In [161... import numpy as np import pandas as pd import seaborn as sns import matplotlib.pyplot as plt In [162... data = pd.read\_csv('restaurant-scores-lives-standard.csv',sep=",") In [163... data.shape (53973, 23) Out[163... In [164... data.dtypes business\_id int64 Out[164... object business\_name business\_address object business\_city object business\_state object object business\_postal\_code business\_latitude float64 business\_longitude float64 business\_location object business\_phone\_number float64 inspection\_id object inspection\_date object float64 inspection\_score object inspection\_type violation\_id object violation description object object risk\_category Neighborhoods (old) float64 Police Districts float64 float64 Supervisor Districts Fire Prevention Districts float64 Zip Codes float64 Analysis Neighborhoods float64 dtype: object In [165... # проверим есть ли пропущенные значения isnull = data.isnull().sum() print (isnull) business\_id 0 business\_name 0 business\_address 0 business\_city 0 business\_state business\_postal\_code 1018 19556 business\_latitude business\_longitude 19556 business\_location 19556 36938 business\_phone\_number inspection\_id 0 inspection\_date inspection\_score 13610 inspection\_type violation\_id 12870 violation\_description 12870 12870 risk\_category Neighborhoods (old) 19594 Police Districts 19594 19594 Supervisor Districts Fire Prevention Districts 19646 19576 Zip Codes Analysis Neighborhoods 19594 dtype: int64 Для обработки пропусков возьмём один числовой признак(inspection\_score) и категориальный признак (violation\_description) Будем использовать методику импьютации, т.к. были выбраны столбцы с не огромными пропусками данных (максимум 25%) In [166... dataint = data['inspection\_score'] dataint.head() NaN 96.0 NaN 3 NaN NaN Name: inspection\_score, dtype: float64 In [167... num\_cols=[] for col in data.columns: if col=='inspection\_score': num\_cols.append(col) In [168... sort\_null\_data = data[num\_cols] data\_inspection = sort\_null\_data[['inspection\_score']] In [169... from sklearn.impute import SimpleImputer from sklearn.impute import MissingIndicator In [170... indicator = MissingIndicator() mask\_missing\_values\_only = indicator.fit\_transform(data\_inspection) mask\_missing\_values\_only Out[170... array([[ True], [False], [ True], [False], [False], [False]]) In [171... for col in num cols: plt.hist(data[col], 100) plt.xlabel(col) plt.show() 4000 3500 3000 -2500 2000 1500 1000 500 50 60 70 80 90 100 inspection\_score In [172... num\_cols\_inspection\_score = sort\_null\_data num\_cols\_inspection\_score.head() Out[172... inspection\_score NaN 96.0 NaN NaN NaN In [173... indicator = MissingIndicator() mask\_missing\_values\_only = indicator.fit\_transform(num\_cols\_inspection\_score) mask\_missing\_values\_only Out[173... array([[ True], [False], [True], ..., [False], [False], [False]]) In [174... strategies=['mean', 'median', 'most\_frequent'] In [175... def test\_num\_impute(strategy\_param): imp\_num = SimpleImputer(strategy=strategy\_param) data\_num\_imp = imp\_num.fit\_transform(num\_cols\_inspection\_score) return data\_num\_imp[mask\_missing\_values\_only] In [176... strategies[0], test\_num\_impute(strategies[0]) Out[176... ('mean' array([86.22679186, 86.22679186, 86.22679186, ..., 86.22679186, 86.22679186, 86.22679186])) In [177... strategies[1], test\_num\_impute(strategies[1]) ('median', array([87., 87., 87., 87., 87., 87., 87.])) Out[177... In [178... strategies[2], test\_num\_impute(strategies[2]) ('most\_frequent', array([90., 90., 90., ..., 90., 90., 90.])) Возьмём моду для значений пропусков In [179... imp\_num = SimpleImputer(strategy=strategies[2]) data\_int\_full = imp\_num.fit\_transform(num\_cols\_inspection\_score) data\_int\_full Out[179... array([[90.], [96.], [90.], [92.], [76.], [80.]]) num\_cols\_inspection\_score\_new = num\_cols\_inspection\_score num\_cols\_inspection\_score\_new = data\_int\_full nmp = num\_cols\_inspection\_score.to\_numpy data.iloc[:,12] = num\_cols\_inspection\_score\_new ##data.replace({'inspection\_score' : { nmp : num\_cols\_inspection\_score\_new}}) data['inspection\_score'] 90.0 Out[180... 90.0 90.0 90.0 53968 80.0 53969 90.0 53970 92.0 53971 76.0 53972 80.0 Name: inspection\_score, Length: 53973, dtype: float64 In [181... isnull = data.isnull().sum() print (isnull) business\_id business\_name business\_address business\_city business\_state business\_postal\_code 1018 19556 business\_latitude business\_longitude 19556 business\_location 19556 business\_phone\_number 36938 inspection\_id inspection\_date inspection\_score inspection\_type 12870 violation\_id 12870 violation\_description risk\_category 12870 Neighborhoods (old) 19594 19594 Police Districts Supervisor Districts 19594 Fire Prevention Districts 19646 Zip Codes 19576 Analysis Neighborhoods 19594 dtype: int64 Получается мы заполнили столбец новыми данными Обработка категориальных признаков In [182... num\_cols=[] for col in data.columns: if col=='violation\_description': num\_cols.append(col) In [183... sort\_null\_data\_obj = data[num\_cols] data\_Desc = sort\_null\_data\_obj[['violation\_description']] data\_Desc.head() Out[183... violation\_description NaN **1** Inadequately cleaned or sanitized food contact... 2 NaN 3 High risk vermin infestation In [184... implicator = SimpleImputer(missing\_values=np.nan, strategy='constant', fill\_value='NA') Desc\_values = implicator.fit\_transform(data\_Desc) Desc\_values Out[184... array([['NA'], 'Inadequately cleaned or sanitized food contact surfaces'], ['NA'], ['Foods not protected from contamination'], ['Inadequate food safety knowledge or lack of certified food safety manager'], ['Food safety certificate or food handler card not available']], dtype=object) In [185... data\_Desc['violation\_description'].unique() Out[185... array([nan, 'Inadequately cleaned or sanitized food contact surfaces', 'High risk vermin infestation', 'Moderate risk food holding temperature', 'Improper storage use or identification of toxic substances', 'Improper or defective plumbing', 'No hot water or running water', 'Inadequate and inaccessible handwashing facilities', 'Low risk vermin infestation', 'Insufficient hot water or running water', 'Foods not protected from contamination', 'Improper food storage', 'Unapproved or unmaintained equipment or utensils', 'Wiping cloths not clean or properly stored or inadequate sanitizer', 'Food safety certificate or food handler card not available', 'Improper thawing methods', 'Improper cooling methods', 'Unclean or degraded floors walls or ceilings', 'Inadequate food safety knowledge or lack of certified food safety manager', 'Other low risk violation', 'Inadequate procedures or records for time as a public health control', 'Unclean nonfood contact surfaces', 'High risk food holding temperature', 'Unclean or unsanitary food contact surfaces', 'Noncompliance with HAACP plan or variance', 'Permit license or inspection report not posted', 'Inadequate dressing rooms or improper storage of personal items', 'Unapproved living quarters in food facility', 'Inadequate or unsanitary refuse containers or area or no garbage service', 'Inadequate HACCP plan record keeping', 'Moderate risk vermin infestation', 'Inadequate warewashing facilities or equipment', 'Mobile food facility not operating with an approved commissary', 'Inadequate ventilation or lighting', 'No person in charge of food facility', 'Unauthorized or unsafe use of time as a public health control measure', 'Improper storage of equipment utensils or linens', 'Unclean unmaintained or improperly constructed toilet facilities', 'Improper reheating of food', 'Sewage or wastewater contamination', 'Unclean hands or improper use of gloves', 'Unpermitted food facility', 'No thermometers or uncalibrated thermometers', 'Inadequate sewage or wastewater disposal', 'Contaminated or adulterated food', 'Mobile food facility with unapproved operating conditions', 'Employee eating or smoking', 'Improper food labeling or menu misrepresentation', 'Food in poor condition', 'Worker safety hazards', 'Noncompliance with shell fish tags or display', 'Other moderate risk violation', 'Unsanitary employee garments hair or nails', 'Unapproved food source', 'Improperly washed fruits and vegetables', 'No plan review or Building Permit', 'No restroom facility within 200 feet of mobile food facility', 'Improperly displayed mobile food permit or signage', 'Consumer advisory not provided for raw or undercooked foods', 'Noncompliance with Gulf Coast oyster regulation', 'Other high risk violation', 'Reservice of previously served foods', 'Improper cooking time or temperatures', 'Non service animal', 'Mobile food facility stored in unapproved location', 'Discharge from employee nose mouth or eye'], dtype=object) In [186... data\_frame = pd.DataFrame({'violation\_description': Desc\_values.T[0]}) from sklearn.preprocessing import LabelEncoder, OneHotEncoder In [188... data\_frame.head(10) Out[188... violation\_description 1 Inadequately cleaned or sanitized food contact... 2 NA 3 High risk vermin infestation NA 7 Inadequately cleaned or sanitized food contact... Moderate risk food holding temperature In [189... imp2 = SimpleImputer(missing\_values=np.nan, strategy='most\_frequent') data\_imp2 = imp2.fit\_transform(data\_Desc) data\_imp2 Out[189... array([['Unclean or degraded floors walls or ceilings'], ['Inadequately cleaned or sanitized food contact surfaces'], ['Unclean or degraded floors walls or ceilings'], ['Foods not protected from contamination'], ['Inadequate food safety knowledge or lack of certified food safety manager'], ['Food safety certificate or food handler card not available']], dtype=object) In [190... data.iloc[:,15] = data\_imp2 data['violation\_description'] Unclean or degraded floors walls or ceilings Out[190... 0 Inadequately cleaned or sanitized food contact... Unclean or degraded floors walls or ceilings Unclean or degraded floors walls or ceilings High risk vermin infestation Inadequately cleaned or sanitized food contact... 53968 53969 Unclean or degraded floors walls or ceilings 53970 Foods not protected from contamination Inadequate food safety knowledge or lack of ce... 53971 Food safety certificate or food handler card n... Name: violation\_description, Length: 53973, dtype: object In [191... one = OneHotEncoder() data label hot = one.fit transform(data frame) In [192... data\_label\_hot Out[192... <53973x66 sparse matrix of type '<class 'numpy.float64'>' with 53973 stored elements in Compressed Sparse Row format> In [193... data\_label\_hot.todense()[0:5] Out[193... 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.], 0., 0.], 0., 0.], [0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.]0., 0.]]) Для дальнейшего построения моделей я буду использовать карелирующие признаки, потому что на их основе можно построить более качественные модели машинного обучения. Также приоритет получат те признаки, в которых отсутствует пропуск данных или он минимален. In [ ]: