

PONTIFÍCIA UNIVERSIDADE CATÓLICA DE MINAS GERAIS
NÚCLEO DE EDUCAÇÃO A DISTÂNCIA
Pós-graduação *Lato Sensu* em Ciência de Dados e Big Data

Marcelo Joventino Freitas

CLASSIFICAÇÃO DA ATIVIDADE ECONÔMICA SPED EFD-CONTRIBUIÇÕES
UTILIZANDO A DESCRIÇÃO DOS SERVIÇOS PRESTADOS

Vitória
2021

Marcelo Joventino Freitas

**CLASSIFICAÇÃO DA ATIVIDADE ECONÔMICA SPED EFD-CONTRIBUIÇÕES
UTILIZANDO A DESCRIÇÃO DOS SERVIÇOS PRESTADOS**

Trabalho de Conclusão de Curso apresentado
ao Curso de Especialização em Ciência de
Dados e Big Data como requisito parcial à
obtenção do título de especialista.

Vitória

2021

SUMÁRIO

1. Introdução.....	4
1.1. Contextualização	4
1.2. O problema proposto	4
2. Coleta de Dados	4
3. Processamento/Tratamento de Dados	5
4. Análise e Exploração dos Dados	7
5. Criação de Modelos de Machine Learning	13
6. Apresentação dos Resultados	15
7. Links	18
APÊNDICE.....	19

1. Introdução

1.1. Contextualização

A Escrituração Fiscal Digital Contribuições (EFD Contribuições) trata de arquivo digital instituído no Sistema Público de Escrituração Digital (SPED), a ser utilizado pelas pessoas jurídicas de direito privado na escrituração da Contribuição Previdenciária sobre a Receita Bruta (CPRB), da Contribuição para o PIS/Pasep e da Cofins, com base nas operações representativas de receitas auferidas, custos e despesas.

Foi instituído um código de atividade para cada tipo de serviço ou produto sujeito à CPRB, utilizado na apuração do valor da CPRB no bloco P100 da EFD-Contribuições. As tabelas de alíquotas da EFD-Contribuições são montadas a partir deste código.

A apuração no bloco P100 é feita com base nas receitas de venda de produtos e serviços contidos nos blocos C170 e A170 de forma manual.

Há dificuldade na associação dos serviços contidos no bloco A170 com a atividade econômica do bloco P100, principalmente quando há mais de um tipo de serviço prestado no período analisado, pois essa etapa é feita manualmente.

1.2. O problema proposto

O problema proposto é aferir o código de atividade econômica a ser informado no bloco P100 da EFD-Contribuições a partir da descrição do serviço informada no bloco A170 da EFD-Contribuições.

A descrição do serviço no bloco A170 também pode conter a classificação do serviço com base na Lei Complementar 116 de 31/07/2003.

2. Coleta de Dados

A coleta de dados foi feita na base A170 (descrição dos serviços) para os declarantes de 1 código de atividade na base P100 do SPED EFD-Contribuições no

período entre os anos 2016 e 2020. Os registros foram rotulados com o código de atividade informado na base P100.

Foi escolhido o mês de janeiro, abril e setembro para montagem do dataset de treino, e o mês de junho para o dataset de validação, aplicando-se amostragem aleatória para obtenção de classes balanceadas.

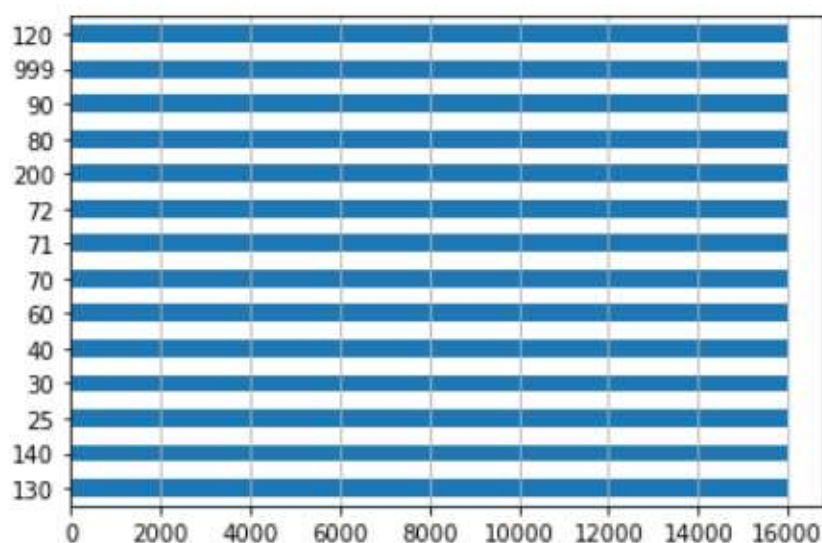
Os datasets possuem a seguinte estrutura:

Nome da coluna	Descrição	Tipo
ativ	Código da atividade econômica (alvo)	Inteiro
cd_lc	Código do serviço pela LC 116	Texto
desc_it	Descrição do serviço	Texto

3. Processamento/Tratamento de Dados

Primeiramente identificou-se os declarantes na base P100 da EFD-Contribuições com somente 1 código de atividade econômica informado no ano em análise. Esse código de atividade foi usado como rótulo na base A170 da EFD-Contribuições com as descrições de serviço informadas.

Foram então consultados os registros de descrição do serviço da base A170 para aqueles declarantes selecionados inicialmente. As consultas foram realizadas nos meses de janeiro, abril, junho e setembro de 2016 a 2020, separando-se o mês junho para validação. Foi aplicada uma amostragem aleatória com o objetivo de separar 16 mil amostras para cada classe no dataset de treino:



Os registros das classes 91 a 98 foram agrupados na classe 90 (construção civil). Foi criada a classe 999 de registros que não pertencem a nenhuma das outras classes em estudo.

Foram separados 1000 registros de cada classe para o dataset de validação no mês de junho.

As descrições de serviço com erro de leitura por conta de acentos e caracteres especiais foram excluídos em virtude da quantidade não ser significativa.

Foi efetuada a limpeza do campo descrição com as seguintes ações:

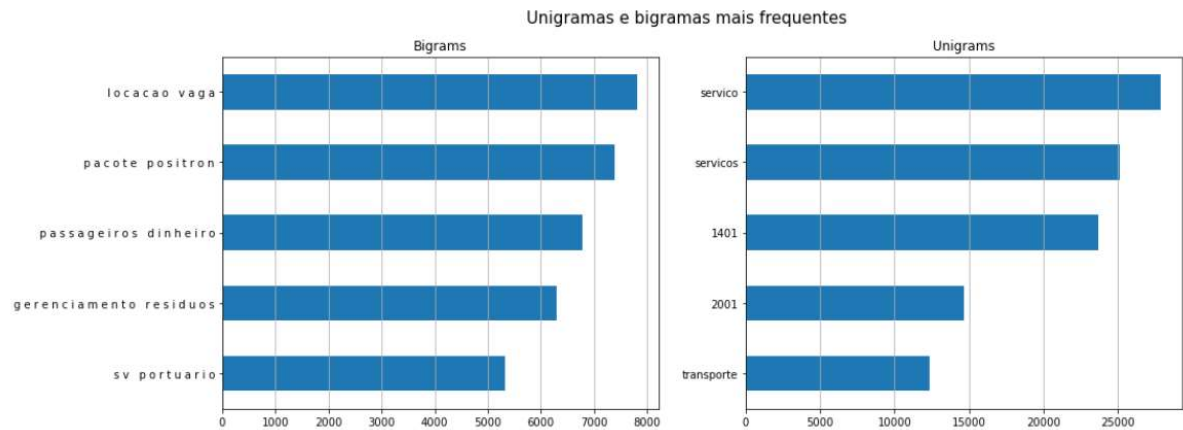
- a) Remoção de pontuação;
- b) Remoção de caracteres numéricos;
- c) Remoção de stopwords;
- d) Remoção de palavras com 1 caracter

O campo com o código da lei complementar 116 teve apenas a remoção de pontuação e a manutenção do código numérico.

Foi criado o campo 'texto_limpo' com o código da lei complementar 116 + descrição após a limpeza.

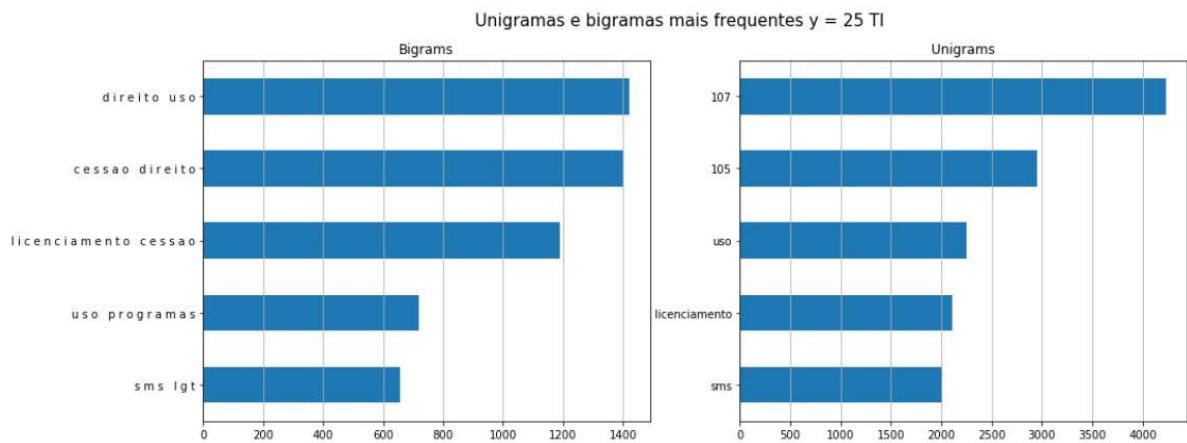
Foram excluídos os registros do dataset de treino quando:

- a) 'texto_limpo' contém 'call center' e 'ativ' é diferente de 30;
- b) 'texto_limpo' contém 'cobranca' e 'ativ' é diferente de 30;
- c) 'texto_limpo' contém 'empreitada' e 'ativ' não está em '90,120';
- d) 'texto_limpo' contém 'subempreitada' e 'ativ' não está em '90,120';
- e) 'texto_limpo' contém 'construcao' e 'ativ' não está em '90,120';
- f) 'texto_limpo' contém 'carga' e 'ativ' é igual a 999;
- g) 'texto_limpo' contém 'descarga' e 'ativ' é igual a 999;
- h) 'texto_limpo' contém 'portuario' e 'ativ' é igual a 999;

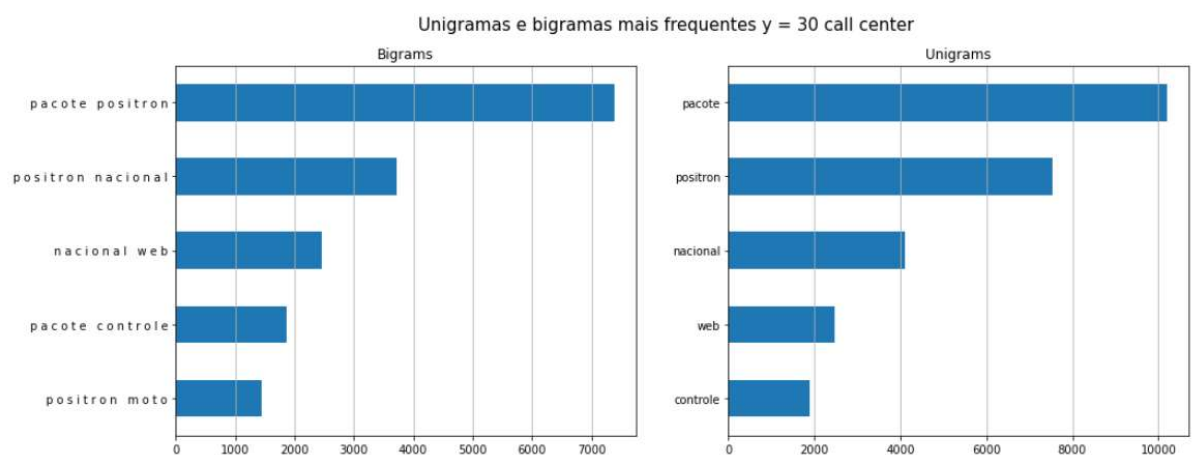


Também foram identificados os unigramas e bigramas mais frequentes por classe:

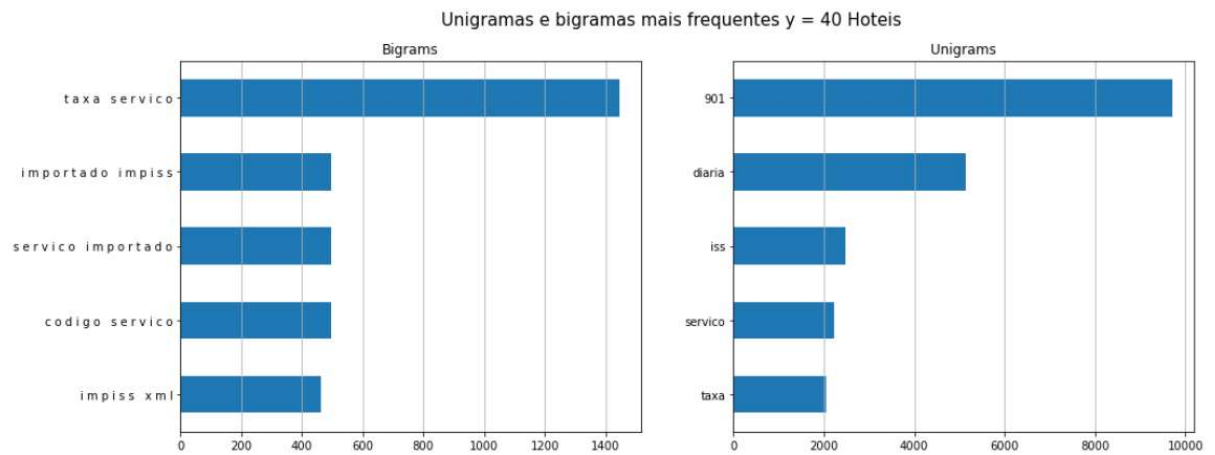
a) Classe 25 – serviços de TI:



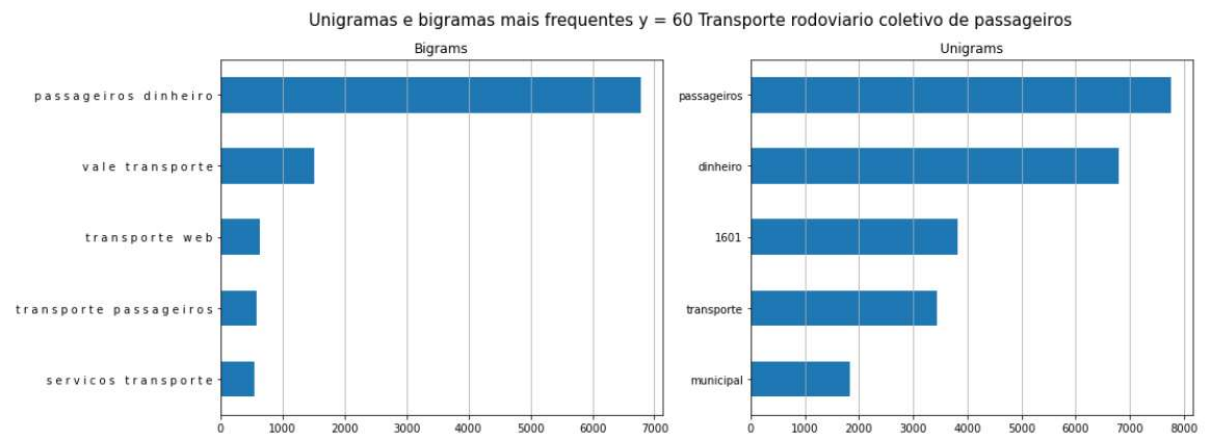
a) Classe 30 – serviços de call center:



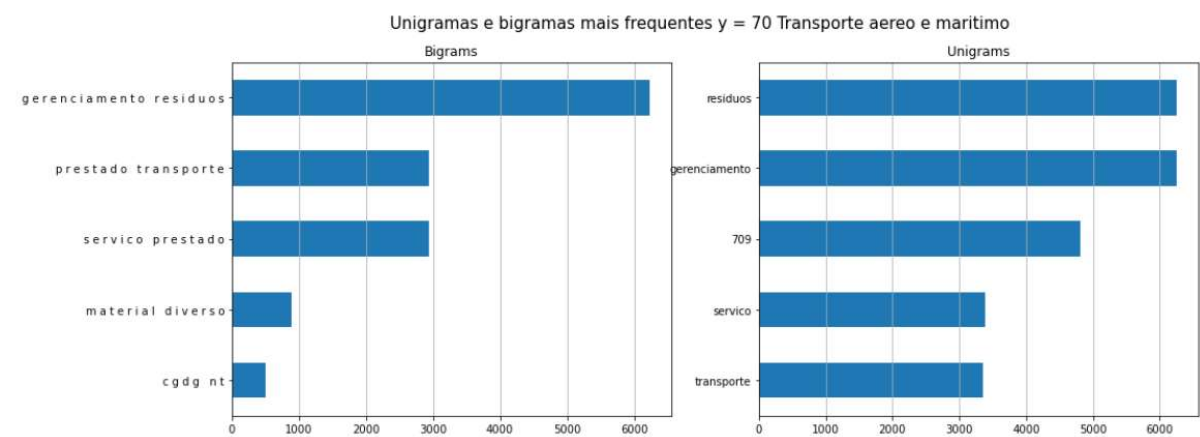
b) Classe 40 – setor hoteleiro:



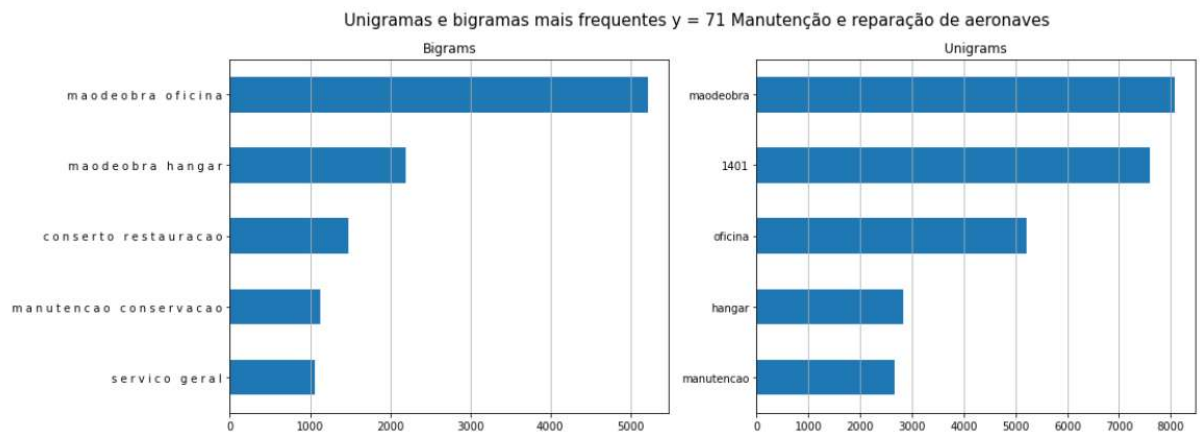
c) Classe 60 – transporte rodoviário coletivo de passageiros:



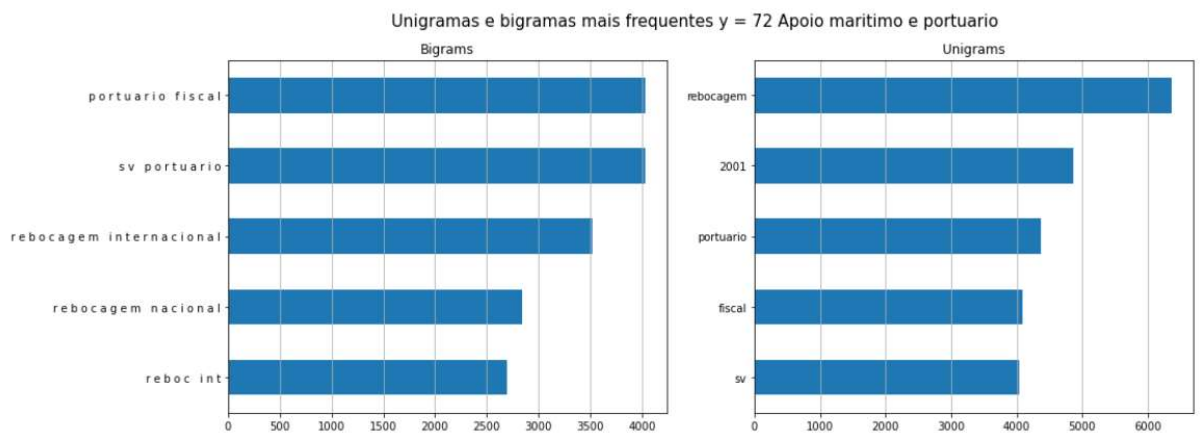
d) Classe 70 – transporte aéreo e marítimo:



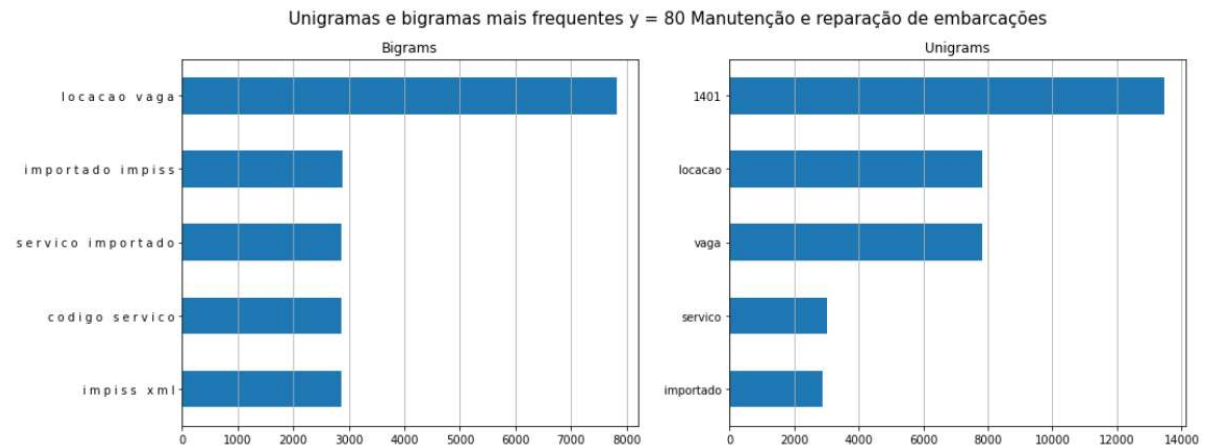
e) Classe 71 – manutenção e reparação de aeronaves:



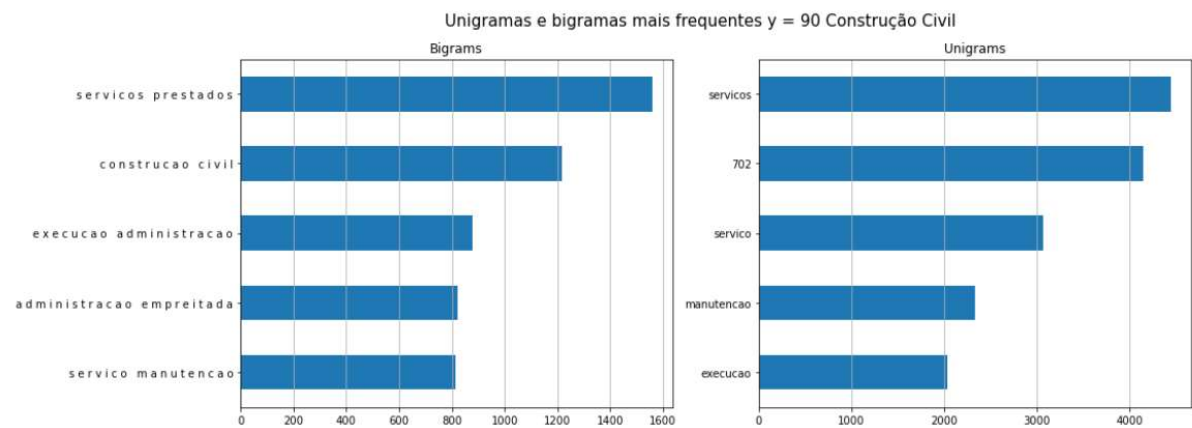
f) Classe 72 – apoio marítimo e portuário:



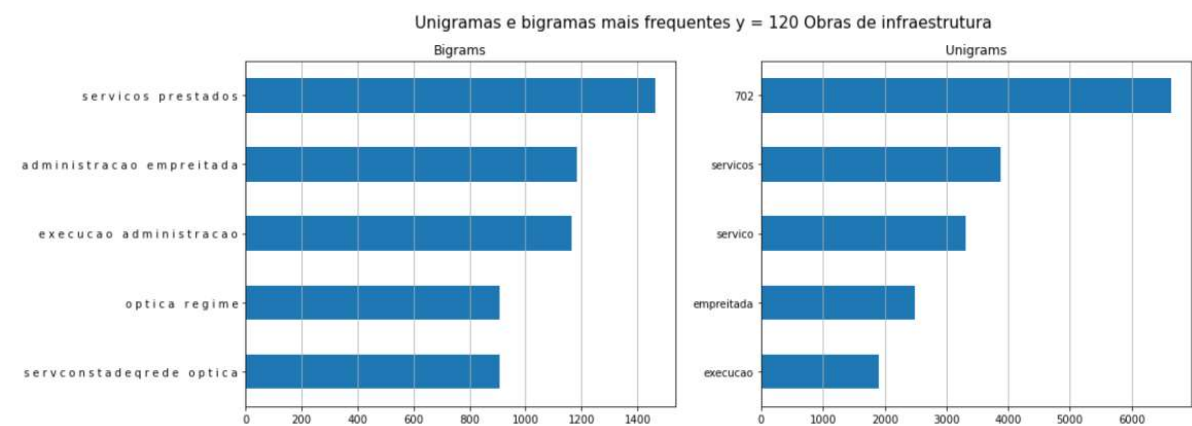
g) Classe 80 – manutenção e reparação de embarcações:



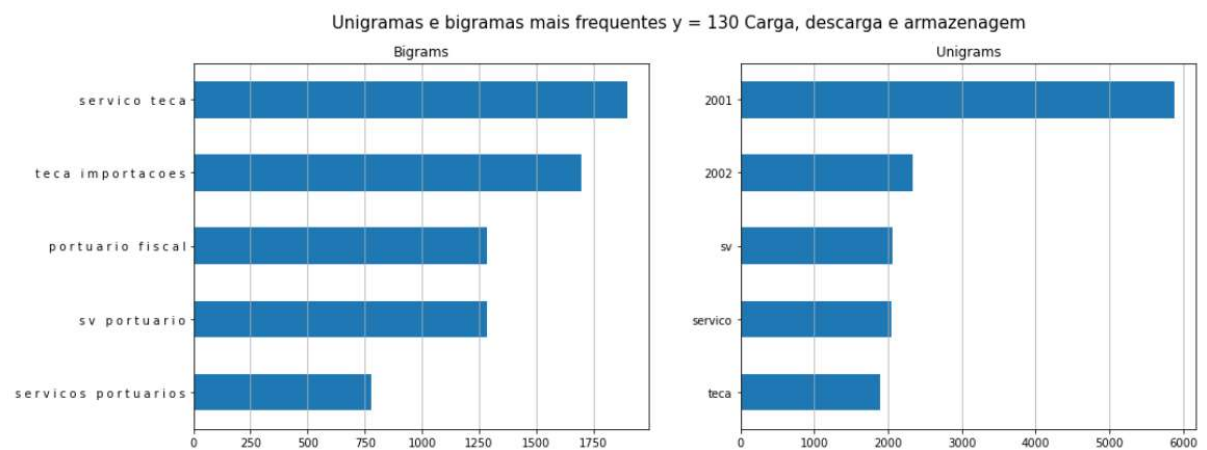
h) Classe 90 – construção civil:



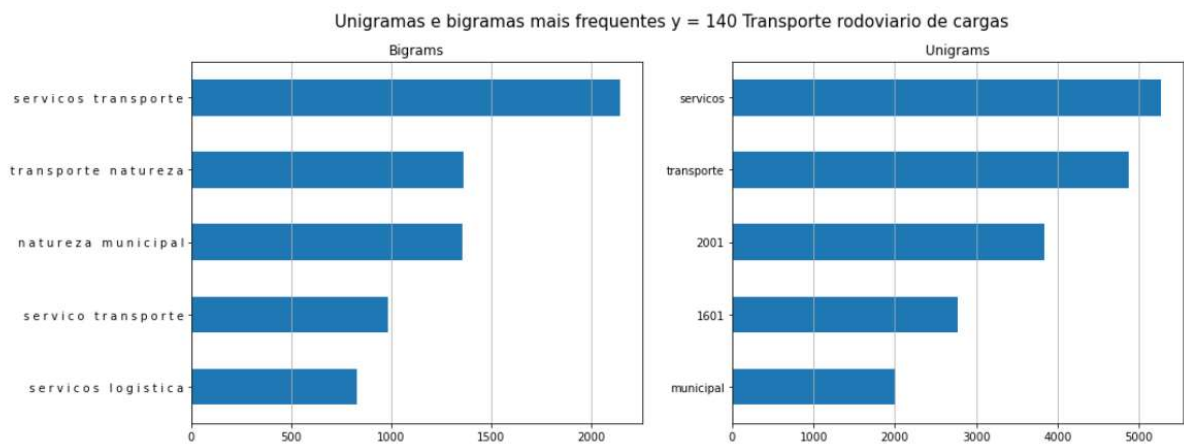
i) Classe 120 – obras de infraestrutura:



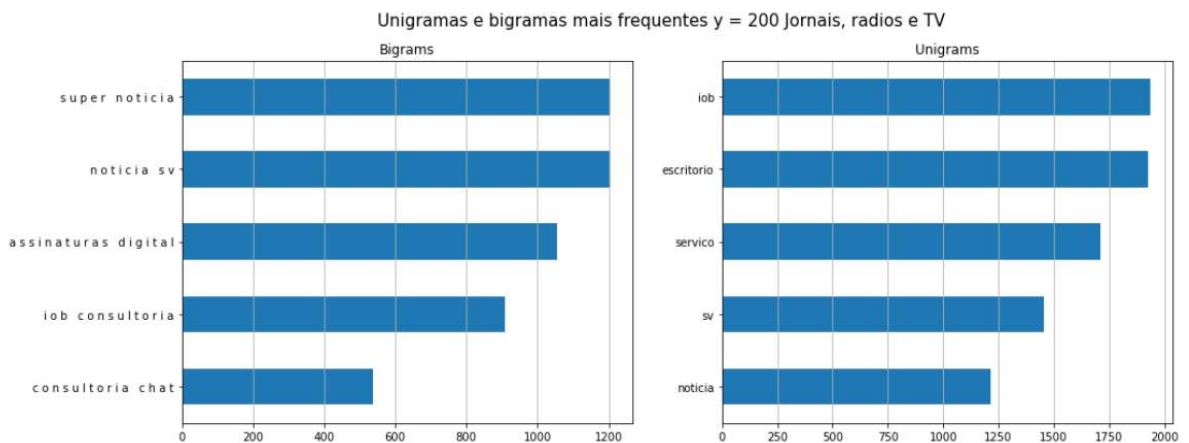
j) Classe 130 – carga, descarga e armazenagem:



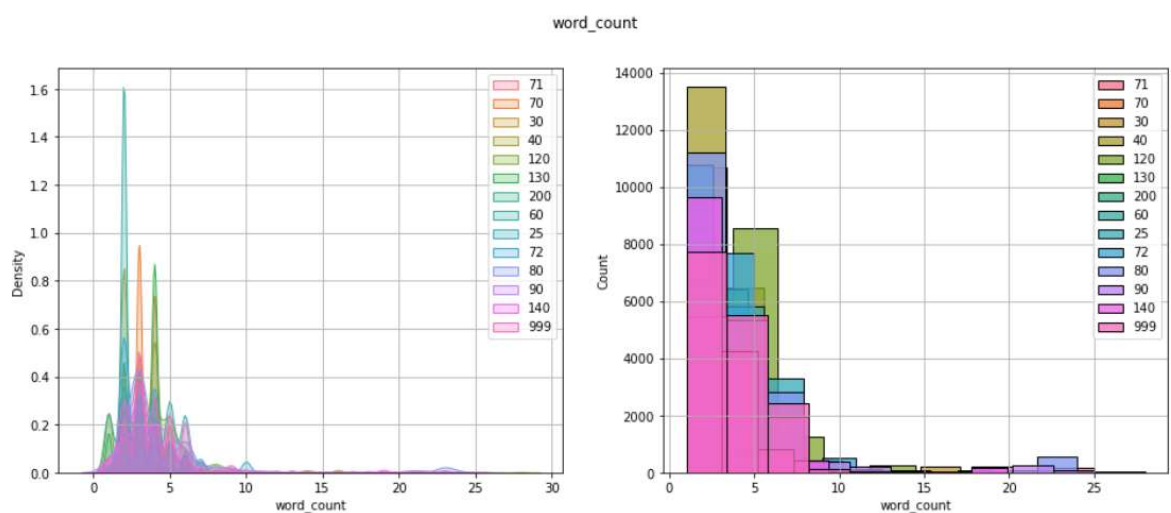
k) Classe 140 – transporte rodoviário de cargas:



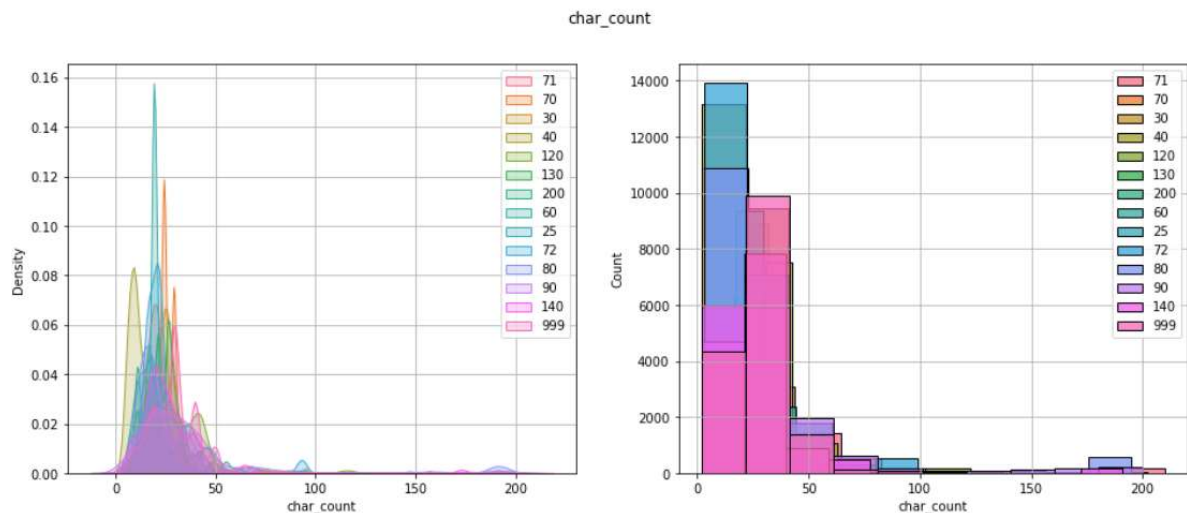
I) Classe 200 – jornais, rádios e televisão:



Distribuição da contagem de palavras por registro do dataset de treino:



Distribuição da contagem de caracteres por registro do dataset de treino:



5. Criação de Modelos de Machine Learning

Foram instanciados vários métodos de vetorização de documentos para avaliação da melhor adequação ao problema analisado para unigramas e bigramas:

- Binária: palavras são marcadas como presentes (1) ou ausentes (0);
- Contagem: contagem de ocorrências de cada palavra como um número inteiro;
- Frequência: frequência de cada palavra no documento;
- Tfidf: frequência das palavras com penalização das palavras mais comuns (frequência + grau de raridade);
- Hash: frequência das palavras com limitação das features (colunas) na matriz de vetores (versão sklearn);

```
In [38]: 1 # unigramas
2 binary_vectorizer = CountVectorizer(binary=True)
3 count_vectorizer = CountVectorizer()
4 freq_vectorizer = TfidfVectorizer(use_idf=False)
5 tfidf_vectorizer = TfidfVectorizer()
6 hashing_vectorizer = HashingVectorizer(n_features = 2 ** 11)
7
8 # bigramas
9 binary_vectorizer2 = CountVectorizer(binary=True, ngram_range=(1,2))
10 count_vectorizer2 = CountVectorizer(ngram_range=(1,2))
11 freq_vectorizer2 = TfidfVectorizer(use_idf=False, ngram_range=(1,2))
12 tfidf_vectorizer2 = TfidfVectorizer(ngram_range=(1,2))
13 hashing_vectorizer2 = HashingVectorizer(n_features = 2 ** 11, ngram_range=(1,2))
```

- Teste qui-quadrado: função que determina quais palavras são estatisticamente relevantes, mantendo aquelas com significância de 95%:

```

91 def teste_chi2(docs=None, y=None, vectorizer=None, p_value_limit = 0.95):
92     ''' chi-squared test to determine whether a feature and target are independent - statistically relevant
93         keep only features with p_value_limit = 0.95 '''
94     X = vectorizer.fit_transform(docs)
95     X_names = vectorizer.get_feature_names()
96     print('tamanho_vocabulario = ', len(vectorizer.vocabulary_))
97     print(X.shape)
98     df_features = pd.DataFrame()
99     for cat in np.unique(y):
100         chi2_, p = chi2(X, y == cat)
101         df_features = df_features.append(pd.DataFrame({'feature': X_names, 'score': 1-p, 'y': cat}))
102     df_features = df_features.sort_values(['y', 'score'], ascending=[True, False])
103     df_features = df_features[df_features['score'] > p_value_limit]
104     print(df_features.sample(10))
105     print('vocab_testado:', df_features.feature.nunique())
106     return df_features.feature.unique().tolist()

```

A função seguinte recebe o algoritmo de aprendizagem de máquina e os vetorizadores para montar um DataFrame com os resultados dos testes aplicados. A função também possibilita aplicar uma redução de componentes pelo TruncatedSVD do sklearn:

```

In [47]: 1 results_ml = pd.DataFrame(columns = ['acc', 'precision_macro', 'recall_macro', 'f1_macro', 'mcc'])
2 def s_results(docs_train, ytrain, vectorizer=None, model=None, svd=None, docs_val=None, yval=None):
3     if svd:
4         pipe = Pipeline(steps=[('vectorizer', vectorizer),
5                                 ('svd', TruncatedSVD(n_components=100, n_iter=10, random_state=0)),
6                                 ('model', model)])
7     else:
8         pipe = Pipeline(steps=[('vectorizer', vectorizer),
9                                 ('model', model)])
10    pipe.fit(docs_train, ytrain)
11    yv_pred = pipe.predict(docs_val)
12    acc = accuracy_score(yval, yv_pred)
13    p_m = precision_score(yval, yv_pred, average='macro')
14    r_m = recall_score(yval, yv_pred, average='macro')
15    f1_m = f1_score(yval, yv_pred, average='macro')
16    mcc = matthews_corrcoef(yval, yv_pred)
17    return pd.Series(data=[acc, p_m, r_m, f1_m, mcc], index=['acc', 'precision_macro', 'recall_macro', 'f1_macro', 'mcc'])

```

Foram testados os seguintes algoritmos de aprendizagem de máquina:

- a) NB: Naive Bayes MultinomialNB;
- b) RF: Random Forest;
- c) LG: Regressão Logística;
- d) SVM: Suport Vector Machines.

E utilizadas as seguintes abreviaturas nas chamadas da função:

- a) uni: unigramas;
- b) bi: bigramas;
- c) chi2: redução de variáveis pelo método qui-quadrado;
- d) SVD: redução de variáveis pelo método SVD;
- e) Binary: vetorização binária;
- f) Count: vetorização por contagem;
- g) Freq: vetorização por frequência;
- h) Tfidf: vetorização por tfidf;
- i) Hash: vetorização hash.

Foram acrescentados blocos de código organizados por: algoritmo, vetorização, unigramas, bigramas, unigramas com teste qui-quadrado, unigramas com SVD, bigramas com teste qui-quadrado, bigramas com SVD. Segue como exemplo o bloco de código com o algoritmo Random Forest vetorizado por frequências:

```
In [64]: 1 # freq
2 results_ml.loc['RF_freq_uni'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
3                                           model= RandomForestClassifier(random_state=0), docs_val=docs_val, yval=yval)
4
5 results_ml.loc['RF_freq_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
6                                           svd=True,
7                                           model= RandomForestClassifier(random_state=0), docs_val=docs_val, yval=yval)
8
9 results_ml.loc['RF_freq_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer_chi2,
10                                              model= RandomForestClassifier(random_state=0), docs_val=docs_val, yval=yval)
11
12 results_ml.loc['RF_freq_bi'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer2,
13                                         model= RandomForestClassifier(random_state=0), docs_val=docs_val, yval=yval)
14
15 results_ml.loc['RF_freq_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer2,
16                                         svd=True,
17                                         model= RandomForestClassifier(random_state=0), docs_val=docs_val, yval=yval)
18
19 results_ml.loc['RF_freq_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer2_chi2,
20                                              model= RandomForestClassifier(random_state=0), docs_val=docs_val, yval=yval)
```

6. Apresentação dos Resultados

O dataset de validação foi montado com 1000 amostras por classe em mês distinto dos registros utilizados para treino do modelo. O DataFrame seguinte foi montado com os resultados dos testes aplicados, tendo sido utilizadas as seguintes métricas do sklearn para comparação dos resultados:

- a) acc: accuracy score;
- b) precision_macro: precision score com o parâmetro average = macro (calcula a métrica para cada classe e depois a média simples);
- c) recall_macro: recall score com o parâmetro average = macro;
- d) f1_macro: f1 score com o parâmetro average = macro;
- e) mcc: Matthews correlation coefficient.

modelo ML com melhor resultado

```
In [70]: 1 results_ml.sort_values(by=['mcc'], ascending=False).head(20)
```

```
Out[70]:
```

	acc	precision_macro	recall_macro	f1_macro	mcc
LG_binary_bi	0.8841	0.8903	0.8841	0.8850	0.8757
LG_count_bi	0.8837	0.8899	0.8837	0.8846	0.8752
SVM_tfidf_bi	0.8822	0.8907	0.8822	0.8835	0.8737
SVM_tfidf_uni_chi2	0.8816	0.8896	0.8816	0.8829	0.8731
SVM_freq_bi	0.8816	0.8900	0.8816	0.8828	0.8731
SVM_tfidf_bi_chi2	0.8816	0.8895	0.8816	0.8828	0.8730
SVM_tfidf_uni	0.8814	0.8894	0.8814	0.8826	0.8729

Regressão Logística com bigramas e vetorização binária teve o melhor resultado. Suport Vector Classification com bigramas vetorizado por tfidf ficou em 3º lugar.

O modelo com a regressão logística foi ainda ajustado com os seguintes hiperparâmetros obtendo-se uma pequena melhora nos scores de avaliação:

```
In [74]: 1 grid_search.best_params_
```

```
Out[74]: {'C': 10.0, 'solver': 'sag'}
```

Segue o relatório de classificação do sklearn evidenciando o desempenho por classe do melhor modelo:


```
1 print(classification_report(yval, yv_pred, digits=4))
```

	precision	recall	f1-score	support
25	0.8300	0.9030	0.8649	1000
30	0.9850	0.9880	0.9865	1000
40	0.9834	0.9490	0.9659	1000
60	0.9692	0.9750	0.9721	1000
70	0.9739	0.8970	0.9339	1000
71	0.9836	0.9600	0.9717	1000
72	0.9385	0.9920	0.9645	1000
80	0.9930	0.9870	0.9900	1000
90	0.7807	0.7440	0.7619	1000
120	0.8961	0.8190	0.8558	1000
130	0.8312	0.7240	0.7739	1000
140	0.7120	0.6800	0.6957	1000
200	0.9690	0.9060	0.9364	1000
999	0.6535	0.8900	0.7536	1000
accuracy			0.8867	14000
macro avg	0.8928	0.8867	0.8876	14000
weighted avg	0.8928	0.8867	0.8876	14000

Complementado pela matriz de confusão (colunas com valores reais e linhas com valores preditos):

```
1 # matriz de confuso: columns: True, index: Predict
2 cm = confusion_matrix(yval, yv_pred)
3 pd.DataFrame(data=cm, columns=np.sort(val.y.unique()), index=np.sort(val.y.unique()))
```

	25	30	40	60	70	71	72	80	90	120	130	140	200	999
25	903	6	2	6	0	0	0	0	17	2	0	5	6	53
30	0	988	1	0	0	0	0	0	2	0	0	0	1	8
40	6	0	949	0	0	1	0	0	5	1	0	4	1	33
60	0	0	1	975	0	0	1	0	3	0	0	17	1	2
70	0	0	0	3	897	1	2	0	5	7	21	33	1	30
71	1	0	0	0	0	960	0	0	26	3	0	0	1	9
72	1	0	0	0	0	0	992	0	0	1	4	0	0	2
80	0	0	0	0	2	0	1	987	0	1	0	4	0	5
90	10	0	2	0	1	0	0	6	744	62	0	3	1	171
120	4	1	1	0	0	3	0	0	99	819	1	5	12	55
130	0	0	2	0	7	4	54	0	0	0	724	183	0	26
140	125	5	2	19	7	1	5	0	8	1	119	680	1	27
200	24	0	0	0	1	0	1	0	4	2	0	11	906	51
999	14	3	5	3	6	6	1	1	40	15	2	10	4	890

Percebe-se acima dificuldade de classificação entre as classes 90 e 120; e entre as classes 130 e 140, que pode ser explicado pela própria composição das classes.

As classes 90 e 120 são ligadas à atividade de obras de construção civil, sendo a classe 120 específica na parte de infraestrutura.

As classes 130 e 140 tem em comum o transporte de carga, sendo a classe 130 específica em carga, descarga e armazenagem; e a classe 140 específica em transporte rodoviário de carga.

Os erros de classificação acima também ficam minimizados quando o objetivo for buscar a alíquota aplicável à classe, pois a alíquota da classe 90 é igual à alíquota da classe 120, o mesmo ocorrendo com as classes 130 e 140.

A classe 999 constitui-se de uma amostragem de registros não pertencentes às classes anteriores. A classe 999 foi atribuída quando a limpeza do texto deixou a linha vazia no dataset de validação. Falhas na predição da classe também tem relação ao regime opcional da CPRB. Assim, podemos ter registros de empresas cujo serviço enquadra-se numa das classes acima, mas que não fez opção pelo regime da CPRB no ano.

Dessa forma, considero a acurácia obtida pelo modelo satisfatória, principalmente porque o objetivo do modelo é encontrar a classe do serviço (não é escopo deste trabalho classificar a opção pelo regime da CPRB).

7. Links

Os dados utilizados são de propriedade da Receita Federal do Brasil, e a autorização de utilização não permitiu a disponibilização dos datasets gerados. Dessa forma, serão disponibilizados: o script com criação do modelo, arquivo com o modelo 'ML_model.pickle', vocabulário utilizado e script exemplificando o uso do modelo construído para classificar as descrições de serviço como proposto.

Link para o vídeo: <https://youtu.be/AzIus0IaVQI>

Link para o repositório: https://github.com/mjoventino/cod_atividade_cprb

APÊNDICE

Programação/Scripts

Script em Python com a implementação do modelo preditivo do código de atividade econômica.

Script em Python utilizando o modelo construído para predição da atividade econômica a partir de descrições de serviço digitadas.

Processamento da descrição do serviço de nota fiscal para predição da atividade econômica

```
1 Foram utilizadas amostras de dias coletados em 4 meses dos anos 2016 a 2020
2 O mês de junho foi escolhido para teste
3 Foi utilizada a descrição do serviço + código de serviço da Lei Complementar 116/2003 quando informado
```

```
1 Conda Env:
2 conda install -c anaconda pydot          # and graphviz: graph models
3 conda install -c conda-forge fastparquet  # compression gzip rodou no linux
4
5 NLP
6 conda install -c anaconda nltk
7 conda install -c anaconda seaborn
8 conda install -c conda-forge wordcloud
9 conda install -c anaconda scikit-learn
10
11 conda install jupyter
12 # criar as pastas models e best_models
```

```
1 # na 1ª execução: download de stopwords
2 nltk.download('stopwords')
```

In [1]:

```
1  ## for data
2  import numpy as np
3  import pandas as pd
4  import pickle
5  ## for plotting
6  import matplotlib.pyplot as plt
7  import seaborn as sns
8  ## for processing
9  from nltk.corpus import stopwords
10 from nltk import bigrams
11 import re
12 import string
13 from collections import Counter
14 from sklearn.decomposition import TruncatedSVD
15 from sklearn.feature_selection import chi2
16 from sklearn.preprocessing import OneHotEncoder, LabelBinarizer
17 from sklearn.pipeline import Pipeline
18 ## for bag-of-words
19 from wordcloud import WordCloud
20 #from sklearn.model_selection import train_test_split
21 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer, HashingVectorizer
22 ## for ML
23 from sklearn.naive_bayes import MultinomialNB
24 from sklearn.svm import SVC
25 from sklearn.ensemble import RandomForestClassifier
26 from sklearn.linear_model import LogisticRegression
27 from sklearn.metrics import accuracy_score, balanced_accuracy_score, f1_score, precision_score, recall_score
28 from joblib import dump, load
```

In [2]:

```
1  d_configuracao = {'display.max_columns': 30,
2                    'display.expand_frame_repr': True,
3                    'display.max_rows': 120,
4                    'display.precision': 2,
5                    'display.show_dimensions': True,
6                    'display.float_format': '{:,.4f}'.format}
7
8  for op, value in d_configuracao.items():
9      pd.set_option(op, value)
10     print(op, value)
```

```
display.max_columns 30
display.expand_frame_repr True
display.max_rows 120
display.precision 2
display.show_dimensions True
display.float_format <built-in method format of str object at 0x7f48ddd5f730>
```

In [3]:

```

1 def add2vocab(doc, vocab):
2     tokens = clean_doc(doc, False)
3     vocab.update(tokens)
4
5 def lc2vocab(doc, vocab):
6     l_lc = []
7     if doc is not None:
8         doc = re.sub(r'^\w\s', '', str(doc))
9         l_lc.append(doc)
10        vocab.update(l_lc)
11
12 def add2vocab2(bigramas, vocab2):
13     tokens = bigramas.split(',')
14     vocab2.update(tokens)
15
16 def clean_doc(doc=None, string_=True):
17     ''' # turn a doc into clean tokens '''
18     # split into tokens by white space
19     tokens = doc.split()
20     # prepare regex for char filtering
21     re_punc = re.compile('[%s]' % re.escape(string.punctuation))
22     # remove punctuation from each word
23     tokens = [re_punc.sub('', w) for w in tokens]
24     # remove remaining tokens that are not alphabetic
25     tokens = [word.lower() for word in tokens if not word.isdigit()]
26     # filter out stop words
27     stop_words = set(stopwords.words('portuguese'))
28     stop_words.update(['<br />'])
29     tokens = [w for w in tokens if not w in stop_words]
30     # filter out short tokens
31     tokens = [word for word in tokens if len(word) > 1]
32     if string_: return " ".join(tokens)
33     else: return tokens
34
35 def save_list(lines, filename):
36     ''' convert lines to a single blob of text '''
37     data = '\n'.join(lines)
38     # open file
39     file = open(filename, 'w')
40     # write text
41     file.write(data)
42     # close file
43     file.close()
44
45 def doc_to_line(cols):
46     ''' limpa e filtra vocab '''
47     lc = cols[0]
48     text = cols[1]
49     l_text_clean = clean_doc(text, False)
50     l_text_clean = [w for w in l_text_clean if w in l_vocab]
51     l_lc = list()
52     lc_clean = re.sub(r'^\w\s', '', str(lc))
53     if lc_clean in l_vocab:
54         l_lc.append(lc_clean)
55         l_lc.extend(l_text_clean)
56         return " ".join(l_lc)
57     else: return " ".join(l_text_clean)
58
59 def doc_to_list(lc, text):

```

```

60     ''' limpa e filtra vocab '''
61     l_text_clean = clean_doc(text, False)
62     l_text_clean = [w for w in l_text_clean if w in l_vocab]
63     l_lc = list()
64     lc_clean = re.sub(r'^\w\s', '', str(lc))
65     if lc_clean in l_vocab:
66         l_lc.append(lc_clean)
67         l_lc.extend(l_text_clean)
68     return l_lc
69 else: return l_text_clean
70
71 def bigram_deep(text_clean):
72     l_text = text_clean.split()
73     l_bi = list(bigrams(l_text))
74     l_bi2 = []
75     for tup in l_bi:
76         l_bi2.append(tup[0] + ' ' + tup[1])
77     l_text.extend(l_bi2)
78     return ','.join(l_text)
79
80 def text2bigram(text):
81     l_text = clean_doc(text, False)
82     l_bi = list(bigrams(l_text))
83     l_bi2 = []
84     for tup in l_bi:
85         l_bi2.append(tup[0] + ' ' + tup[1])
86     return ','.join(l_bi2)
87
88 def drop_dfindex(df, index):
89     return df.drop(index=index)
90
91 def teste_chi2(docs=None, y=None, vectorizer=None, p_value_limit = 0.95):
92     ''' chi-squared test to dermine whether a feature and target are independent -
93         keep only features with p_value_limit = 0.95 '''
94     #vectorizer = TfidfVectorizer(ngram_range=(1,2))
95     #vectorizer.fit(docs)
96     #X = vectorizer.transform(docs)
97     X = vectorizer.fit_transform(docs)
98     X_names = vectorizer.get_feature_names()
99     print('tamanho_vocabulario = ', len(vectorizer.vocabulary_))
100    print(X.shape)
101    df_features = pd.DataFrame()
102    for cat in np.unique(y):
103        chi2_, p = chi2(X, y == cat)
104        df_features = df_features.append(pd.DataFrame({'feature': X_names, 'score'
105df_features = df_features.sort_values(['y', 'score'], ascending=[True, False])
106df_features = df_features[df_features['score'] > p_value_limit]
107print(df_features.sample(10))
108print('vocab_testado:', df_features.feature.nunique())
109return df_features.feature.unique().tolist()
110#return df_features
111
112 def load_doc(filename):
113     ''' # load doc into memory '''
114     # open the file as read only
115     file = open(filename, 'r')
116     # read all text
117     text = file.read()
118     # close the file
119     file.close()
120     return text

```

In [4]:

```
1 path = r'../data/dataset_ML_treino.parq'
2 treino = pd.read_parquet(path, engine='auto')
3 treino.shape
```

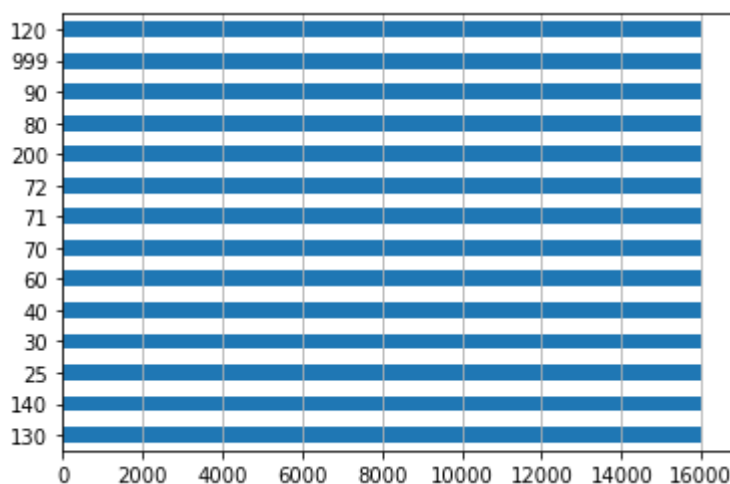
Out[4]:

(224000, 3)

In [5]:

```
1 # atividades econômicas classificadas pelo código do bloco P100 da EFD-Contribuições (S
2 # 999: amostras de atividades não listadas
3 i_classes = treino.y.unique()
4 print(treino.shape)
5 treino.y.value_counts().sort_values().plot(kind='barh', legend=False).grid(axis='x');
```

(224000, 3)



In [6]:

```
1 # retirar registros do dataset de treino
2 #treino.loc[treino.text.str.contains('\x80')]
```

In [7]:

```
1 treino = drop_dffindex(treino, treino.loc[treino.text.str.contains('\x80')].index)
```

In [8]:

```
1 #treino.loc[treino.text.str.contains('\x87')]
```

In [9]:

```
1 treino = drop_dffindex(treino, treino.loc[treino.text.str.contains('\x87')].index)
```


In [10]:

```
1 # montagem do vocabulário unigramas
2 # limpeza do texto: remoção de pontuação, stopwords, caracteres isolados e números
3 # foi mantido o nr correspondente ao código da atividade econômica quando informado
4 vocab = Counter()
5 _ = treino.text.apply(lambda x: add2vocab(x, vocab))
6 _ = treino.cd_lc.apply(lambda x: lc2vocab(x, vocab))
7 print(len(vocab))
8 print(vocab.most_common(10))
```

10359

```
[('servico', 27908), ('servicos', 25137), ('1401', 23711), ('2001', 14670),
('transporte', 12317), ('702', 10867), ('pacote', 10224), ('1601', 10023),
('901', 9743), ('manutencao', 9247)]
```

In [11]:

```
1 # mantém token com uma ocorrência mínima
2 min_occurrence = 2
3 l_vocab = [k for k,c in vocab.items() if c >= min_occurrence]
4 # salva tokens em arquivo texto
5 save_list(l_vocab, r'../best_models/vocab.txt')
6
7 print(len(l_vocab))
```

6350

In [12]:

```
1 # bigramas
2 treino['bigrams'] = treino['text'].apply(lambda x: text2bigram(x))
```

In [13]:

```
1 # montagem do vocabulário bigramas
2 # limpeza do texto: remoção de pontuação, stopwords, caracteres isolados e números
3 vocab2 = Counter()
4 _ = treino.bigrams.apply(lambda x: add2vocab2(x, vocab2))
5 del vocab2['']
6 print(len(vocab2))
7 print(vocab2.most_common(10))
```

21960

```
[('locacao vaga', 7820), ('pacote positron', 7386), ('passageiros dinheiro',
6787), ('gerenciamento residuos', 6297), ('sv portuario', 5318), ('portuario
fiscal', 5318), ('servicos prestados', 5257), ('maodeobra oficina', 5217),
('codigo servico', 4144), ('servico importado', 4143)]
```

In [14]:

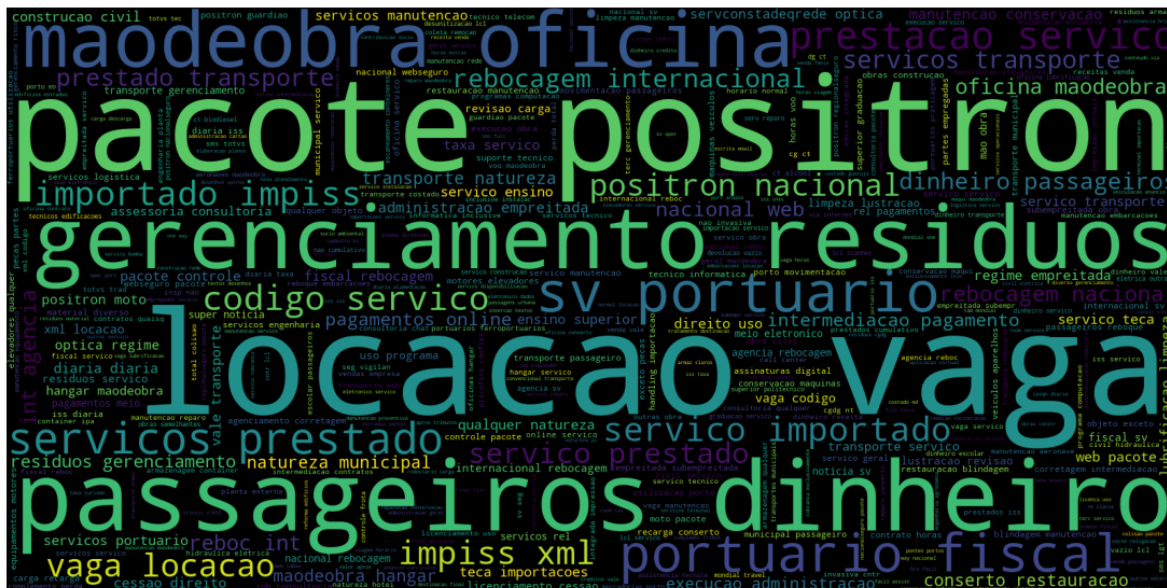
```
1 # unigramas com limpeza do texto
2 treino['text_clean'] = treino[['cd_lc', 'text']].apply(doc_to_line, axis=1)
```

```
1 # bow: bag of words
2 unigrams = treino.text_clean
3 bow = " ".join(unigrams)
```

```
1 # gerando WordCloud
2 wordcloud = WordCloud(background_color="black",
3                        width=1600, height=800, max_words=400, min_word_length=2).generate
```

```
1 # plotando WordCloud
2 fig, ax = plt.subplots(figsize=(20,10))
3 ax.imshow(wordcloud, interpolation='bilinear')
4 ax.set_axis_off()
5
6 plt.imshow(wordcloud);
7 wordcloud.to_file("wordcloud_desc_it.png")
```

```
<wordcloud.wordcloud.WordCloud at 0x7f48a62fb190>
```



In [18]:

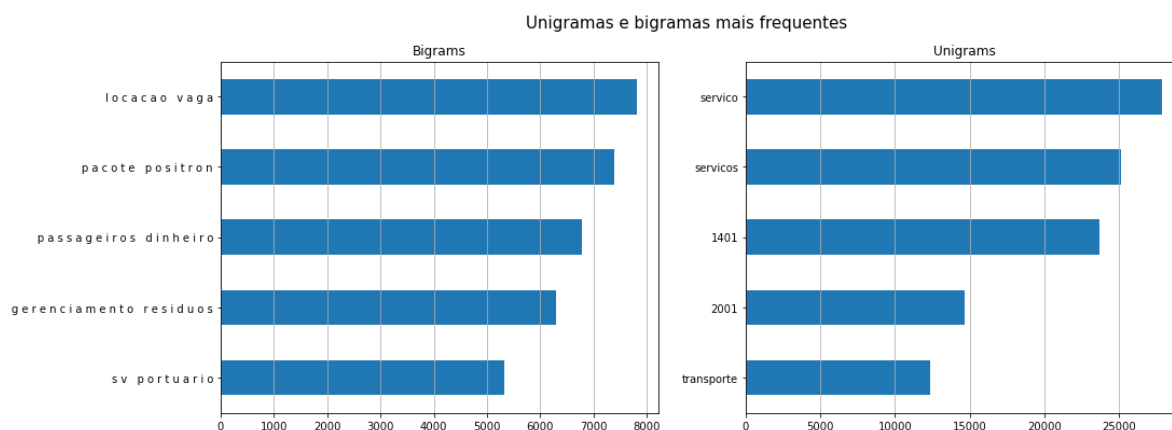
```

1  # unigramas e bigramas mais frequentes
2  top = 5
3  fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(16,6))
4  fig.suptitle('Unigramas e bigramas mais frequentes', fontsize=15)
5  ## unigrams
6  df_uni = pd.DataFrame(vocab.most_common(), columns=['word', 'freq'])
7  df_uni.set_index('word').iloc[:top,:].sort_values(by='freq').plot(kind='barh', title='Unigrams',
8  legend=False).grid(axes=[0])
9  ax[1].set(ylabel=None)
10 ## bigrams
11 df_bi = pd.DataFrame(vocab2.most_common(), columns=['word', 'freq'])
12 df_bi.word = df_bi.word.apply(lambda x: ' '.join(string for string in x))
13 df_bi.set_index('word').iloc[:top,:].sort_values(by='freq').plot(kind="barh", title="Bigrams",
14 legend=False).grid(axes=[0])
15 ax[0].set(ylabel=None)

```

Out[18]:

[Text(0, 0.5, '')]



In [19]:

```

1 def plot_uni_bigramas(y, desc):
2     ''' plota unigramas e bigramas de classe '''
3     # vocab unigramas
4     vocab_y = Counter()
5     _ = treino.loc[treino.y == y].text.apply(lambda x: add2vocab(x, vocab_y))
6     _ = treino.loc[treino.y == y].cd_lc.apply(lambda x: lc2vocab(x, vocab_y))
7     print('Tamanho vocabulario unigramas:', len(vocab_y))
8     print('Unigramas mais comuns:', vocab_y.most_common(10))
9     # vocab bigramas
10    vocab_y2 = Counter()
11    _ = treino.loc[treino.y == y].bigrams.apply(lambda x: add2vocab2(x, vocab_y2))
12    del vocab_y2['']
13    print('Tamanho vocabulario bigramas:', len(vocab_y2))
14    print('Bigramas mais comuns', vocab_y2.most_common(10))
15    fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(16,6))
16    fig.suptitle('Unigramas e bigramas mais frequentes y = ' + str(y) + ' ' + desc , fontweight='bold')
17    ## unigrams
18    df_uni = pd.DataFrame(vocab_y.most_common(), columns=['word', 'freq'])
19    df_uni.set_index('word').iloc[:top,:].sort_values(by='freq').plot(kind='barh', title='Unigrams',
20                                                                    legend=False).grid(True)
21    ax[1].set(ylabel=None)
22    ## bigrams
23    df_bi = pd.DataFrame(vocab_y2.most_common(), columns=['word', 'freq'])
24    df_bi.word = df_bi.word.apply(lambda x: ' '.join(string for string in x))
25    df_bi.set_index('word').iloc[:top,:].sort_values(by='freq').plot(kind="barh", title='Bigrams',
26                                                                    legend=False).grid(True)
27    ax[0].set(ylabel=None)

```

In [20]:

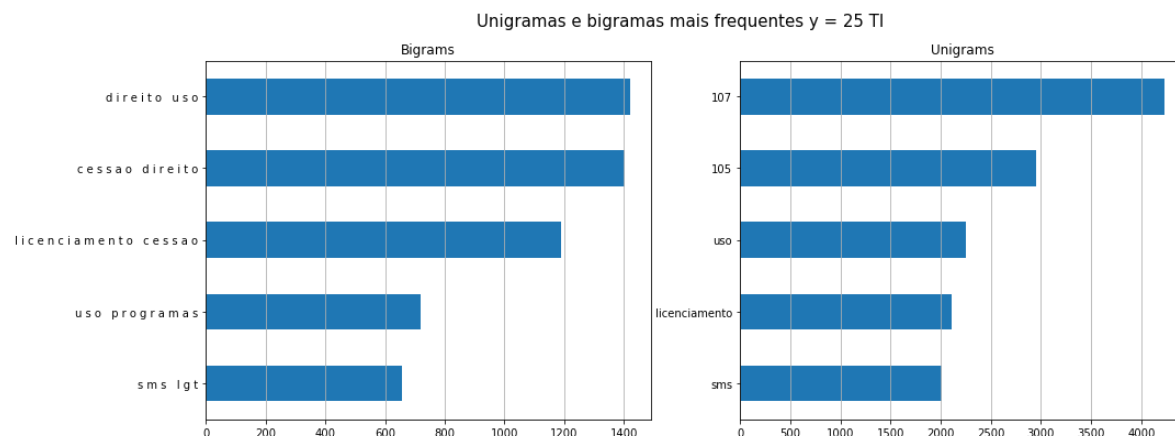
```
1 plot_uni_bigramas(25, 'TI')
```

Tamanho vocabulario unigramas: 2520

Unigramas mais comuns: [('107', 4228), ('105', 2957), ('uso', 2250), ('licenciamento', 2108), ('sms', 2008), ('servicos', 1944), ('103', 1747), ('cessao', 1715), ('101', 1492), ('direito', 1443)]

Tamanho vocabulario bigramas: 4440

Bigramas mais comuns [('direito uso', 1420), ('cessao direito', 1401), ('licenciamento cessao', 1189), ('uso programas', 720), ('sms lgt', 658), ('programas computacao', 651), ('servicos prestados', 634), ('uso programa', 623), ('programa computacao', 588), ('sms totvs', 531)]



In [21]:

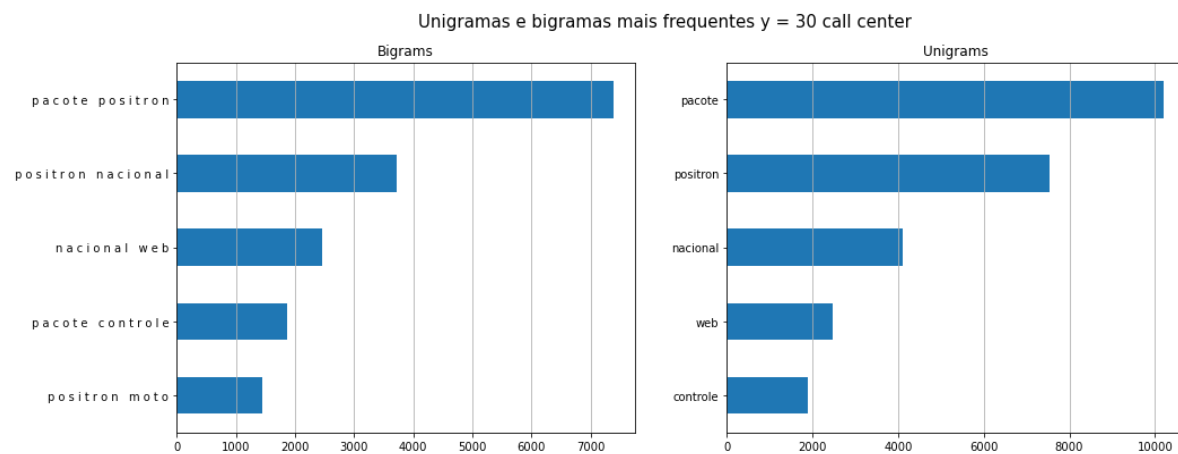
```
1 plot_uni_bigramas(30, 'call center')
```

Tamanho vocabulario unigramas: 820

Unigramas mais comuns: [('pacote', 10196), ('positron', 7536), ('nacional', 4110), ('web', 2484), ('controle', 1895), ('moto', 1456), ('1701', 1274), ('webseguro', 1097), ('guardiao', 903), ('servico', 872)]

Tamanho vocabulario bigramas: 1203

Bigramas mais comuns [('pacote positron', 7386), ('positron nacional', 3726), ('nacional web', 2457), ('pacote controle', 1872), ('positron moto', 1444), ('nacional webseguro', 1097), ('positron guardiao', 903), ('assessoria consultoria', 599), ('positron regionalseguro', 557), ('complemento perda', 526)]



In [22]:

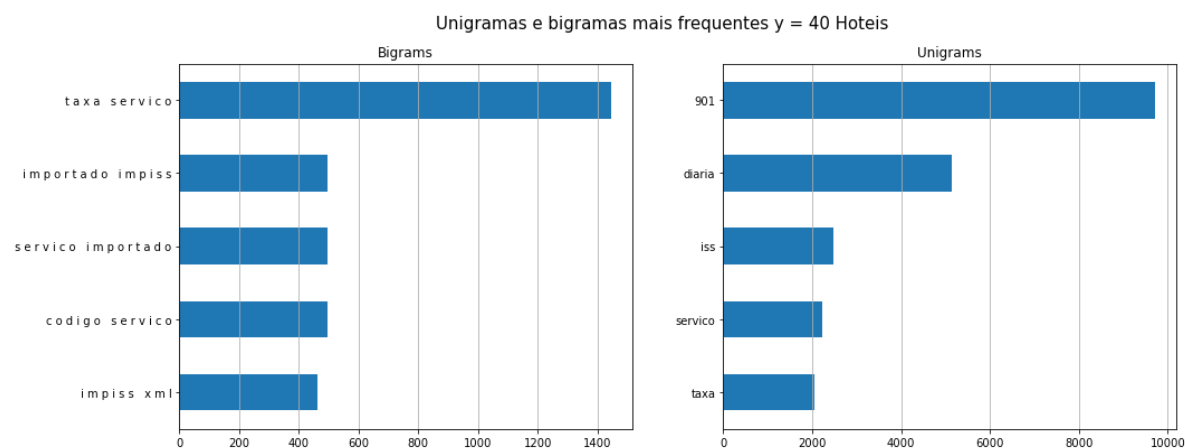
```
1 plot_uni_bigramas(40, 'Hoteis')
```

Tamanho vocabulario unigramas: 694

Unigramas mais comuns: [('901', 9707), ('diaria', 5152), ('iss', 2489), ('servico', 2228), ('taxa', 2054), ('hospedagem', 1015), ('servicos', 719), ('diarias', 698), ('room', 653), ('turismo', 535)]

Tamanho vocabulario bigramas: 830

Bigramas mais comuns [('taxa servico', 1445), ('codigo servico', 496), ('servico importado', 496), ('importado impiss', 496), ('impiss xml', 463), ('hospedagem qualquer', 449), ('qualquer natureza', 449), ('natureza hoteis', 350), ('room tax', 277), ('contribuicao socio', 270)]



In [23]:

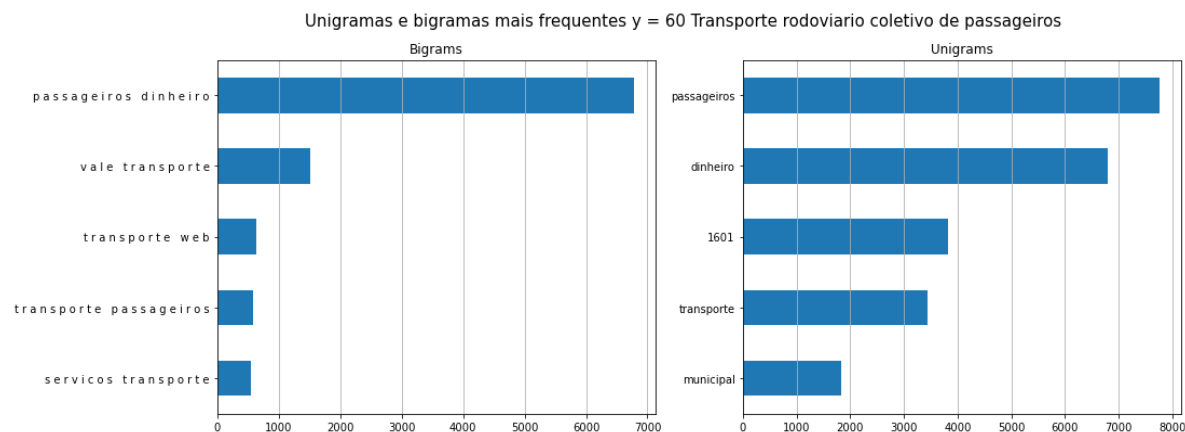
```
1 plot_uni_bigramas(60, 'Transporte rodoviario coletivo de passageiros')
```

Tamanho vocabulario unigramas: 430

Unigramas mais comuns: [('passageiros', 7764), ('dinheiro', 6788), ('1601', 3817), ('transporte', 3433), ('municipal', 1829), ('vale', 1826), ('escolar', 1075), ('servicos', 902), ('venda', 875), ('web', 668)]

Tamanho vocabulario bigramas: 616

Bigramas mais comuns [('passageiros dinheiro', 6787), ('vale transporte', 1522), ('transporte web', 639), ('transporte passageiros', 586), ('servicos transporte', 549), ('passageiros municipal', 455), ('transporte municipal', 451), ('receitas vendas', 449), ('vendas empresa', 440), ('natureza municipal', 409)]



In [24]:

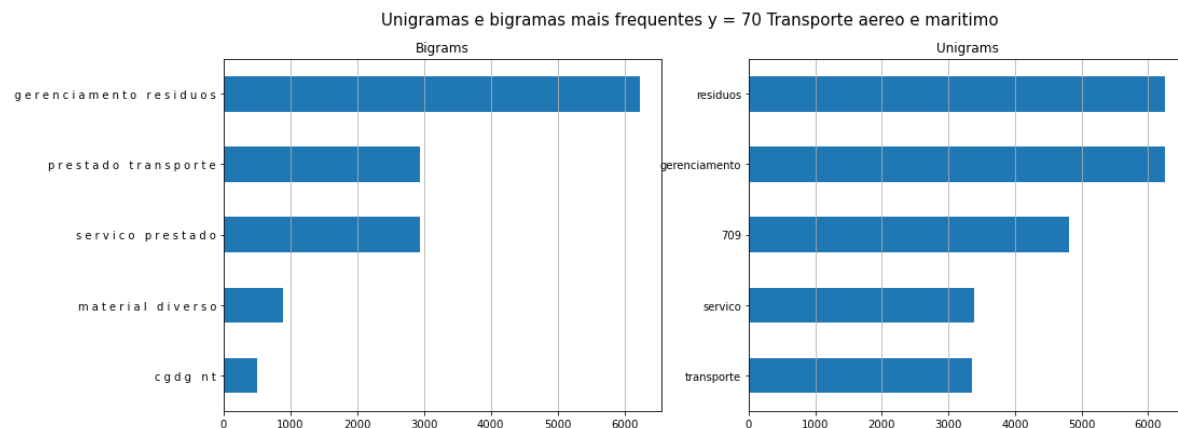
```
1 plot_uni_bigramas(70, 'Transporte aereo e maritimo')
```

Tamanho vocabulario unigramas: 607

Unigramas mais comuns: [('residuos', 6260), ('gerenciamento', 6252), ('709', 4814), ('servico', 3391), ('transporte', 3347), ('1601', 3252), ('prestado', 2940), ('armaz', 1243), ('claros', 1231), ('terc', 1171)]

Tamanho vocabulario bigramas: 970

Bigramas mais comuns [('gerenciamento residuos', 6240), ('servico prestado', 2940), ('prestado transporte', 2940), ('material diverso', 885), ('cgdg nt', 507), ('cg ct', 454), ('ct alcool', 453), ('ct biodiesel', 388), ('dg ct', 366), ('natureza municipal', 314)]



In [25]:

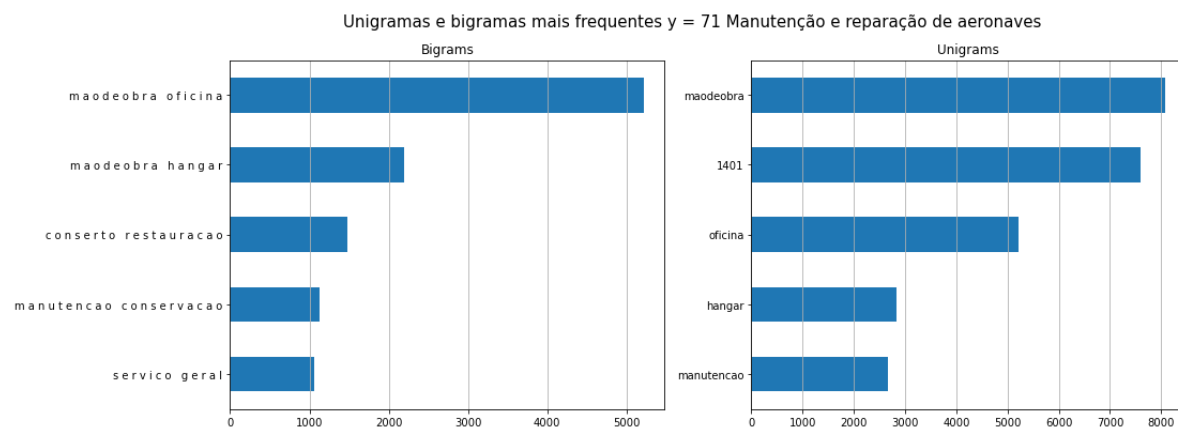
```
1 plot_uni_bigramas(71, 'Manutenção e reparação de aeronaves')
```

Tamanho vocabulário unigramas: 266

Unigramas mais comuns: [('maodeobra', 8075), ('1401', 7606), ('oficina', 5217), ('hangar', 2838), ('manutencao', 2670), ('servico', 1934), ('conserto', 1669), ('restauracao', 1488), ('reparo', 1463), ('servicos', 1444)]

Tamanho vocabulário bigramas: 332

Bigramas mais comuns [('maodeobra oficina', 5217), ('maodeobra hangar', 2193), ('conserto restauracao', 1485), ('manutencao conservacao', 1125), ('servico geral', 1060), ('lubrificacao limpeza', 1039), ('contrato horas', 971), ('horas voo', 971), ('restauracao manutencao', 917), ('lustracao revisao', 659)]



In [26]:

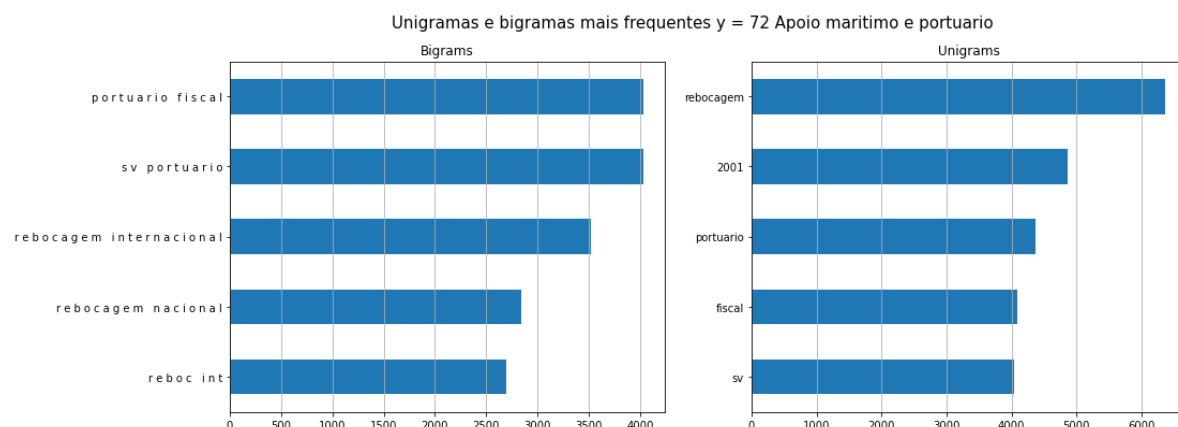
```
1 plot_uni_bigramas(72, 'Apoio marítimo e portuario')
```

Tamanho vocabulário unigramas: 433

Unigramas mais comuns: [('rebocagem', 6368), ('2001', 4862), ('portuario', 4369), ('fiscal', 4093), ('sv', 4037), ('internacional', 3520), ('nacional', 2860), ('agencia', 2721), ('reboc', 2696), ('int', 2695)]

Tamanho vocabulário bigramas: 491

Bigramas mais comuns [('sv portuario', 4036), ('portuario fiscal', 4036), ('rebocagem internacional', 3520), ('rebocagem nacional', 2837), ('reboc int', 2695), ('int agencia', 2695), ('servicos portuarios', 677), ('utilizacao porto', 579), ('portuarios ferroportuarios', 576), ('ferroportuarios utiliza cao', 576)]



In [27]:

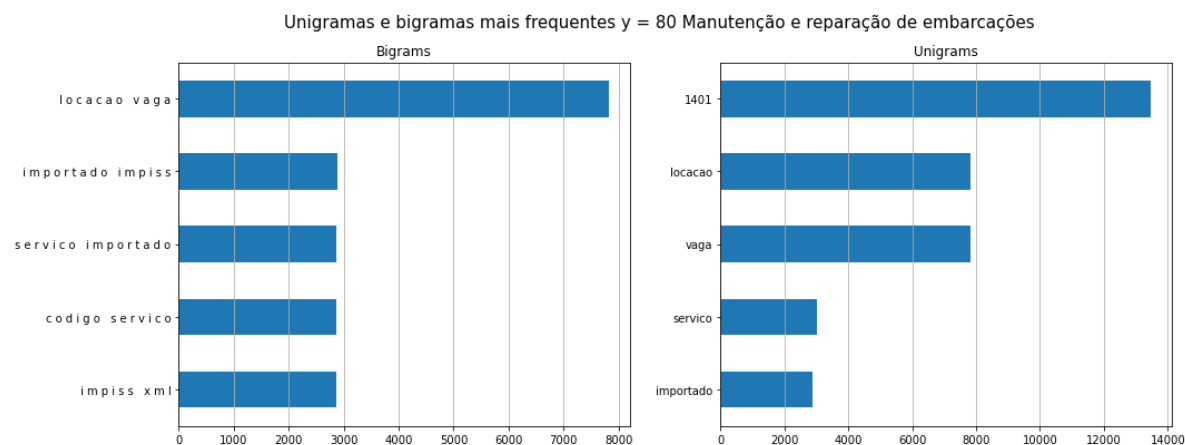
```
1 plot_uni_bigramas(80, 'Manutenção e reparação de embarcações')
```

Tamanho vocabulário unigramas: 386

Unigramas mais comuns: [('1401', 13459), ('locacao', 7824), ('vaga', 7820), ('servico', 3014), ('importado', 2878), ('impiss', 2878), ('codigo', 2875), ('xml', 2864), ('manutencao', 1856), ('equipamentos', 981)]

Tamanho vocabulário bigramas: 451

Bigramas mais comuns [('locacao vaga', 7820), ('importado impiss', 2878), ('codigo servico', 2875), ('servico importado', 2875), ('impiss xml', 2864), ('horario normal', 672), ('manutencao embarcacoes', 609), ('lubrificacao limpeza', 599), ('limpeza lustracao', 599), ('lustracao revisao', 599)]



In [28]:

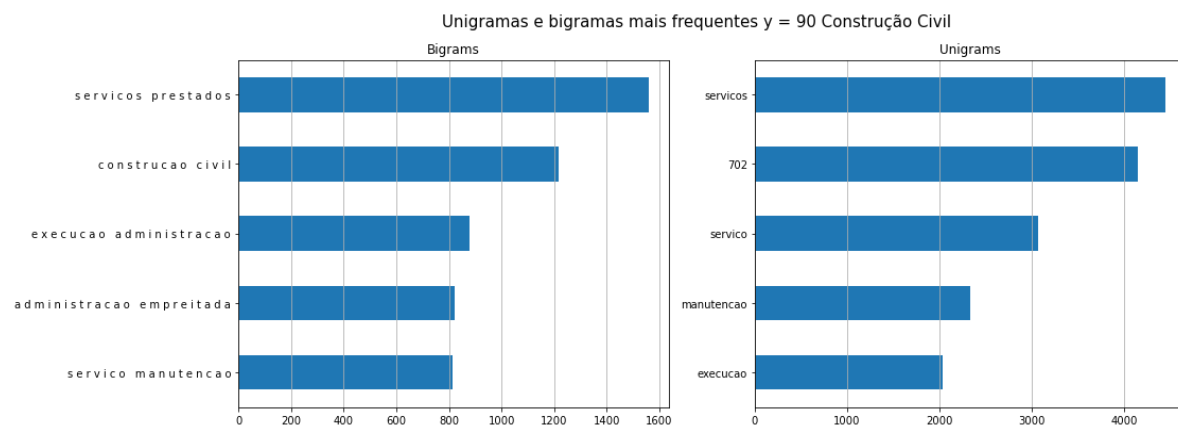
```
1 plot_uni_bigramas(90, 'Construção Civil')
```

Tamanho vocabulário unigramas: 2611

Unigramas mais comuns: [('servicos', 4442), ('702', 4153), ('servico', 3072), ('manutencao', 2328), ('execucao', 2038), ('construcao', 1864), ('civil', 1836), ('empreitada', 1681), ('prestados', 1594), ('obras', 1526)]

Tamanho vocabulário bigramas: 4605

Bigramas mais comuns [('servicos prestados', 1558), ('construcao civil', 1218), ('execucao administracao', 878), ('administracao empreitada', 821), ('servico manutencao', 815), ('prestacao servicos', 759), ('manutencao elevadores', 695), ('mao obra', 690), ('empreitada subempreitada', 606), ('regime empreitada', 591)]



In [29]:

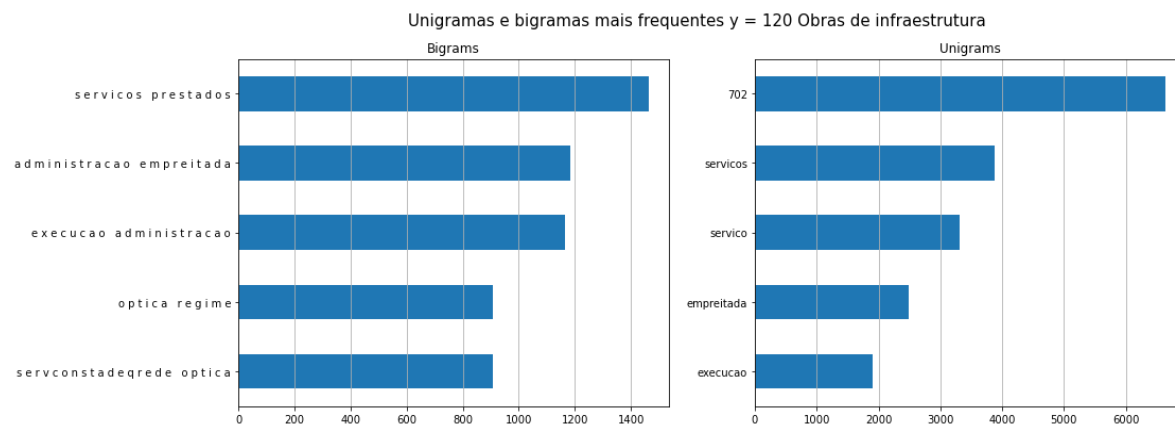
```
1 plot_uni_bigramas(120, 'Obras de infraestrutura')
```

Tamanho vocabulário unigramas: 1760

Unigramas mais comuns: [('702', 6635), ('servicos', 3872), ('servico', 3305), ('empreitada', 2490), ('execucao', 1903), ('prestados', 1553), ('civil', 1342), ('3101', 1243), ('administracao', 1219), ('obra', 1166)]

Tamanho vocabulário bigramas: 2713

Bigramas mais comuns [('servicos prestados', 1462), ('administracao empreitada', 1184), ('execucao administracao', 1164), ('servconstadeqrede optica', 906), ('optica regime', 906), ('regime empreitada', 906), ('obra civil', 704), ('servicos engenharia', 675), ('servico obra', 654), ('empreitada subemp', 647)]



In [30]:

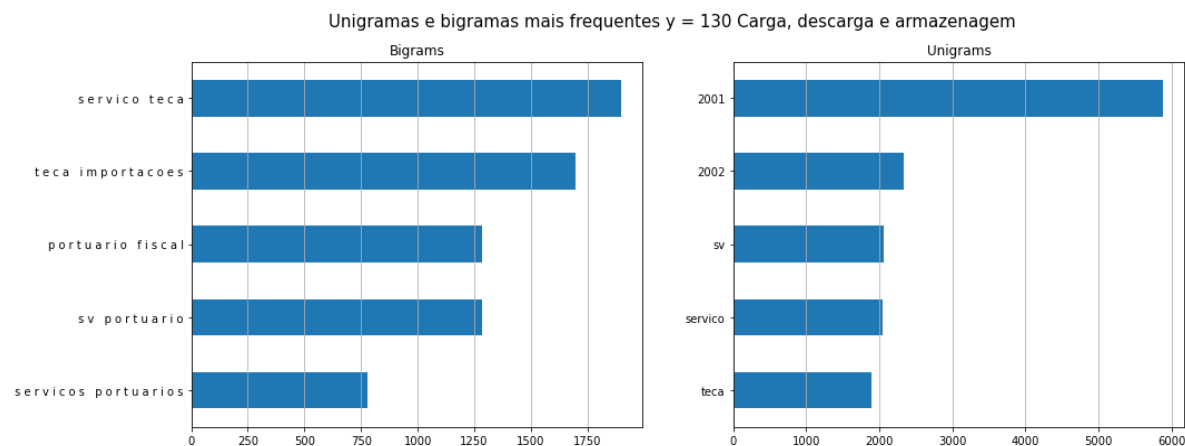
```
1 plot_uni_bigramas(130, 'Carga, descarga e armazenagem')
```

Tamanho vocabulário unigramas: 658

Unigramas mais comuns: [('2001', 5882), ('2002', 2329), ('sv', 2066), ('servico', 2043), ('teca', 1898), ('importacoes', 1694), ('fiscal', 1524), ('armazenagem', 1461), ('lcl', 1417), ('portuario', 1287)]

Tamanho vocabulário bigramas: 935

Bigramas mais comuns [('servico teca', 1898), ('teca importacoes', 1694), ('sv portuario', 1282), ('portuario fiscal', 1282), ('servicos portuarios', 778), ('sv seg', 542), ('seg vigilan', 542), ('portuarios praticagem', 512), ('nao invasiva', 493), ('movimentacao mercadorias', 400)]



In [31]:

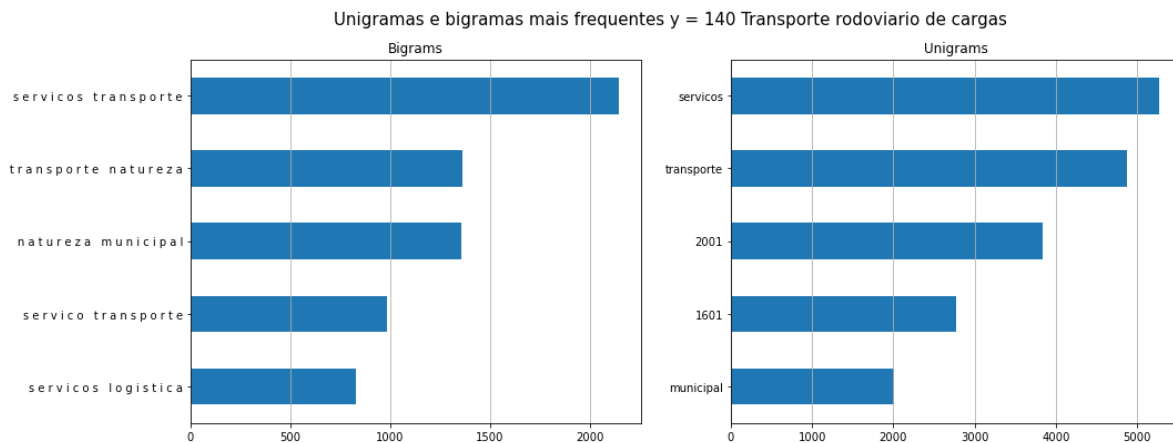
```
1 plot_uni_bigramas(140, 'Transporte rodoviario de cargas')
```

Tamanho vocabulario unigramas: 1207

Unigramas mais comuns: [('servicos', 5277), ('transporte', 4879), ('2001', 3837), ('1601', 2778), ('municipal', 2008), ('servico', 1888), ('natureza', 1559), ('1104', 915), ('armazenagem', 911), ('logistica', 848)]

Tamanho vocabulario bigramas: 1721

Bigramas mais comuns [('servicos transporte', 2149), ('transporte natureza', 1361), ('natureza municipal', 1359), ('servico transporte', 985), ('servicos logistica', 828), ('servicos prestados', 521), ('handling importacao', 516), ('prestacao servicos', 478), ('armazenagem container', 411), ('container ipa', 409)]



In [32]:

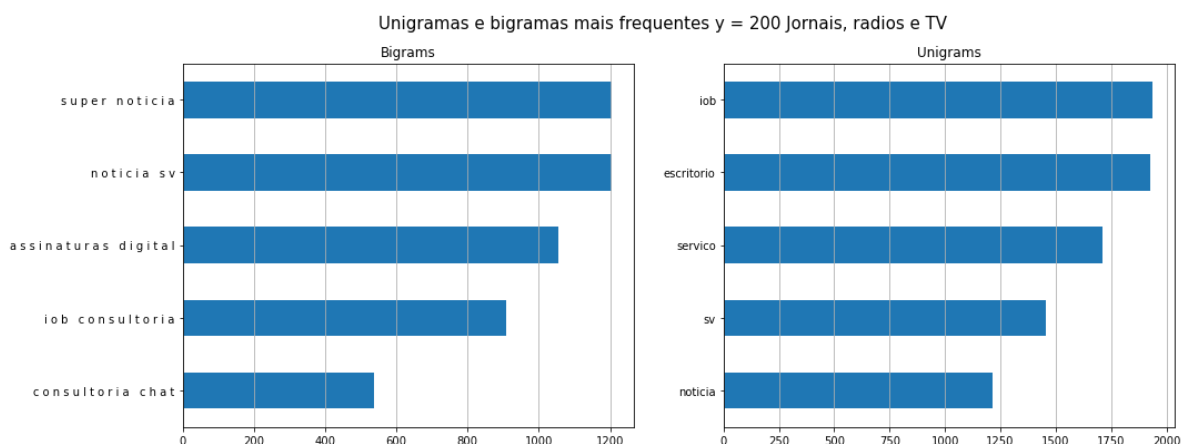
```
1 plot_uni_bigramas(200, 'Jornais, radios e TV')
```

Tamanho vocabulario unigramas: 1217

Unigramas mais comuns: [('iob', 1937), ('escritorio', 1926), ('servico', 1709), ('sv', 1455), ('noticia', 1213), ('super', 1206), ('digital', 1202), ('assinaturas', 1056), ('1706', 1011), ('anuncio', 964)]

Tamanho vocabulario bigramas: 1503

Bigramas mais comuns [('super noticia', 1205), ('noticia sv', 1204), ('assinaturas digital', 1055), ('iob consultoria', 908), ('consultoria chat', 537), ('edicao integrada', 433), ('integrada impressao', 433), ('impressao jornais', 433), ('servico disponibilizacao', 354), ('disponibilizacao temporaria', 354)]



In [33]:

```
1 # call center
2 drop_index = treino.loc[(treino.text_clean.str.contains('call center')) & (treino.y != 30)]
3 treino = treino.drop(index = drop_index)
4 drop_index = treino.loc[(treino.text_clean.str.contains('cobranca')) & (treino.y != 30)]
5 treino = treino.drop(index = drop_index)
```

In [34]:

```
1 # construcao_civil: empreitada, subempreitada, construcao
2 drop_index = treino.loc[(treino.text_clean.str.contains('empreitada')) & (~ treino.y.isin([999]))]
3 treino = treino.drop(index = drop_index)
4
5 drop_index = treino.loc[(treino.text_clean.str.contains('subempreitada')) & (~ treino.y.isin([999]))]
6 treino = treino.drop(index = drop_index)
7
8 drop_index = treino.loc[(treino.text_clean.str.contains('construcao')) & (~ treino.y.isin([999]))]
9 treino = treino.drop(index = drop_index)
```

In [35]:

```
1 # transporte e construcao_civil 999 (nao optantes)
2 # carga, descarga, armazenagem, portuario, praticagem, manutencao, obra, obras, engenharia
3 drop_index = treino.loc[(treino.text_clean.str.contains('carga')) & (treino.y == 999)].index
4 treino = treino.drop(index = drop_index)
5
6 drop_index = treino.loc[(treino.text_clean.str.contains('descarga')) & (treino.y == 999)].index
7 treino = treino.drop(index = drop_index)
8
9 drop_index = treino.loc[(treino.text_clean.str.contains('armazenagem')) & (treino.y == 999)].index
10 treino = treino.drop(index = drop_index)
11
12 drop_index = treino.loc[(treino.text_clean.str.contains('portuario')) & (treino.y == 999)].index
13 treino = treino.drop(index = drop_index)
14
15 drop_index = treino.loc[(treino.text_clean.str.contains('praticagem')) & (treino.y == 999)].index
16 treino = treino.drop(index = drop_index)
17
18 drop_index = treino.loc[(treino.text_clean.str.contains('manutencao')) & (treino.y == 999)].index
19 treino = treino.drop(index = drop_index)
20
21 drop_index = treino.loc[(treino.text_clean.str.contains('obra')) & (treino.y == 999)].index
22 treino = treino.drop(index = drop_index)
23
24 drop_index = treino.loc[(treino.text_clean.str.contains('obras')) & (treino.y == 999)].index
25 treino = treino.drop(index = drop_index)
26
27 drop_index = treino.loc[(treino.text_clean.str.contains('engenharia')) & (treino.y == 999)].index
28 treino = treino.drop(index = drop_index)
29
30 drop_index = treino.loc[(treino.text_clean.str.contains('transporte')) & (treino.y == 999)].index
31 treino = treino.drop(index = drop_index)
32
33 drop_index = treino.loc[(treino.text_clean.str.contains('logistica')) & (treino.y == 999)].index
34 treino = treino.drop(index = drop_index)
```

In [36]:

```

1 # excluindo da analise linhas com text_clean = ''
2 treino.loc[treino.text_clean == '', ['text_clean']] = 'servico'
3
4 # retirar linhas com y<>999 E text_clean: 'prestacao servico','servico','servicos','ser
5 treino = drop_dffindex(treino, treino.loc[(treino.y != 999) & (treino.text_clean == 'pre
6 treino = drop_dffindex(treino, treino.loc[(treino.y != 999) & (treino.text_clean == 'ser
7 treino = drop_dffindex(treino, treino.loc[(treino.y != 999) & (treino.text_clean == 'ser
8 treino = drop_dffindex(treino, treino.loc[(treino.y != 999) & (treino.text_clean == 'ser
9
10 # Linhas (cada linha com 1 documento) para composicao de vetores com unigramas
11 docs_train = treino.text_clean.values # unigramas
12 ytrain = treino.y.values

```

In [37]:

```

1 def plot_distribuiçao_classe(df, x, y):
2     palette = sns.husl_palette(df[y].nunique())
3     fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=(16,6))
4     fig.suptitle(x, fontsize=12)
5     for nr,cat in enumerate(df[y].unique()):
6         sns.kdeplot(data= df[df[y]==cat], x=x,
7                     cbar_kws={"alpha":0.8}, color=palette[nr], shade=True, ax=ax1)
8         sns.histplot(data= df[df[y]==cat], x=x,
9                     bins=10, cbar_kws={"alpha":0.8}, color=palette[nr], ax=ax2)
10    ax1.grid(True)
11    ax2.grid(True)
12    ax1.legend(df[y].unique())
13    ax2.legend(df[y].unique())

```

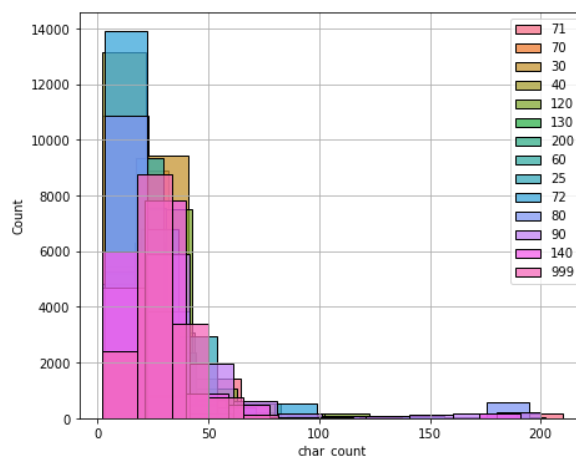
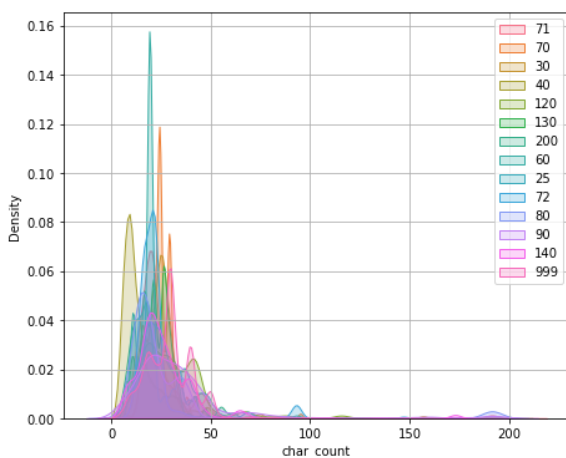
In [38]:

```

1 # distribuicao da soma de caracteres por documento
2 treino['char_count'] = treino["text_clean"].apply(lambda x: sum(len(word) for word in s
3 plot_distribuiçao_classe(treino, 'char_count', 'y')

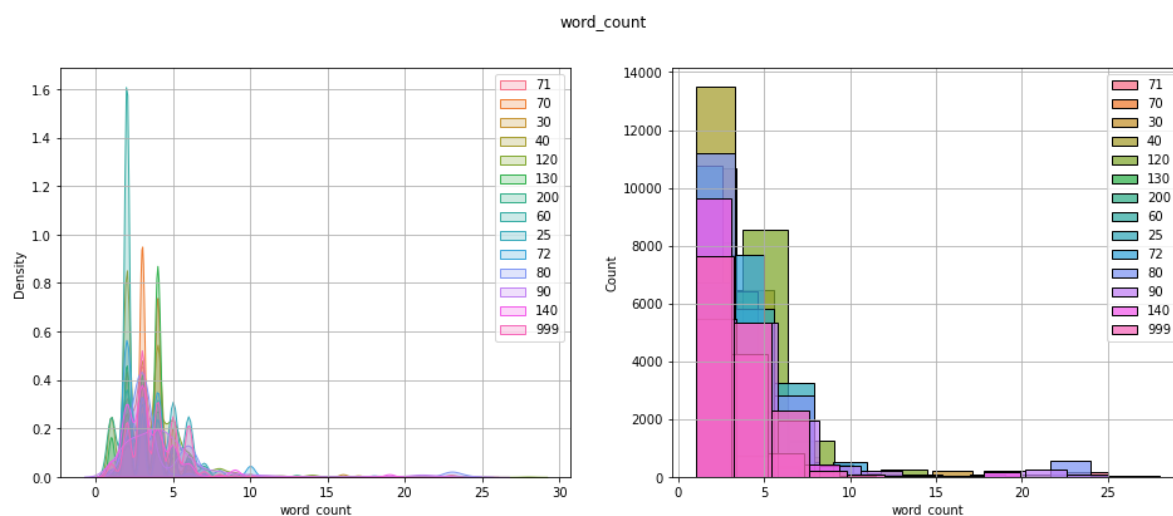
```

char_count



In [39]:

```
1 # distribuicao da contagem de palavras por documento
2 treino['word_count'] = treino["text_clean"].apply(lambda x: len(str(x).split(" "))) #
3 plot_distribuicao_classe(treino, 'word_count', 'y')
```



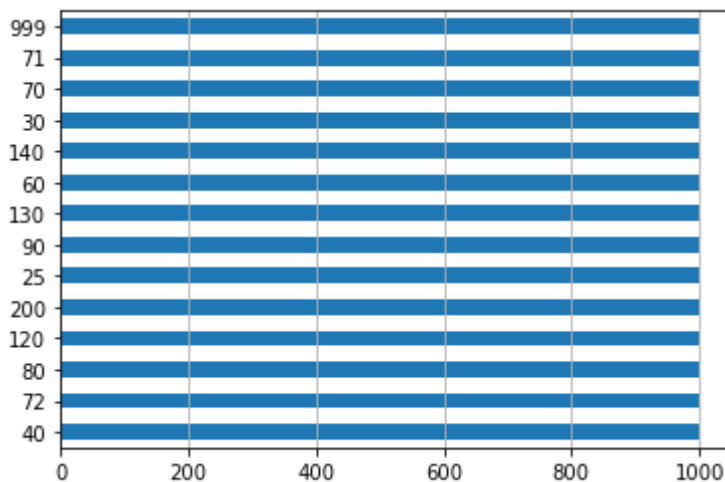
In [40]:

```

1 # DataSet de validação
2 path = r'../data/dataset_ML_val.parq'
3 val = pd.read_parquet(path, engine='auto')
4
5 # unigramas com limpeza do texto no DataSet de validação
6 val['text_clean'] = val[['cd_lc', 'text']].apply(doc_to_line, axis=1)
7 val.loc[val.text_clean == '', ['text_clean']] = 'servico' # classe 999 quando l
8
9 docs_val = val.text_clean.values
10 yval = val.y.values
11
12 print(val.shape)
13 val.y.value_counts().sort_values().plot(kind='barh', legend=False).grid(axis='x');

```

(14000, 4)



Comparando resultados na vetorização dos documentos:

- binary = palavras são marcadas como presentes (1) ou ausentes (0)
- count = contagem de ocorrências de cada palavra como um número inteiro
- freq = frequência de cada palavra no documento
- tfidf = frequência das palavras com penalização das palavras mais comuns (frequência + grau de raridade)
- hash = frequência das palavras com limitação das features (colunas) na matriz de vetores

In [41]:

```

1 # unigramas
2 binary_vectorizer = CountVectorizer(binary=True)
3 count_vectorizer = CountVectorizer()
4 freq_vectorizer = TfidfVectorizer(use_idf=False)
5 tfidf_vectorizer = TfidfVectorizer()
6 hashing_vectorizer = HashingVectorizer(n_features = 2 ** 11)
7
8 # bigramas
9 binary_vectorizer2 = CountVectorizer(binary=True, ngram_range=(1,2))
10 count_vectorizer2 = CountVectorizer(ngram_range=(1,2))
11 freq_vectorizer2 = TfidfVectorizer(use_idf=False, ngram_range=(1,2))
12 tfidf_vectorizer2 = TfidfVectorizer(ngram_range=(1,2))
13 hashing_vectorizer2 = HashingVectorizer(n_features = 2 ** 11, ngram_range=(1,2))

```

In [42]:

```

1 # teste chi2 para seleção de features relevantes com unigramas
2 Xtrain_names_binary = teste_chi2(docs=docs_train, y=ytrain, vectorizer = binary_vectorizer)
3 binary_vectorizer_chi2 = CountVectorizer(binary=True, vocabulary = Xtrain_names_binary)

```

```

tamanho_vocabulario = 6315
(215298, 6315)

```

	feature	score	y
2869	graduacao	1.0000	80
2612	flat	0.9997	70
808	backbone	0.9873	130
5720	talao	0.9874	40
6047	unisys	1.0000	90
4565	portuario	1.0000	200
1580	contabilidade	0.9874	140
4913	reaventrega	1.0000	120
2220	encomperiodo	0.9908	80
710	assist	1.0000	40

```

[10 rows x 3 columns]
vocab_testado: 6301

```

In [43]:

```

1 Xtrain_names_count = teste_chi2(docs=docs_train, y=ytrain, vectorizer = count_vectorizer)
2 count_vectorizer_chi2 = CountVectorizer(vocabulary = Xtrain_names_count)

```

```

tamanho_vocabulario = 6315
(215298, 6315)

```

	feature	score	y
5168	rio	0.9713	140
4203	outras	1.0000	71
1466	concreta	1.0000	90
1288	civileletrica	1.0000	90
5834	terminal	1.0000	140
1079	campos	0.9599	999
5609	stc	1.0000	120
653	armazenagem	1.0000	200
5781	tef	0.9820	72
4488	plataforma	1.0000	70

```

[10 rows x 3 columns]
vocab_testado: 6301

```

In [44]:

```
1 Xtrain_names_freq = teste_chi2(docs=docs_train, y=ytrain, vectorizer = freq_vectorizer)
2 freq_vectorizer_chi2 = TfidfVectorizer(use_idf=False, vocabulary = Xtrain_names_freq)
```

```
tamanho_vocabulario = 6315
(215298, 6315)
```

	feature	score	y
1128	carregamento	0.9996	200
3075	imp	0.9809	200
4354	ped	1.0000	90
5745	taxdoacao	0.9597	71
1069	cambagem	1.0000	999
4422	petrolina	1.0000	999
1129	carregamentosaida	1.0000	80
4225	pacote	1.0000	200
1978	distribuicao	0.9998	72
962	bookingcom	1.0000	40

```
[10 rows x 3 columns]
vocab_testado: 6175
```

In [45]:

```
1 Xtrain_names_tfidf = teste_chi2(docs=docs_train, y=ytrain, vectorizer = tfidf_vectorizer)
2 tfidf_vectorizer_chi2 = TfidfVectorizer(use_idf=True, vocabulary = Xtrain_names_tfidf)
```

```
tamanho_vocabulario = 6315
(215298, 6315)
```

	feature	score	y
5000	refrigerada	1.0000	130
428	agenciaepcm	1.0000	200
4439	picking	1.0000	130
4013	nh	0.9556	60
1724	cursogenerico	1.0000	200
1038	cabotagem	0.9925	71
3444	lgt	1.0000	30
326	acompanhamento	0.9981	60
5138	retido	1.0000	120
3658	maritimo	0.9885	120

```
[10 rows x 3 columns]
vocab_testado: 6247
```


In [46]:

```

1 # teste chi2 para seleção de features relevantes com unigramas e bigramas
2 Xtrain_names_binary2 = teste_chi2(docs=docs_train, y=ytrain, vectorizer = binary_vectorizer)
3 binary_vectorizer2_chi2 = CountVectorizer(binary=True, vocabulary = Xtrain_names_binary2.get_vocab())

```

```

tamanho_vocabulario = 26580
(215298, 26580)

```

	feature	score	y
21699	sapata freio	1.0000	90
22267	servico balanceamento	0.9997	999
11147	fed	0.9878	130
22632	servicos	1.0000	140
14657	locsupgespam compras	0.9997	25
20407	redacao	0.9929	25
16959	obra ferro	1.0000	90
3493	aeronaves parte	0.9996	130
19804	psubst roldana	0.9997	80
18225	permanente	0.9694	130

```

[10 rows x 3 columns]
vocab_testado: 26553

```

In [47]:

```

1 Xtrain_names_count2 = teste_chi2(docs=docs_train, y=ytrain, vectorizer = count_vectorizer)
2 count_vectorizer2_chi2 = CountVectorizer(vocabulary = Xtrain_names_count2.get_vocab())

```

```

tamanho_vocabulario = 26580
(215298, 26580)

```

	feature	score	y
17070	obras frs	1.0000	90
24332	taxa servico	1.0000	25
14759	lubrificacao	1.0000	71
15179	manutencao sistemas	0.9776	999
22052	serv emendas	0.9767	71
3919	analise area	0.9997	999
1145	1401 manut	0.9756	71
16524	nao perigosos	1.0000	70
16573	natureza servicos	0.9995	70
13185	infraestrutura administrativa	0.9851	140

```

[10 rows x 3 columns]
vocab_testado: 26553

```

In [48]:

```
1 Xtrain_names_freq2 = teste_chi2(docs=docs_train, y=ytrain, vectorizer = freq_vectorizer)
2 freq_vectorizer2_chi2 = TfidfVectorizer(use_idf=False, vocabulary = Xtrain_names_freq2)
```

```
tamanho_vocabulario = 26580
(215298, 26580)
```

	feature	score	y
17260	oper	0.9893	25
1929	2002 ataero	0.9998	90
26181	versoes	0.9953	25
2687	705 servicos	0.9505	71
6073	cartao	0.9984	90
4944	autorizacao	1.0000	200
4790	ato cooperado	1.0000	999
25373	treinamento orientacao	1.0000	999
13526	intermediacao	1.0000	130
9418	ecomp aracruz	1.0000	70

```
[10 rows x 3 columns]
vocab_testado: 24243
```

In [49]:

```
1 Xtrain_names_tfidf2 = teste_chi2(docs=docs_train, y=ytrain, vectorizer = tfidf_vectorizer)
2 tfidf_vectorizer2_chi2 = TfidfVectorizer(use_idf=True, vocabulary = Xtrain_names_tfidf2)
```

```
tamanho_vocabulario = 26580
(215298, 26580)
```

	feature	score	y
25906	vale transportes	0.9999	60
24399	tec gaspaj	0.9993	70
19313	prestados703	0.9999	90
12691	hotel	0.9959	72
566	106 programa	1.0000	999
25878	vaga	1.0000	60
7430	congeneres	0.9997	200
15623	mensagem natal	0.9980	200
14538	lmi6114 rec	0.9996	72
20452	redes districao	0.9896	90

```
[10 rows x 3 columns]
vocab_testado: 25240
```

In [50]:

```

1 results_ml = pd.DataFrame(columns = ['acc', 'precision_macro', 'recall_macro', 'f1_macro'],
2 def s_results(docs_train, ytrain, vectorizer=None, model=None, svd=None, docs_val=None,
3     if svd:
4         pipe = Pipeline(steps=[('vectorizer', vectorizer),
5                                 ('svd', TruncatedSVD(n_components=100, n_iter=10, random
6                                 ('model', model))])
7     else:
8         pipe = Pipeline(steps=[('vectorizer', vectorizer),
9                                 ('model', model)])
10    pipe.fit(docs_train, ytrain)
11    yv_pred = pipe.predict(docs_val)
12    acc = accuracy_score(yval, yv_pred)
13    p_m = precision_score(yval, yv_pred, average='macro')
14    r_m = recall_score(yval, yv_pred, average='macro')
15    f1_m = f1_score(yval, yv_pred, average='macro')
16    mcc = matthews_corrcoef(yval, yv_pred)
17    return pd.Series(data=[acc, p_m, r_m, f1_m, mcc], index=['acc', 'precision_macro', 'r

```

ML teste de algoritmos

legenda utilizada:

- NB: Naive Bayes MultinomialNB
- RF: Random Forest
- LG: Regressão Logística
- SVM: Suport Vector Machines
- uni: unigramas
- bi: bigramas
- chi2: redução de features pelo método chi_quadrado
- SVD: redução de features pelo método SVD

Naive Bayes MultinomialNB

- parâmetro padrão: alpha = 1.0

In [51]:

```

1 # binary
2 results_ml.loc['NB_binary_uni'] = s_results(docs_train, ytrain, vectorizer= binary_vect
3                                     model= MultinomialNB(), docs_val=docs_val, yval=yv
4
5 results_ml.loc['NB_binary_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= binary_vect
6                                     model= MultinomialNB(), docs_val=docs_val, yval=yv
7
8 results_ml.loc['NB_binary_bi'] = s_results(docs_train, ytrain, vectorizer= binary_vect
9                                     model= MultinomialNB(), docs_val=docs_val, yval=yv
10
11 results_ml.loc['NB_binary_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= binary_vect
12                                     model= MultinomialNB(), docs_val=docs_val, yval=yv

```

In [52]:

```

1 # count
2 results_ml.loc['NB_count_uni'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
3                                           model= MultinomialNB(), docs_val=docs_val, yval=yval)
4
5 results_ml.loc['NB_count_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
6                                                model= MultinomialNB(), docs_val=docs_val, yval=yval)
7
8 results_ml.loc['NB_count_bi'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
9                                           model= MultinomialNB(), docs_val=docs_val, yval=yval)
10
11 results_ml.loc['NB_count_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
12                                                model= MultinomialNB(), docs_val=docs_val, yval=yval)

```

In [53]:

```

1 # freq
2 results_ml.loc['NB_freq_uni'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
3                                           model= MultinomialNB(), docs_val=docs_val, yval=yval)
4
5 results_ml.loc['NB_freq_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
6                                                model= MultinomialNB(), docs_val=docs_val, yval=yval)
7
8 results_ml.loc['NB_freq_bi'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
9                                           model= MultinomialNB(), docs_val=docs_val, yval=yval)
10
11 results_ml.loc['NB_freq_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
12                                                model= MultinomialNB(), docs_val=docs_val, yval=yval)

```

In [54]:

```

1 # tfidf
2 results_ml.loc['NB_tfidf_uni'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
3                                           model= MultinomialNB(), docs_val=docs_val, yval=yval)
4
5 results_ml.loc['NB_tfidf_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
6                                                model= MultinomialNB(), docs_val=docs_val, yval=yval)
7
8 results_ml.loc['NB_tfidf_bi'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
9                                           model= MultinomialNB(), docs_val=docs_val, yval=yval)
10
11 results_ml.loc['NB_tfidf_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
12                                                model= MultinomialNB(), docs_val=docs_val, yval=yval)

```

Regressão Logística

In [55]:

```
1 # binary
2 results_ml.loc['LG_binary_uni'] = s_results(docs_train, ytrain, vectorizer= binary_vect
3                                     model= LogisticRegression(solver='liblinear', rand
4
5 results_ml.loc['LG_binary_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= binary_
6                                     svd=True,
7                                     model= LogisticRegression(solver='liblinear', rand
8
9 results_ml.loc['LG_binary_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= binary
10                                     model= LogisticRegression(solver='liblinear', rand
11
12 results_ml.loc['LG_binary_bi'] = s_results(docs_train, ytrain, vectorizer= binary_vecto
13                                     model= LogisticRegression(solver='liblinear', rand
14
15 results_ml.loc['LG_binary_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= binary_v
16                                     svd=True,
17                                     model= LogisticRegression(solver='liblinear', rand
18
19 results_ml.loc['LG_binary_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= binary_
20                                     model= LogisticRegression(solver='liblinear', rand
```

In [56]:

```
1 # count
2 results_ml.loc['LG_count_uni'] = s_results(docs_train, ytrain, vectorizer= count_vector
3                                     model= LogisticRegression(solver='liblinear', rand
4
5 results_ml.loc['LG_count_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= count_ve
6                                     svd=True,
7                                     model= LogisticRegression(solver='liblinear', rand
8
9 results_ml.loc['LG_count_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= count_v
10                                     model= LogisticRegression(solver='liblinear', rand
11
12 results_ml.loc['LG_count_bi'] = s_results(docs_train, ytrain, vectorizer= count_vectori
13                                     model= LogisticRegression(solver='liblinear', rand
14
15 results_ml.loc['LG_count_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= count_vec
16                                     svd=True,
17                                     model= LogisticRegression(solver='liblinear', rand
18
19 results_ml.loc['LG_count_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= count_ve
20                                     model= LogisticRegression(solver='liblinear', rand
```

In [57]:

```
1 # freq
2 results_ml.loc['LG_freq_uni'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
3                                         model= LogisticRegression(solver='liblinear', random_state=42))
4
5 results_ml.loc['LG_freq_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
6                                             svd=True,
7                                             model= LogisticRegression(solver='liblinear', random_state=42))
8
9 results_ml.loc['LG_freq_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
10                                              model= LogisticRegression(solver='liblinear', random_state=42))
11
12 results_ml.loc['LG_freq_bi'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
13                                         model= LogisticRegression(solver='liblinear', random_state=42))
14
15 results_ml.loc['LG_freq_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
16                                             svd=True,
17                                             model= LogisticRegression(solver='liblinear', random_state=42))
18
19 results_ml.loc['LG_freq_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
20                                              model= LogisticRegression(solver='liblinear', random_state=42))
```

In [58]:

```
1 # tfidf
2 results_ml.loc['LG_tfidf_uni'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
3                                         model= LogisticRegression(solver='liblinear', random_state=42))
4
5 results_ml.loc['LG_tfidf_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
6                                             svd=True,
7                                             model= LogisticRegression(solver='liblinear', random_state=42))
8
9 results_ml.loc['LG_tfidf_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
10                                              model= LogisticRegression(solver='liblinear', random_state=42))
11
12 results_ml.loc['LG_tfidf_bi'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
13                                         model= LogisticRegression(solver='liblinear', random_state=42))
14
15 results_ml.loc['LG_tfidf_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
16                                             svd=True,
17                                             model= LogisticRegression(solver='liblinear', random_state=42))
18
19 results_ml.loc['LG_tfidf_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
20                                              model= LogisticRegression(solver='liblinear', random_state=42))
```

In [59]:

```

1 # hash
2 results_ml.loc['LG_hash_uni'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
3                                         model= LogisticRegression(solver='liblinear', random_state=0))
4
5 results_ml.loc['LG_hash_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
6                                             svd=True,
7                                             model= LogisticRegression(solver='liblinear', random_state=0))
8
9 results_ml.loc['LG_hash_bi'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
10                                         model= LogisticRegression(solver='liblinear', random_state=0))
11
12 results_ml.loc['LG_hash_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
13                                             svd=True,
14                                             model= LogisticRegression(solver='liblinear', random_state=0))

```

SVM: Support Vector Machines

parâmetro padrão: kernel rbf, C = 1.0, gamma = scale

In [60]:

```

1 # binary
2 results_ml.loc['SVM_binary_uni'] = s_results(docs_train, ytrain, vectorizer= binary_vectorizer,
3                                             model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
4
5 results_ml.loc['SVM_binary_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= binary_vectorizer,
6                                             svd=True,
7                                             model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
8
9 results_ml.loc['SVM_binary_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= binary_vectorizer,
10                                                  model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
11
12 results_ml.loc['SVM_binary_bi'] = s_results(docs_train, ytrain, vectorizer= binary_vectorizer,
13                                             model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
14
15 results_ml.loc['SVM_binary_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= binary_vectorizer,
16                                             svd=True,
17                                             model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
18
19 results_ml.loc['SVM_binary_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= binary_vectorizer,
20                                                  model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)

```

In [61]:

```

1 # count
2 results_ml.loc['SVM_count_uni'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
3                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
4
5 results_ml.loc['SVM_count_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
6                                           svd=True,
7                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
8
9 results_ml.loc['SVM_count_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
10                                                  model= SVC(random_state=0), docs_val=docs_val, yval=yval)
11
12 results_ml.loc['SVM_count_bi'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
13                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
14
15 results_ml.loc['SVM_count_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
16                                           svd=True,
17                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
18
19 results_ml.loc['SVM_count_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= count_vectorizer,
20                                                  model= SVC(random_state=0), docs_val=docs_val, yval=yval)

```

In [62]:

```

1 # freq
2 results_ml.loc['SVM_freq_uni'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
3                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
4
5 results_ml.loc['SVM_freq_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
6                                           svd=True,
7                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
8
9 results_ml.loc['SVM_freq_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
10                                                  model= SVC(random_state=0), docs_val=docs_val, yval=yval)
11
12 results_ml.loc['SVM_freq_bi'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
13                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
14
15 results_ml.loc['SVM_freq_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
16                                           svd=True,
17                                           model= SVC(random_state=0), docs_val=docs_val, yval=yval)
18
19 results_ml.loc['SVM_freq_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
20                                                  model= SVC(random_state=0), docs_val=docs_val, yval=yval)

```


In [63]:

```

1 # tfidf
2 results_ml.loc['SVM_tfidf_uni'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
3                                           model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
4
5 results_ml.loc['SVM_tfidf_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
6                                               svd=True,
7                                               model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
8
9 results_ml.loc['SVM_tfidf_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
10                                                model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
11
12 results_ml.loc['SVM_tfidf_bi'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
13                                           model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
14
15 results_ml.loc['SVM_tfidf_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
16                                               svd=True,
17                                               model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
18
19 results_ml.loc['SVM_tfidf_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
20                                                model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)

```

In [64]:

```

1 # hash
2 results_ml.loc['SVM_hash_uni'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
3                                           model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
4
5 results_ml.loc['SVM_hash_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
6                                               svd=True,
7                                               model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
8
9 results_ml.loc['SVM_hash_bi'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
10                                           model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)
11
12 results_ml.loc['SVM_hash_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
13                                               svd=True,
14                                               model= SVC(random_state=0), docs_val=docs_val, y_val=y_val)

```

Random Forest

In [65]:

```
1 # binary
2 results_ml.loc['RF_binary_uni'] = s_results(docs_train, ytrain, vectorizer= binary_vect
3                                     model= RandomForestClassifier(random_state=0), doc
4
5 results_ml.loc['RF_binary_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= binary_
6                                     svd=True,
7                                     model= RandomForestClassifier(random_state=0), doc
8
9 results_ml.loc['RF_binary_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= binary
10                                     model= RandomForestClassifier(random_state=0), doc
11
12 results_ml.loc['RF_binary_bi'] = s_results(docs_train, ytrain, vectorizer= binary_vect
13                                     model= RandomForestClassifier(random_state=0), doc
14
15 results_ml.loc['RF_binary_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= binary_v
16                                     svd=True,
17                                     model= RandomForestClassifier(random_state=0), doc
18
19 results_ml.loc['RF_binary_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= binary_
20                                     model= RandomForestClassifier(random_state=0), doc
```

In [66]:

```
1 # count
2 results_ml.loc['RF_count_uni'] = s_results(docs_train, ytrain, vectorizer= count_vecto
3                                     model= RandomForestClassifier(random_state=0), doc
4
5 results_ml.loc['RF_count_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= count_ve
6                                     svd=True,
7                                     model= RandomForestClassifier(random_state=0), doc
8
9 results_ml.loc['RF_count_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= count_v
10                                     model= RandomForestClassifier(random_state=0), doc
11
12 results_ml.loc['RF_count_bi'] = s_results(docs_train, ytrain, vectorizer= count_vectori
13                                     model= RandomForestClassifier(random_state=0), doc
14
15 results_ml.loc['RF_count_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= count_vec
16                                     svd=True,
17                                     model= RandomForestClassifier(random_state=0), doc
18
19 results_ml.loc['RF_count_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= count_ve
20                                     model= RandomForestClassifier(random_state=0), doc
```

In [67]:

```
1 # freq
2 results_ml.loc['RF_freq_uni'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
3                                         model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
4
5 results_ml.loc['RF_freq_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
6                                              svd=True,
7                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
8
9 results_ml.loc['RF_freq_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
10                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
11
12 results_ml.loc['RF_freq_bi'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
13                                         model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
14
15 results_ml.loc['RF_freq_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
16                                              svd=True,
17                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
18
19 results_ml.loc['RF_freq_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= freq_vectorizer,
20                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
```

In [68]:

```
1 # tfidf
2 results_ml.loc['RF_tfidf_uni'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
3                                         model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
4
5 results_ml.loc['RF_tfidf_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
6                                              svd=True,
7                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
8
9 results_ml.loc['RF_tfidf_uni_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
10                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
11
12 results_ml.loc['RF_tfidf_bi'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
13                                         model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
14
15 results_ml.loc['RF_tfidf_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
16                                              svd=True,
17                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
18
19 results_ml.loc['RF_tfidf_bi_chi2'] = s_results(docs_train, ytrain, vectorizer= tfidf_vectorizer,
20                                              model= RandomForestClassifier(random_state=0), docs_train= docs_train, y_train= ytrain)
```

In [69]:

```
1 # hash
2 results_ml.loc['RF_hash_uni'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
3                                         model= RandomForestClassifier(random_state=0), doc
4
5 results_ml.loc['RF_hash_uni_SVD'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
6                                              svd=True,
7                                              model= RandomForestClassifier(random_state=0), doc
8
9 results_ml.loc['RF_hash_bi'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
10                                         model= RandomForestClassifier(random_state=0), doc
11
12 results_ml.loc['RF_hash_bi_SVD'] = s_results(docs_train, ytrain, vectorizer= hashing_vectorizer,
13                                              svd=True,
14                                              model= RandomForestClassifier(random_state=0), doc
```

modelo ML com melhor resultado

In [70]:

```
1 results_ml.sort_values(by=['mcc'], ascending=False).head(20)
```

Out[70]:

	acc	precision_macro	recall_macro	f1_macro	mcc
LG_binary_bi	0.8841	0.8903	0.8841	0.8850	0.8757
LG_count_bi	0.8837	0.8899	0.8837	0.8846	0.8752
SVM_tfidf_bi	0.8822	0.8907	0.8822	0.8835	0.8737
SVM_tfidf_uni_chi2	0.8816	0.8896	0.8816	0.8829	0.8731
SVM_freq_bi	0.8816	0.8900	0.8816	0.8828	0.8731
SVM_tfidf_bi_chi2	0.8816	0.8895	0.8816	0.8828	0.8730
SVM_tfidf_uni	0.8814	0.8894	0.8814	0.8826	0.8729
SVM_freq_uni	0.8811	0.8897	0.8811	0.8825	0.8726
SVM_freq_uni_chi2	0.8811	0.8896	0.8811	0.8825	0.8726
SVM_freq_bi_chi2	0.8811	0.8896	0.8811	0.8825	0.8726
RF_tfidf_uni	0.8806	0.8913	0.8806	0.8823	0.8723
RF_freq_uni	0.8804	0.8909	0.8804	0.8820	0.8721
RF_freq_bi	0.8801	0.8911	0.8801	0.8818	0.8717
RF_tfidf_bi_chi2	0.8803	0.8880	0.8803	0.8812	0.8717
RF_freq_bi_chi2	0.8801	0.8898	0.8801	0.8817	0.8717
RF_tfidf_uni_chi2	0.8801	0.8888	0.8801	0.8814	0.8716
SVM_count_bi	0.8799	0.8893	0.8799	0.8814	0.8714
RF_freq_uni_chi2	0.8796	0.8903	0.8796	0.8813	0.8712
SVM_binary_bi	0.8798	0.8890	0.8798	0.8813	0.8712
RF_tfidf_bi	0.8797	0.8887	0.8797	0.8810	0.8712

20 rows × 5 columns

Model Tuning

In [71]:

```
1 from sklearn.model_selection import GridSearchCV
```

In [72]:

```
1 vectorizer= binary_vectorizer2
2 docs_vect = vectorizer.fit_transform(docs_train)
```

In [73]:

```

1 grid = {'solver': ['liblinear', 'lbfgs', 'sag', 'saga'],
2         'C': [10., 1.]}
3 grid_search = GridSearchCV(LogisticRegression(),
4                             grid,
5                             n_jobs=-1)
6 grid_search.fit(docs_vect, ytrain)

```

/home/00084776714/miniconda3/envs/tf/lib/python3.7/site-packages/sklearn/linear_model/_sag.py:330: ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
 "the coef_ did not converge", ConvergenceWarning)

Out[73]:

```

GridSearchCV(estimator=LogisticRegression(), n_jobs=-1,
              param_grid={'C': [10.0, 1.0],
                           'solver': ['liblinear', 'lbfgs', 'sag', 'saga']})

```

In [74]:

```
1 grid_search.best_params_
```

Out[74]:

```
{'C': 10.0, 'solver': 'sag'}
```

In [75]:

```

1 model= LogisticRegression(**grid_search.best_params_)
2 pipe = Pipeline(steps=[('vectorizer', vectorizer),
3                          ('model', model)])
4 pipe.fit(docs_train, ytrain)

```

/home/00084776714/miniconda3/envs/tf/lib/python3.7/site-packages/sklearn/linear_model/_sag.py:330: ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
 "the coef_ did not converge", ConvergenceWarning)

Out[75]:

```

Pipeline(steps=[('vectorizer',
                  CountVectorizer(binary=True, ngram_range=(1, 2))),
                ('model', LogisticRegression(C=10.0, solver='sag'))])

```

In [76]:

```

1 yv_pred = pipe.predict(docs_val)
2 matthews_corrcoef(yval, yv_pred)

```

Out[76]:

```
0.8784409142350387
```

Salvando o melhor modelo de ML

In [77]:

```
1 # save ML best model
2 with open(r"../best_models/ML_model.pickle", 'wb') as handle:
3     pickle.dump(pipe, handle)
```

Carregando o melhor modelo de ML

In [78]:

```
1 def load_doc(filename):
2     ''' # load doc into memory '''
3     # open the file as read only
4     file = open(filename, 'r')
5     # read all text
6     text = file.read()
7     # close the file
8     file.close()
9     return text
```

In [79]:

```
1 with open(r"../best_models/ML_model.pickle", 'rb') as handle:
2     loaded_pipe = pickle.load(handle)
3 # vocabulario
4 loaded_vocab = load_doc(r'../best_models/vocab.txt')
```

In [80]:

```
1 yv_pred = loaded_pipe.predict(docs_val)
2 matthews_corrcoef(yval, yv_pred)
```

Out[80]:

0.8784409142350387

In [81]:

```
1 print(classification_report(yval, yv_pred, digits=4))
```

	precision	recall	f1-score	support
25	0.8300	0.9030	0.8649	1000
30	0.9850	0.9880	0.9865	1000
40	0.9834	0.9490	0.9659	1000
60	0.9692	0.9750	0.9721	1000
70	0.9739	0.8970	0.9339	1000
71	0.9836	0.9600	0.9717	1000
72	0.9385	0.9920	0.9645	1000
80	0.9930	0.9870	0.9900	1000
90	0.7807	0.7440	0.7619	1000
120	0.8961	0.8190	0.8558	1000
130	0.8312	0.7240	0.7739	1000
140	0.7120	0.6800	0.6957	1000
200	0.9690	0.9060	0.9364	1000
999	0.6535	0.8900	0.7536	1000
accuracy			0.8867	14000
macro avg	0.8928	0.8867	0.8876	14000
weighted avg	0.8928	0.8867	0.8876	14000

In [82]:

```
1 # matriz de confuso: columns: True, index: Predict
2 cm = confusion_matrix(yval, yv_pred)
3 pd.DataFrame(data=cm, columns=np.sort(val.y.unique()), index=np.sort(val.y.unique()))
```

Out[82]:

	25	30	40	60	70	71	72	80	90	120	130	140	200	999
25	903	6	2	6	0	0	0	0	17	2	0	5	6	53
30	0	988	1	0	0	0	0	0	2	0	0	0	1	8
40	6	0	949	0	0	1	0	0	5	1	0	4	1	33
60	0	0	1	975	0	0	1	0	3	0	0	17	1	2
70	0	0	0	3	897	1	2	0	5	7	21	33	1	30
71	1	0	0	0	0	960	0	0	26	3	0	0	1	9
72	1	0	0	0	0	0	992	0	0	1	4	0	0	2
80	0	0	0	0	2	0	1	987	0	1	0	4	0	5
90	10	0	2	0	1	0	0	6	744	62	0	3	1	171
120	4	1	1	0	0	3	0	0	99	819	1	5	12	55
130	0	0	2	0	7	4	54	0	0	0	724	183	0	26
140	125	5	2	19	7	1	5	0	8	1	119	680	1	27
200	24	0	0	0	1	0	1	0	4	2	0	11	906	51
999	14	3	5	3	6	6	1	1	40	15	2	10	4	890

14 rows × 14 columns

```
1 90 = construção civil CNAE grupos 412, 432, 433, 439
```



```

2 120 = construção de obras de infraestrutura CNAE grupos 421, 422, 429, 431
3 isso explica a dificuldade de classificacao de atividades com servicos comuns
4 130 = operações de carga, descarga e armazenagem de conteineres em portos organizados
5 140 = transporte rodoviario de cargas

```

```

1 Fiz a exclusao no treino dos registros com classe diferente de 999 (outros serviços
  no sujeitos aliquota), sem informacao do codigo de serviço, e na lista: servico,
  servicos, servicos prestados, prestacao servico
2 para forcar a classificacao em 999, cujo resultado perceptvel na matriz de confusao

```

Fazendo predições ML:

In [83]:

```

1 def predict_ml_classe(lc, text, l_vocab, pipe):
2     l_text_clean = clean_doc(text, False)
3     l_text_clean = [w for w in l_text_clean if w in l_vocab]
4     l_lc = list()
5     lc_clean = re.sub(r'^\w\s', '', str(lc))
6     if lc_clean in l_vocab:
7         l_lc.append(lc_clean)
8         l_lc.extend(l_text_clean)
9         l_lc = [w for w in l_lc if w != '']
10    return pipe.predict(l_lc)[0]
11    return pipe.predict(l_text_clean)[0]

```

In [84]:

```

1 # Descrição de serviço de TI - classe 25
2 cd = ''
3 texto = 'desenvolvimento de sistemas'
4 classe = predict_ml_classe(cd, texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')

```

desenvolvimento de sistemas: classe 25

In [86]:

```

1 # Descrição de serviço de hotelaria - classe 40
2 cd = ''
3 texto = 'hospedagem'
4 classe = predict_ml_classe(cd, texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')

```

hospedagem: classe 40

In [87]:

```

1 # Descrição de serviço de obra de construo civil grupo CNAE 412, 432, 433,439 - classe
2 cd = ''
3 texto = 'obra de construção civil'
4 classe = predict_ml_classe(cd, texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')

```

obra de construção civil: classe 90

In [88]:

```
1 # Descrição de serviço de carga, descarga e armazenagem - classe 130
2 cd = ''
3 texto = 'armazenagem'
4 classe = predict_ml_classe(cd,texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')
```

armazenagem: classe 130

In [89]:

```
1 # Descrição de serviço prestado por empresas jornalísticas e de radiodifusão sonora - classe 200
2 cd = ''
3 texto = 'propaganda em jornal'
4 classe = predict_ml_classe(cd,texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')
```

propaganda em jornal: classe 200

In [90]:

```
1 import sklearn
2 print(sklearn.__version__)
```

0.23.2

In []:

```
1
```

In [1]:

```
1 import pickle
2 from joblib import dump, load
3 import re
4 import string
5 from nltk.corpus import stopwords
6 # sklearn.__version__ == 0.23.2
```

```
1 #conda install -c anaconda nltk
2 nltk.download('stopwords')
```

In [2]:

```
1 def clean_doc(doc=None, string_=True):
2     ''' # turn a doc into clean tokens '''
3     # split into tokens by white space
4     tokens = doc.split()
5     # prepare regex for char filtering
6     re_punc = re.compile('[%s]' % re.escape(string.punctuation))
7     # remove punctuation from each word
8     tokens = [re_punc.sub('', w) for w in tokens]
9     # remove remaining tokens that are not alphabetic
10    tokens = [word.lower() for word in tokens if not word.isdigit()]
11    # filter out stop words
12    stop_words = set(stopwords.words('portuguese'))
13    stop_words.update(['<br />'])
14    tokens = [w for w in tokens if not w in stop_words]
15    # filter out short tokens
16    tokens = [word for word in tokens if len(word) > 1]
17    if string_: return " ".join(tokens)
18    else: return tokens
19
20 def load_doc(filename):
21     ''' # load doc into memory '''
22     # open the file as read only
23     file = open(filename, 'r')
24     # read all text
25     text = file.read()
26     # close the file
27     file.close()
28     return text
29
30 def predict_ml_classe(lc, text, l_vocab, pipe):
31     l_text_clean = clean_doc(text, False)
32     l_text_clean = [w for w in l_text_clean if w in l_vocab]
33     l_lc = list()
34     lc_clean = re.sub(r'^\w\s', '', str(lc))
35     if lc_clean in l_vocab:
36         l_lc.append(lc_clean)
37         l_lc.extend(l_text_clean)
38         l_lc = [w for w in l_lc if w != '']
39         return pipe.predict(l_lc)[0]
40     return pipe.predict(l_text_clean)[0]
```

In [3]:

```
1 with open(r"../best_models/ML_model.pickle", 'rb') as handle:
2     loaded_pipe = pickle.load(handle)
3 # vocabulario
4 loaded_vocab = load_doc(r"../best_models/vocab.txt")
```

In [4]:

```
1 # Descrição de serviço de TI - classe 25
2 cd = ''
3 texto = 'desenvolvimento de sistemas'
4 classe = predict_ml_classe(cd,texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')
```

desenvolvimento de sistemas: classe 25

In [5]:

```
1 # Descrição de serviço de hotelaria - classe 40
2 cd = ''
3 texto = 'hospedagem'
4 classe = predict_ml_classe(cd,texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')
```

hospedagem: classe 40

In [6]:

```
1 # Descrição de serviço de obra de construo civil grupo CNAE 412, 432, 433,439 - classe
2 cd = ''
3 texto = 'obra de construção civil'
4 classe = predict_ml_classe(cd,texto, loaded_vocab, loaded_pipe)
5 print(f'{texto}: classe {classe}')
```

obra de construção civil: classe 90

In []:

```
1
```