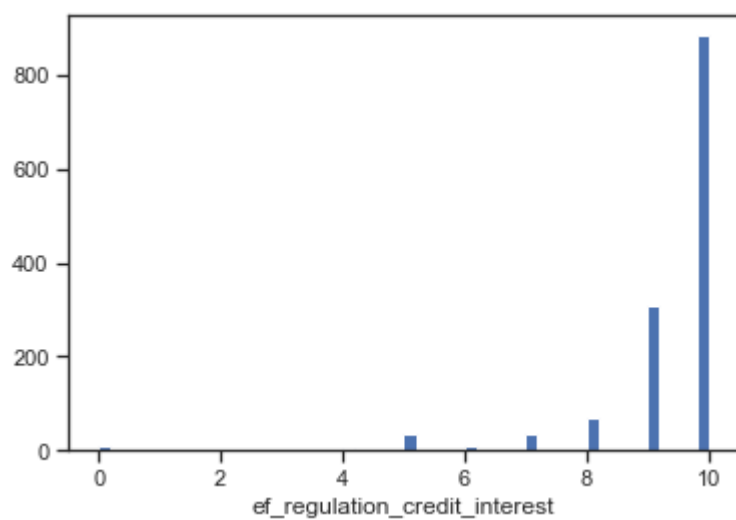
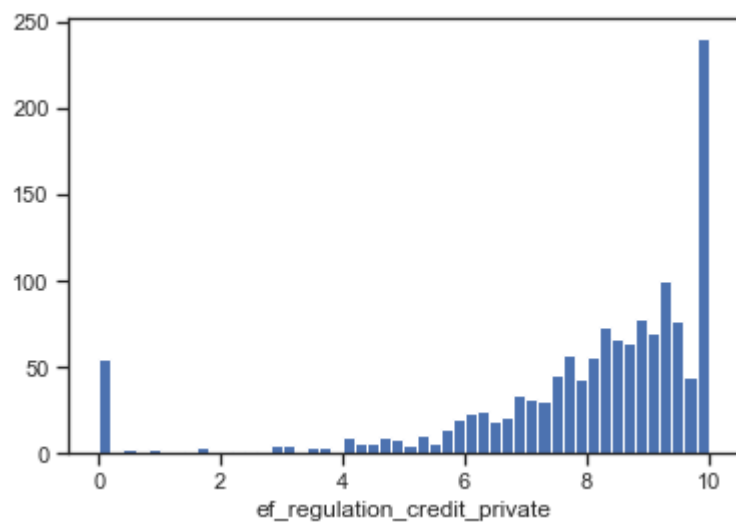




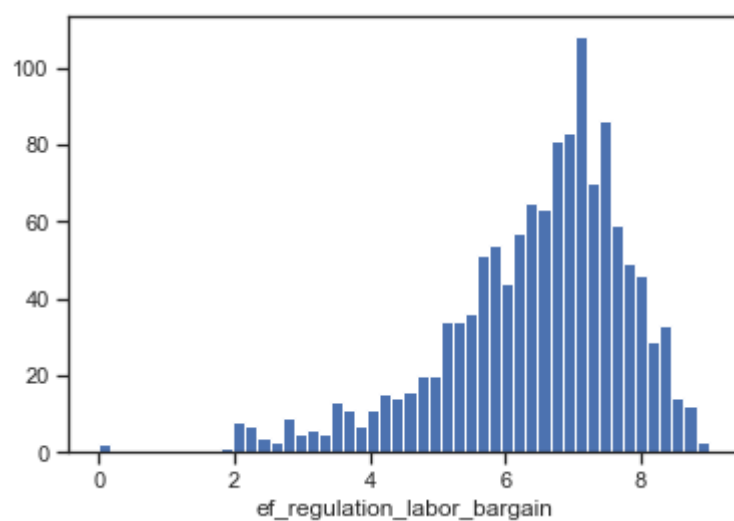
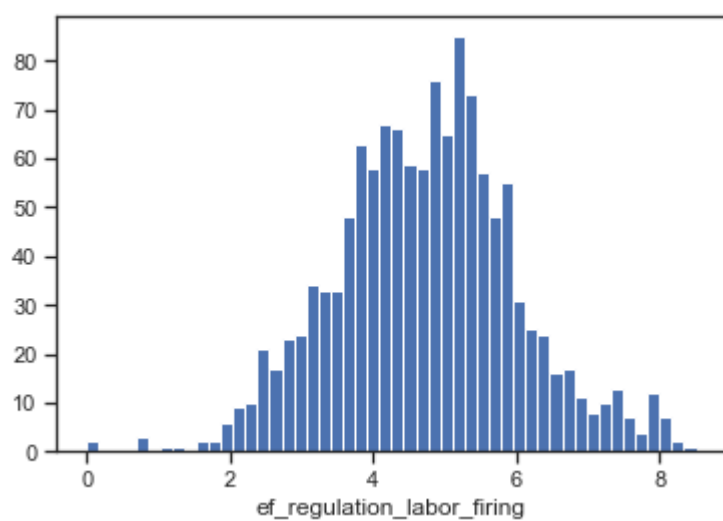
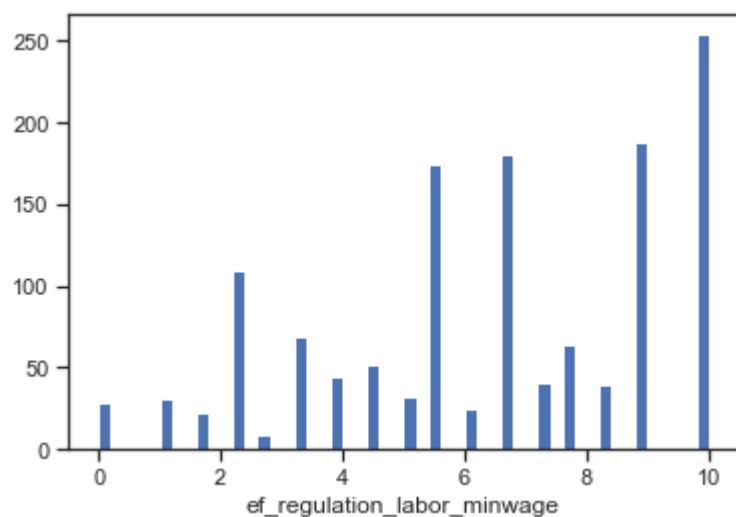


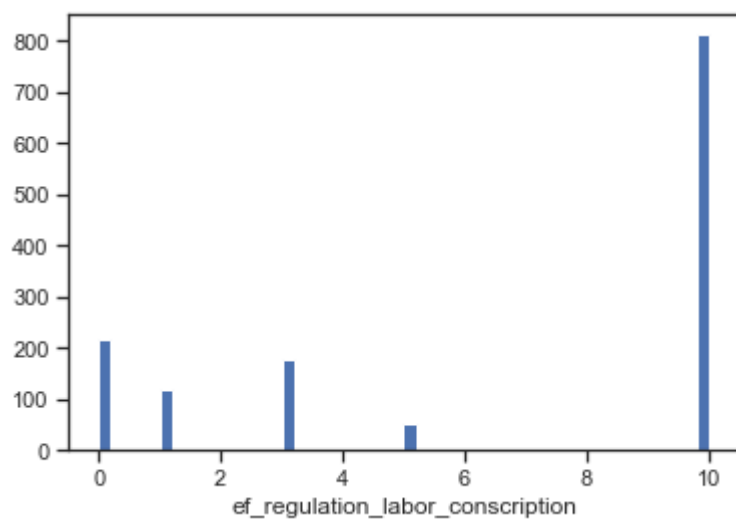
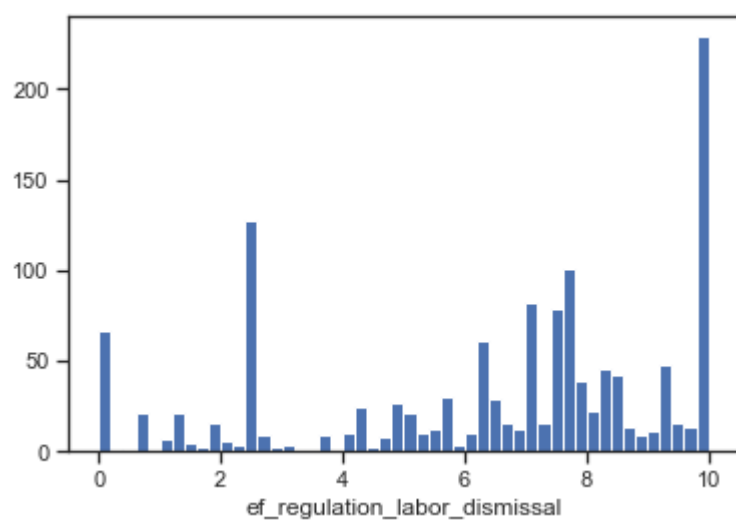
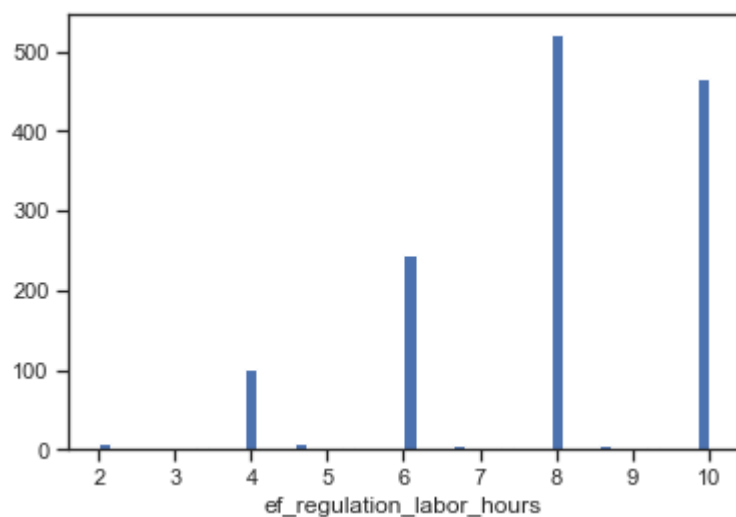


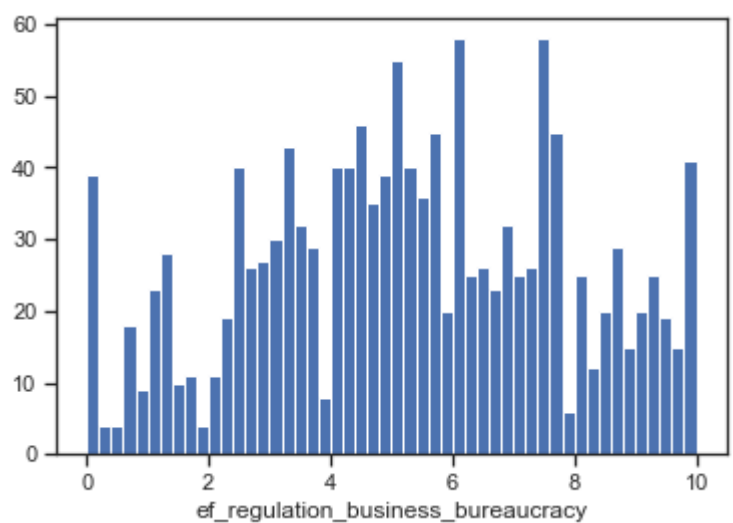
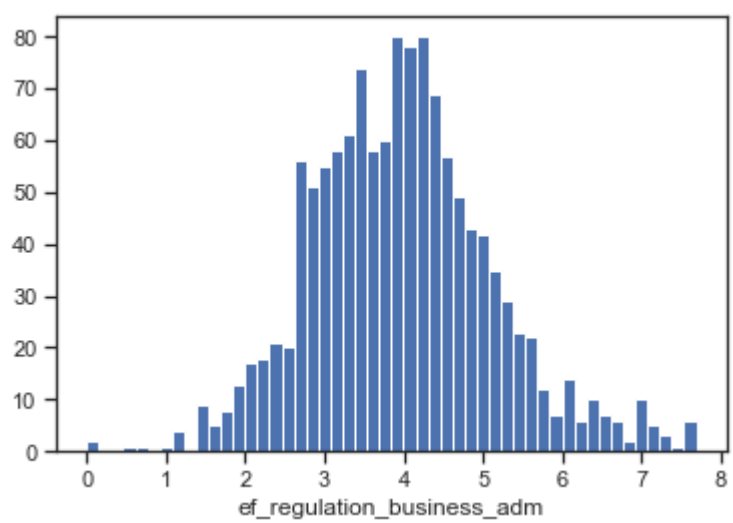
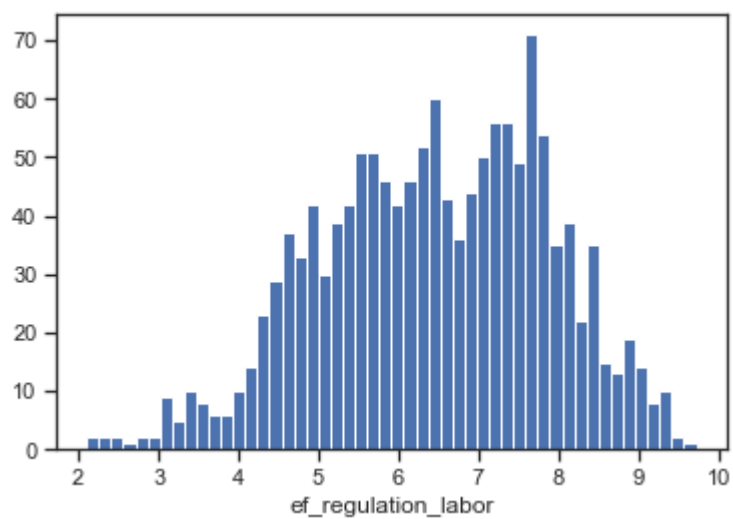


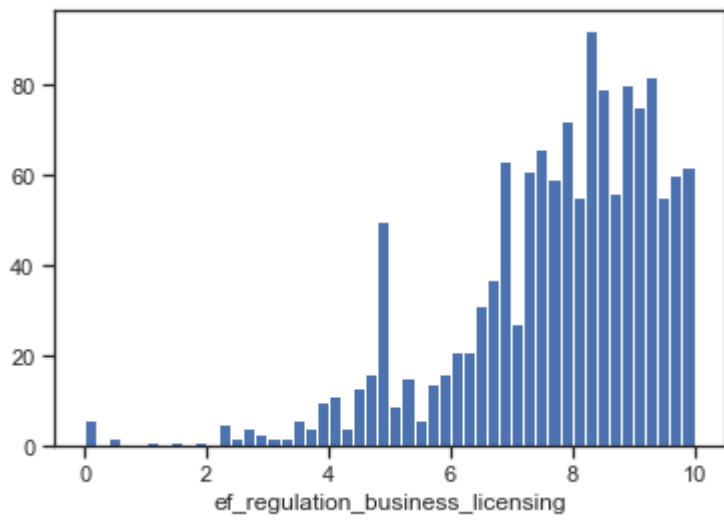
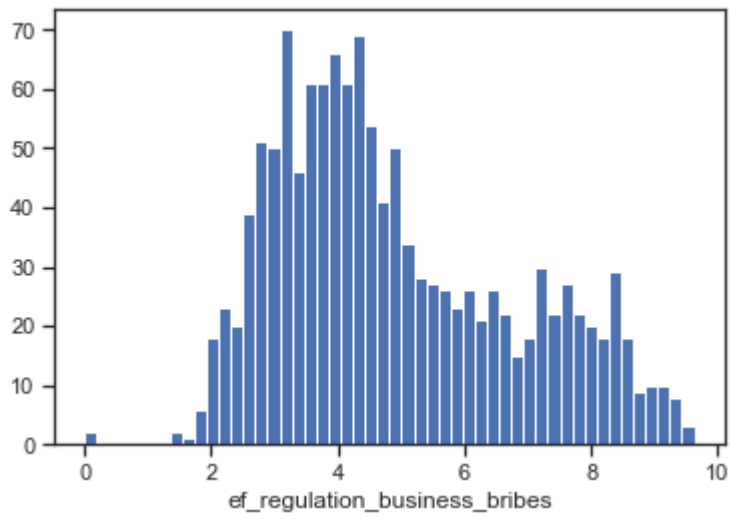
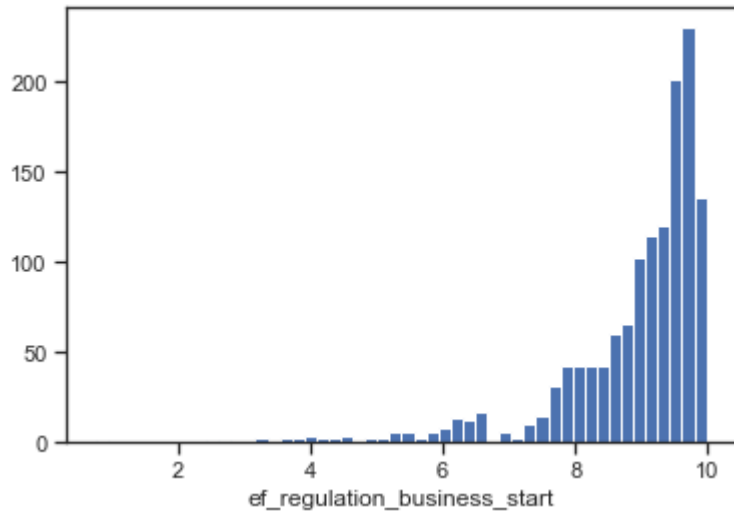


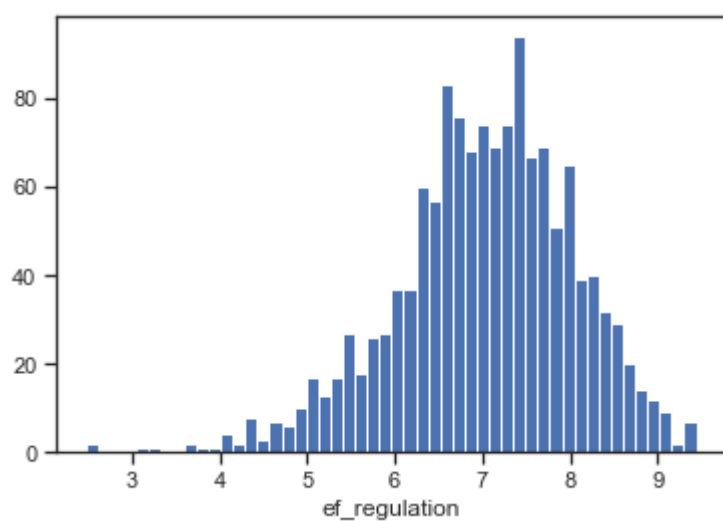
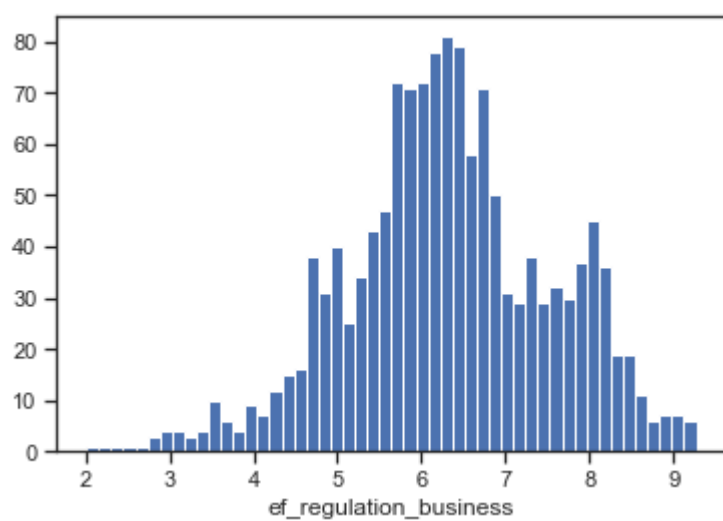
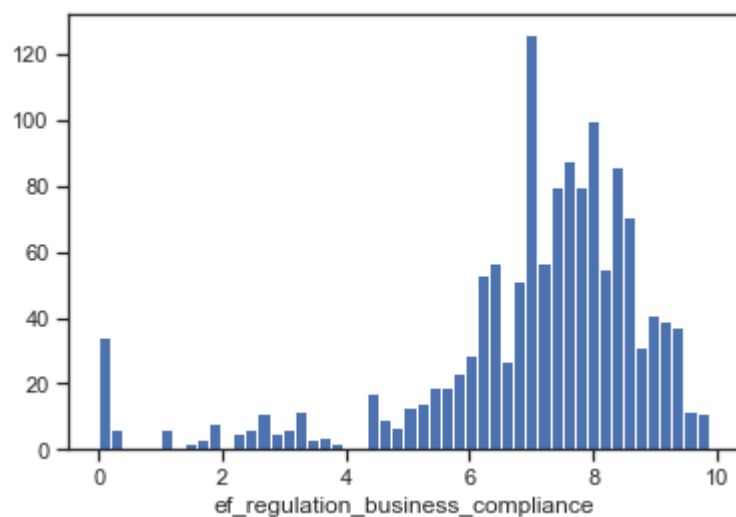


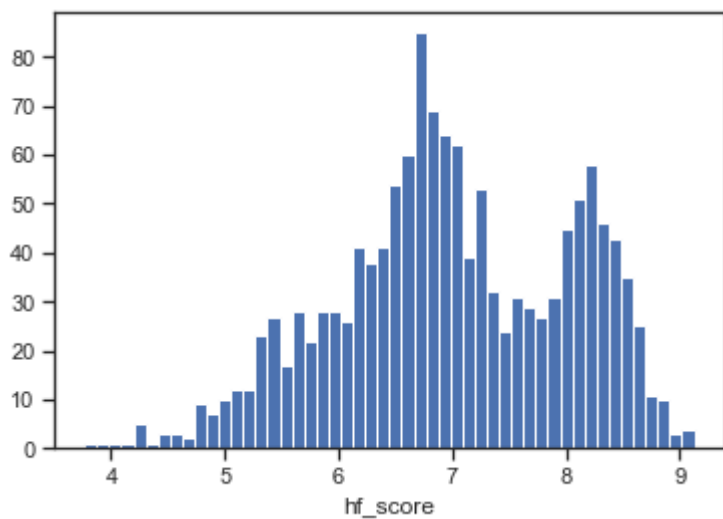
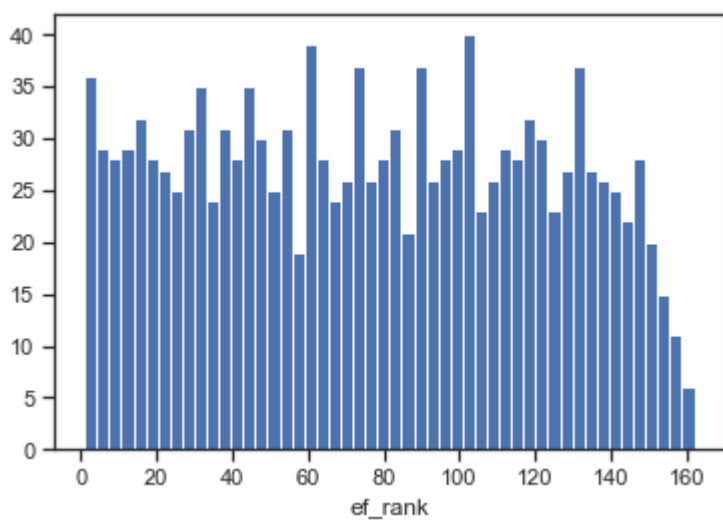
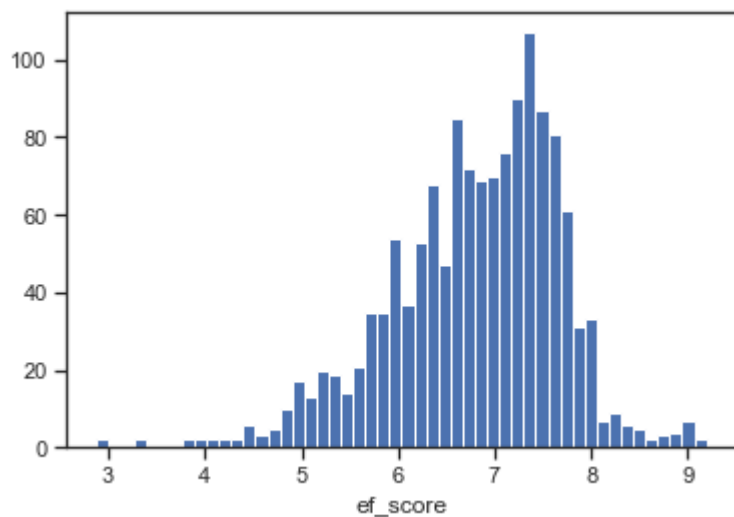


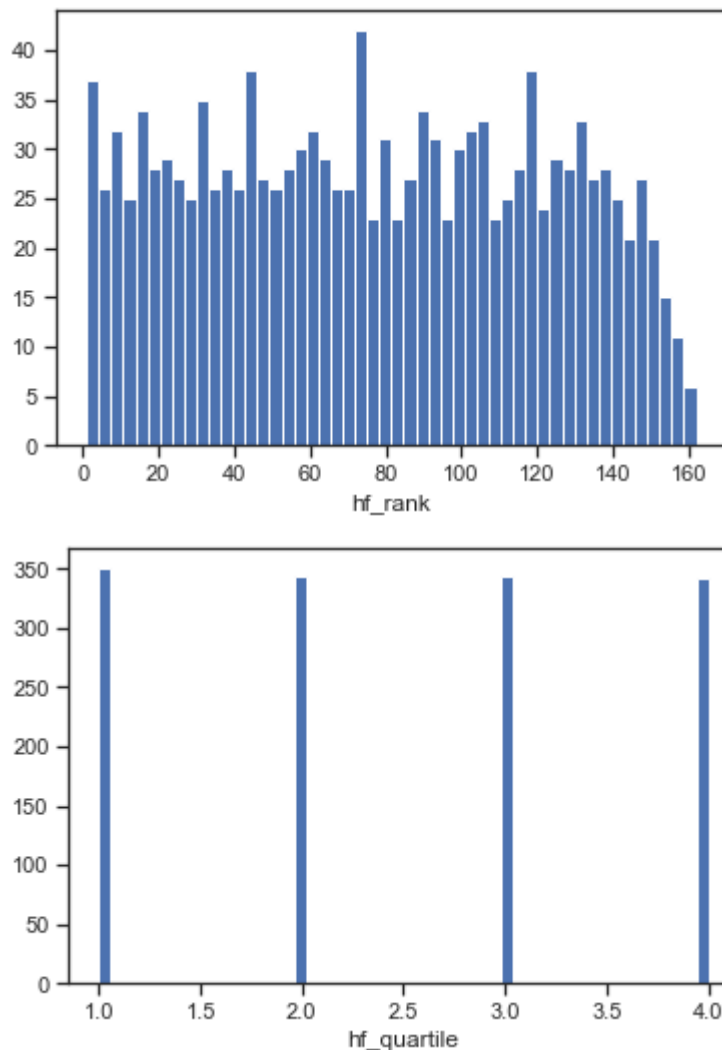












## Обработка категориальных данных

```
In [41]: # Выберем категориальные колонки с пропущенными значениями
# Цикл по колонкам датасета
cat_cols = []
for col in data.columns:
    # Количество пустых значений
    temp_null_count = data[data[col].isnull()].shape[0]
    dt = str(data[col].dtype)
    if temp_null_count > 0 and (dt == 'object'):
        cat_cols.append(col)
        temp_perc = round((temp_null_count / total_count) * 100.0, 2)
        print('Колонка {}. Тип данных {}. Количество пустых значений {},
if len(cat_cols) == 0:
    print("msg: Пропусков значений в категориальных данных не обнаружено")
```

msg: Пропусков значений в категориальных данных не обнаружено

Попробуем вручную удалить некоторые данные:

```
In [46]: # читаем данные набора
data2 = pd.read_csv('JupyterNotebooks/data/hfi_cc_2018_3.csv', sep=',')

for col in data2.columns:
    # Количество пустых значений
    temp_null_count = data2[data2[col].isnull()].shape[0]
    dt = str(data2[col].dtype)
    if temp_null_count > 0 and (dt == 'object'):
        cat_cols.append(col)
```

```
temp_perc = round((temp_null_count / total_count) * 100.0, 2)
print('Колонка {}. Тип данных {}. Количество пустых значений {},
if len(cat_cols) == 0:
print("msg: Пропусков значений в категориальных данных не обнаружен")
```

Колонка ISO\_code. Тип данных object. Количество пустых значений 6, 0.41%.  
Колонка region. Тип данных object. Количество пустых значений 2, 0.14%.

Именно в этом случае можно просто удалить строки с пропущенными данными:

```
In [52]: cat_cols_set = set(cat_cols)
cat_cols_set
```

```
Out[52]: {'ISO_code', 'region'}
```

```
In [62]: # Удаление строк, содержащих пустые категориальные значения
for col in cat_cols_set:
    data2_new = data2[data2[col].notna()]
(data2.shape, data2_new.shape)
```

```
Out[62]: ((1458, 123), (1452, 123))
```

В случае взятого мною набора данных такой способ более подходит: просто выбираем признаки которые в дальнейшем хотим отследить, проверить зависимости и тп, и удаляем строки, у которых у этих признаков значения не установлены.

## Масштабирование данных

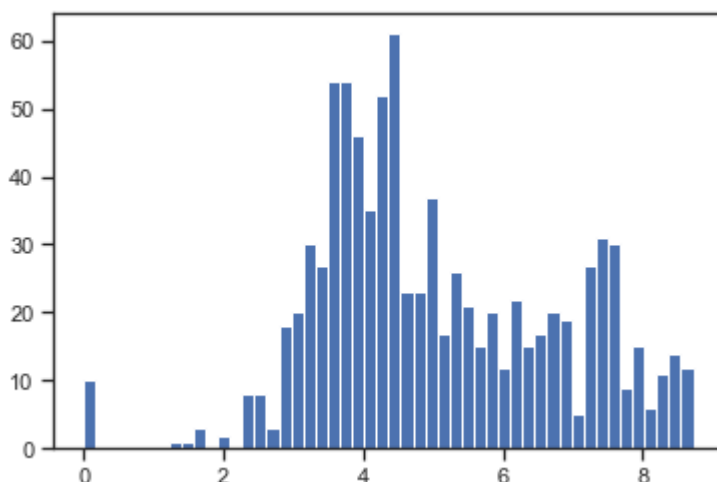
Термины "масштабирование" и "нормализация" часто используются как синонимы. Масштабирование предполагает изменение диапазона измерения величины, а нормализация - изменение распределения этой величины.

```
In [63]: from sklearn.preprocessing import MinMaxScaler, StandardScaler, Normalizer
```

### Масштабирование МиниМакс

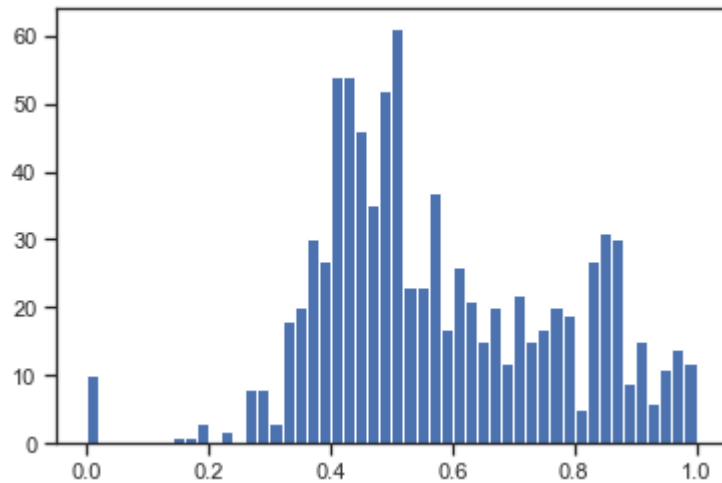
```
In [76]: sc1 = MinMaxScaler()
sc1_data = sc1.fit_transform(data2[['pf_rol_criminal']])
```

```
In [65]: plt.hist(data2['pf_rol_criminal'], 50)
plt.show()
```





```
In [75]: plt.hist(sc1_data, 50)
plt.show()
```



## Масштабирование на основе Z-оценки

```
In [77]: sc2 = StandardScaler()
sc2_data = sc2.fit_transform(data2[['pf_rol_criminal']])
```

```
In [78]: plt.hist(sc2_data, 50)
plt.show()
```

