

1 - Carrega informações do dataframe

2 - Criação do encoder e da matriz binária para cada coluna que iremos usar para classificação

3 - gera o corpus (bag of words) e vetoriza

4 - Faz a transformação TFIDF

5 - Cria bases de treinamento e validação

6 - Cria modelo

7 - Analisa a acurácia e a perda de cada modelo

8 - Cria predições para cada modelo treinado

9 - Aplica o modelo na base completa e analisa resultados

In [1]:

```
import pandas as pd
import numpy as np
import random
import os
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.preprocessing import LabelEncoder, LabelBinarizer, OneHotEncoder
from sklearn.naive_bayes import GaussianNB, MultinomialNB
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfTransformer
from collections import namedtuple
from typing import Dict
```

1 - Carrega informações do dataframe

1.1 - dados de descrição limpa (com stopwords e com afixos)

In [2]:

```
df_itens = pd.read_parquet('itens_desc_limpa.parquet')
df_tec = pd.read_parquet('2_tec_desc_limpa.parquet')
```

In [3]:

```
len(df_itens), len(df_tec)
```

Out[3]:

(26117, 10147)

In [4]:

```
df_itens[df_itens['capitulo'] == '00']
```

Out[4]:

	descricao_limpa	capitulo	posicao	subposicao	item	subitem
--	-----------------	----------	---------	------------	------	---------

In [5]:

```
# Duplicar linhas com somente 1 ou 2 exemplos
df_itens = df_itens.append(df_itens[df_itens['capitulo'].map(df_itens['capitulo'].value_cou
```

In [6]:

```
df_itens.head()
```

Out[6]:

	descricao_limpa	capitulo	posicao	subposicao	item	subitem
0	mascara facial hidratante em embalagem de 25ml...	33	04	99	1	0
1	diocetil ftalato em flexi bag d c 20 toneladas ...	29	17	32	0	0
2	sola de calcado de borracha n 39 a ref xl belf...	64	06	20	0	0
3	sola de calcado de borracha n 40 a ref xl lyon mx	64	06	20	0	0
4	sola de calcado de borracha n 39 a ref xl lyon mx	64	06	20	0	0

In [7]:

```
df_itens[df_itens['capitulo'] == '99']
```

Out[7]:

	descricao_limpa	capitulo	posicao	subposicao	item	subitem
4441889	x	99	99	99	9	9
4337496	seamer drone importado	99	99	99	9	9
4441889	x	99	99	99	9	9
4337496	seamer drone importado	99	99	99	9	9

In [8]:

```
len(df_itens)
```

Out[8]:

26117

In [9]:

```
# apaga linhas vazias
df_itens = df_itens.drop(df_itens[df_itens['descricao_limpa'] == ''].index)
```

In [10]:

```
len(df_itens)
```

Out[10]:

26117

In [11]:

```
df_tec.head()
```

Out[11]:

	descricao	ncm	ncm_str	capitulo	posicao	subposicao	item	subitem	descricao_
0	Reprodutores de raca pura Cavalos Cavalos as...	1012100.0	01012100	01	01	21	0	0	reprodu to raca pura c cavalos
1	Outros Cavalos Cavalos asininos e muares vi...	1012900.0	01012900	01	01	29	0	0	outros c cavalos as e muare:
2	Asininos Cavalos asininos e muares vivos	1013000.0	01013000	01	01	30	0	0	asininos c asir muare:
3	Outros Cavalos asininos e muares vivos	1019000.0	01019000	01	01	90	0	0	outros c asir muare:
4	Prenhes ou com cria ao pe Reprodutores de raca...	1022110.0	01022110	01	02	21	1	0	prenhes c cria reprodu to

In [12]:

```
df_tec = df_tec[['capitulo', 'posicao', 'subposicao', 'item', 'subitem']]
```

In [13]:

```
df_tec.head()
```

Out[13]:

	capitulo	posicao	subposicao	item	subitem
0	01	01	21	0	0
1	01	01	29	0	0
2	01	01	30	0	0
3	01	01	90	0	0
4	01	02	21	1	0

2 - Criação do encoder e da matriz binária para cada coluna que iremos usar para classificação.

In [14]:

```
Encoders = namedtuple('Encoders', 'encoder binarizer')

def encode_fields(df, fields: list) -> Dict[str, Encoders]:
    result = {}
    for i, field in enumerate(fields):
        lblencoder = LabelEncoder() # cria um número para cada categoria
        lblbinarizer = LabelBinarizer() # one hot encoder (ex: 99categorias cria matriz co
        encoded = lblencoder.fit_transform(df[field].values) # transforma os dados do arra
        # dessa forma, retorna um array com as categorias na forma numérica, começando em z
        print(f'field: {field} / encoded shape: {encoded.shape}')
        binarized = lblbinarizer.fit_transform(encoded) # transforma os dados do array "fi
        # dados binários (zeros e uns) para cada categoria, então retorna matriz m x n, ond
        # de linhas do array de entrada e n é a quantidade de colunas de categorias tranfor
        # forma binária
        print(f'field: {field} / binarized shape: {binarized.shape}')
        encoders = Encoders(lblencoder, lblbinarizer)
        result[field] = encoders
    return result
```

In [15]:

```
encoders = encode_fields(df_itens, ['capitulo', 'posicao', 'subposicao', 'item', 'subitem'])

field: capitulo / encoded shape: (26117,)
field: capitulo / binarized shape: (26117, 97)
field: posicao / encoded shape: (26117,)
field: posicao / binarized shape: (26117, 90)
field: subposicao / encoded shape: (26117,)
field: subposicao / binarized shape: (26117, 91)
field: item / encoded shape: (26117,)
field: item / binarized shape: (26117, 10)
field: subitem / encoded shape: (26117,)
field: subitem / binarized shape: (26117, 10)
```

2.1 - encode da coluna "capítulo"

In [17]:

```

y_encoded_cap = encoders['capitulo'].encoder.transform(df_itens.capitulo.values)
y_encoded_pos = encoders['posicao'].encoder.transform(df_itens.posicao.values)
y_encoded_subpos = encoders['subposicao'].encoder.transform(df_itens.subposicao.values)
y_encoded_item = encoders['item'].encoder.transform(df_itens.item.values)
y_encoded_subitem = encoders['subitem'].encoder.transform(df_itens.subitem.values)

```

In [18]:

```

print(f'formato do array "y_encoded_cap": {y_encoded_cap.shape} linhas, \nconteúdo: \n{y_en
print(f'formato do array "y_encoded_pos": {y_encoded_pos.shape} linhas, \nconteúdo: \n{y_en
print(f'formato do array "y_encoded_subpos": {y_encoded_subpos.shape} linhas, \nconteúdo: \
print(f'formato do array "y_encoded_item": {y_encoded_item.shape} linhas, \nconteúdo: \n{y_
print(f'formato do array "y_encoded_subitem": {y_encoded_subitem.shape} linhas, \nconteúdo:

```

```

formato do array "y_encoded_cap": (26117,) linhas,
conteúdo:
[32 28 63 ... 83 96 96]
formato do array "y_encoded_pos": (26117,) linhas,
conteúdo:
[ 3 16  5 ... 27 89 89]
formato do array "y_encoded_subpos": (26117,) linhas,
conteúdo:
[90 23 11 ... 62 90 90]
formato do array "y_encoded_item": (26117,) linhas,
conteúdo:
[1 0 0 ... 9 9 9]
formato do array "y_encoded_subitem": (26117,) linhas,
conteúdo:
[0 0 0 ... 0 9 9]

```

2.2 - encode binário da coluna "capítulo" - gera matriz

In [19]:

```

y_cap = encoders['capitulo'].binarizer.fit_transform(y_encoded_cap)
y_pos = encoders['posicao'].binarizer.fit_transform(y_encoded_pos)
y_subpos = encoders['subposicao'].binarizer.fit_transform(y_encoded_subpos)
y_item = encoders['item'].binarizer.fit_transform(y_encoded_item)
y_subitem = encoders['subitem'].binarizer.fit_transform(y_encoded_subitem)
y_todos = [y_cap, y_pos, y_subpos, y_item, y_subitem]

```

In [20]:

```

print(f'formato da matriz "y_cap": {y_cap.shape} (linhas, colunas),\nconteúdo: \n{y_cap}')
print(f'formato da matriz "y_pos": {y_pos.shape} (linhas, colunas),\nconteúdo: \n{y_pos}')
print(f'formato da matriz "y_subpos": {y_subpos.shape} (linhas, colunas),\nconteúdo: \n{y_s
print(f'formato da matriz "y_item": {y_item.shape} (linhas, colunas),\nconteúdo: \n{y_item}')
print(f'formato da matriz "y_subitem": {y_subitem.shape} (linhas, colunas),\nconteúdo: \n{y

```

formato da matriz "y_cap": (26117, 97) (linhas, colunas),
conteúdo:

```

[[0 0 0 ... 0 0 0]
 [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 1]]

```

formato da matriz "y_pos": (26117, 90) (linhas, colunas),
conteúdo:

```

[[0 0 0 ... 0 0 0]
 [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 1]]

```

formato da matriz "y_subpos": (26117, 91) (linhas, colunas),
conteúdo:

```

[[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 1]]

```

formato da matriz "y_item": (26117, 10) (linhas, colunas),
conteúdo:

```

[[0 1 0 ... 0 0 0]
 [1 0 0 ... 0 0 0]
 [1 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 1]
 [0 0 0 ... 0 0 1]
 [0 0 0 ... 0 0 1]]

```

formato da matriz "y_subitem": (26117, 10) (linhas, colunas),
conteúdo:

```

[[1 0 0 ... 0 0 0]
 [1 0 0 ... 0 0 0]
 [1 0 0 ... 0 0 0]
 ...
 [1 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 1]
 [0 0 0 ... 0 0 1]]

```

3 - gera o corpus (bag of words) e vetoriza

In [21]:

```
corpus = df_itens.descricao_limpa.values # transforma todo o texto em bag of words
vectorizer = CountVectorizer(max_df=0.1, min_df=0.00001) # elimina palavras mto ou pouco f
X_counts = vectorizer.fit_transform(corpus) # aprende o dicionário de vocabulários e gera
X_counts.shape
```

Out[21]:

(26117, 24548)

In [22]:

```
X_pos = y_pos
X_pos.shape
```

Out[22]:

(26117, 90)

In [23]:

```
X_subpos = y_subpos
X_subpos.shape
```

Out[23]:

(26117, 91)

In [24]:

```
X_item = y_item
X_item.shape
```

Out[24]:

(26117, 10)

In [25]:

```
X_subitem = y_subitem
X_subitem.shape
```

Out[25]:

(26117, 10)

limpa memória excluindo variável que não será mais utilizada

In [26]:

```
# del vectorizer
del corpus # apaga o corpus que gerou X_counts
```

4 - Faz a transformação TFIDF - realiza o cálculo da frequência relativa das palavras multiplicando por um peso, de forma a diminuir as palavras muito frequentes e as raras.

In [27]:

```
transformer = TfidfTransformer() # transforma em matrix TFIDF - faz freq relativa multipli  
# peso nas palavras freq ou raras - ou seja, deixa de ser zero e 1.  
X_tf = transformer.fit_transform(X_counts)
```

In [28]:

```
X_tf.shape
```

Out[28]:

```
(26117, 24548)
```

5 - Cria bases de treinamento e validação

5.1 - a base de teste representa 5% do dataset e está estratificada conforme os rótulos de "y" (matriz binária)

In [29]:

```
X_train_cap, X_val_cap, y_train_cap, y_val_cap = train_test_split(X_tf, y_cap, test_size=0.05)
```

In [30]:

```
X_train_pos, X_val_pos, y_train_pos, y_val_pos = train_test_split(X_tf, y_pos, test_size=0.05)
```

In [31]:

```
X_train_subpos, X_val_subpos, y_train_subpos, y_val_subpos = train_test_split(X_tf, y_subpos, test_size=0.05)
```

In [32]:

```
X_train_item, X_val_item, y_train_item, y_val_item = train_test_split(X_tf, y_item, test_size=0.05)
```

In [33]:

```
X_train_subitem, X_val_subitem, y_train_subitem, y_val_subitem = train_test_split(X_tf, y_subitem, test_size=0.05)
```

In [34]:

```
X_train_todos = [X_train_cap, X_train_pos, X_train_subpos, X_train_item, X_train_subitem]  
X_val_todos = [X_val_cap, X_val_pos, X_val_subpos, X_val_item, X_val_subitem]  
y_train_todos = [y_train_cap, y_train_pos, y_train_subpos, y_train_item, y_train_subitem]  
y_val_todos = [y_val_cap, y_val_pos, y_val_subpos, y_val_item, y_val_subitem]
```


In [35]:

```
for X_train in X_train_todos:
    print(f'Treinando com {X_train.shape[0]} exemplos da base e {X_train.shape} palavras di
```

Treinando com 24811 exemplos da base e (24811, 24548) palavras diferentes
 Treinando com 24811 exemplos da base e (24811, 24548) palavras diferentes
 Treinando com 24811 exemplos da base e (24811, 24548) palavras diferentes
 Treinando com 24811 exemplos da base e (24811, 24548) palavras diferentes
 Treinando com 24811 exemplos da base e (24811, 24548) palavras diferentes

In [36]:

```
for X_val in X_val_todos:
    print(f'Validando com {X_val.shape[0]} exemplos da base e {X_val.shape} palavras difere
```

Validando com 1306 exemplos da base e (1306, 24548) palavras diferentes
 Validando com 1306 exemplos da base e (1306, 24548) palavras diferentes
 Validando com 1306 exemplos da base e (1306, 24548) palavras diferentes
 Validando com 1306 exemplos da base e (1306, 24548) palavras diferentes
 Validando com 1306 exemplos da base e (1306, 24548) palavras diferentes

In [37]:

```
for y_train in y_train_todos:
    print(f'y_train com {y_train.shape[0]} exemplos da base e {y_train.shape} palavras dife
```

y_train com 24811 exemplos da base e (24811, 97) palavras diferentes
 y_train com 24811 exemplos da base e (24811, 90) palavras diferentes
 y_train com 24811 exemplos da base e (24811, 91) palavras diferentes
 y_train com 24811 exemplos da base e (24811, 10) palavras diferentes
 y_train com 24811 exemplos da base e (24811, 10) palavras diferentes

In [38]:

```
for y_val in y_val_todos:
    print(f'y_val com {y_val.shape[0]} exemplos da base e {y_val.shape} palavras diferentes
```

y_val com 1306 exemplos da base e (1306, 97) palavras diferentes
 y_val com 1306 exemplos da base e (1306, 90) palavras diferentes
 y_val com 1306 exemplos da base e (1306, 91) palavras diferentes
 y_val com 1306 exemplos da base e (1306, 10) palavras diferentes
 y_val com 1306 exemplos da base e (1306, 10) palavras diferentes

In [39]:

```
# cria dicionário com todos os modelos por coluna
setups = {}
colunas = ['cap', 'pos', 'subpos', 'item', 'subitem']
for i, coluna in enumerate(colunas):
    setups[coluna] = [X_train_todos[i], X_val_todos[i], y_train_todos[i], y_val_todos[i]]
```

In [40]:

```
#### Limpa memória excluindo variáveis que não serão mais utilizadas
```

In [41]:

```
del X_tf # apaga o resultado do TFIDF que foi utilizado para criar a base de teste e valid
```

In [42]:

```
del X_counts # apaga o X_count que originou o X_tf
```

6 - Cria modelo

6.1 - Cria modelo classificador usando 2 camadas full connected (densidade passada por parâmetro 256 ou 512 neurônios e ativador "relu") com dropout passado por parâmetro sendo de 20% ou de 40%, para reduzir overfitting

obs: foram utilizadas 2 camadas pois com duas camadas é suficiente para identificamos relações não lineares, mais de duas camadas teríamos que treinar o modelo muitas vezes o que tornaria mais complexo.

In [43]:

```
import tensorflow as tf
from tensorflow.keras import layers
from tensorflow.keras import backend as K

def model1(input_size, output_size, optimizer='adam', dropout=0.4, dense=128): # "adam" co
    # e tempo, é mais usado (SGD, SGD com momentum)
    model = tf.keras.Sequential()
    model.add(layers.Input(input_size))
    model.add(layers.Dense(dense, activation='relu')) # ativador do neurônio função relu
    model.add(layers.Dropout(dropout)) # a cada passada ignora 40% dos neurônios
    model.add(layers.Dense(dense, activation='relu'))
    model.add(layers.Dropout(dropout))
    model.add(layers.Dense(output_size, activation='softmax')) # coleção de 0, 1 do sigmoide
    model.compile(optimizer=optimizer,
                  loss='categorical_crossentropy', # pega o softmax, onde tá zero penaliza
                  metrics=['accuracy'])

    return model
```

6.2 - Cria modelos com dropouts e densidades diferentes

In [44]:

```
# melhor configuração 256 neuronios e dropout 0.2
models = {}
dense = 256
dropout = 0.2
i = 0
for k, setup in setups.items():
    print(f'model {k}: Xtrain: {setup[0].shape[1]} e y_train: {setup[2].shape[1]} ')
    models[k] = model1(
        setup[0].shape[1],
        setup[2].shape[1],
        optimizer=tf.keras.optimizers.Adam(lr=0.001),
        dropout=dropout, dense=dense)
```

```
model cap: Xtrain: 24548 e y_train: 97
model pos: Xtrain: 24548 e y_train: 90
model subpos: Xtrain: 24548 e y_train: 91
model item: Xtrain: 24548 e y_train: 10
model subitem: Xtrain: 24548 e y_train: 10
```

6.3 - treina esses modelos sendo para cada modelo roda 40 épocas, dividindo a entrada em chunks de tamanho 512, para cada época e com um learning rate decrescente

In [46]:

```

from collections import defaultdict
import math

epochs = 40
# batch_size = 256
batch_size = 512
for key, setup in setups.items():
    rounds = setup[0].shape[0] // batch_size + 1
    history = defaultdict(list)
    X_val_array = setup[1].toarray()

    print(f'\n\n modelo: {key} \n\n')
    for i in range(epochs):
        lr = 0.001 / (math.sqrt(i) + 1)
        print(f'Epoch {i} learning rate {lr}')
        K.set_value(models[key].optimizer.lr, lr)
        for batch_number in range(rounds):
            start = batch_number * batch_size
            X_chunk = setup[0][start: start + batch_size].toarray()
            y_chunk = setup[2][start: start + batch_size]
            models[key].train_on_batch(X_chunk, y_chunk) # treina o modelo efetivamente
            if batch_number % 100 == 0.:
                print(f'Batch n.: {batch_number} de {rounds}')
                loss_acc = models[key].evaluate(X_chunk, y_chunk)
                history['train_loss'].append(loss_acc[0])
                history['train_acc'].append(loss_acc[1])
                val_loss_acc = models[key].evaluate(X_val_array, setup[3])
                history['val_loss'].append(val_loss_acc[0])
                history['val_acc'].append(val_loss_acc[1])
                # print('loss: {:.2f} acc: {:.2f}'.format(val_monitor[0], val_monitor[1]))
        print('#####')
        print(f'Final da época {i}')
        models[key].evaluate(X_chunk, y_chunk)
        models[key].evaluate(setup[1].toarray(), setup[3])
        print('#####')
        del X_chunk
        del y_chunk

```

modelo: cap

Epoch 0 learning rate 0.001

Batch n.: 0 de 49

16/16 [=====] - 0s 6ms/step - loss: 0.0229 - accuracy: 0.9961

41/41 [=====] - 0s 7ms/step - loss: 0.1816 - accuracy: 0.9426

#####

Final da época 0

8/8 [=====] - 0s 9ms/step - loss: 0.0057 - accuracy: 1.0000

41/41 [=====] - 0s 8ms/step - loss: 0.1855 - accuracy: 0.9449

#####

Epoch 1 learning rate 0.0005

Batch n.: 0 de 49

7 - Analisa a acurária e a perda de cada modelo

In [47]:

```
for key, setup in setups.items():
    print(f'model: {key} - metrics: {models[key].metrics_names}')
    loss, acc = models[key].evaluate(setup[1].toarray(), setup[3])
    loss_teste, acc_teste = models[key].evaluate(setup[0][:1000].toarray(), setup[2][:1000])
```

```
model: cap - metrics: ['loss', 'accuracy']
41/41 [=====] - 0s 6ms/step - loss: 0.2031 - accuracy: 0.9479
32/32 [=====] - 0s 6ms/step - loss: 0.0057 - accuracy: 0.9980
(acc_teste - acc): 500.67
(loss_teste - loss): 19.73
```

```
model: pos - metrics: ['loss', 'accuracy']
41/41 [=====] - 0s 6ms/step - loss: 0.2964 - accuracy: 0.9181
32/32 [=====] - 0s 6ms/step - loss: 0.0295 - accuracy: 0.9920
(acc_teste - acc): 739.30
(loss_teste - loss): 26.68
```

```
model: subpos - metrics: ['loss', 'accuracy']
41/41 [=====] - 0s 8ms/step - loss: 1.1148 - accuracy: 0.7282
32/32 [=====] - 0s 7ms/step - loss: 0.3061 - accuracy: 0.9260
(acc_teste - acc): 1978.22
(loss_teste - loss): 80.87
```

```
model: item - metrics: ['loss', 'accuracy']
41/41 [=====] - 0s 6ms/step - loss: 0.9768 - accuracy: 0.7879
32/32 [=====] - 0s 6ms/step - loss: 0.1224 - accuracy: 0.9630
(acc_teste - acc): 1750.98
(loss_teste - loss): 85.44
```

```
model: subitem - metrics: ['loss', 'accuracy']
41/41 [=====] - 0s 6ms/step - loss: 0.6016 - accuracy: 0.8591
32/32 [=====] - 0s 6ms/step - loss: 0.0706 - accuracy: 0.9790
(acc_teste - acc): 1198.88
(loss_teste - loss): 53.11
```

8 - Cria predições para cada modelo treinado anteriormente e salva numa lista

In [48]:

```
preds_list = []
for key, setup in setups.items():
    teste = df_itens.descricao_limpa.values
    preds_list.append(models[key].predict(vectorizer.transform(teste)))
```

9 - Aplica o modelo na base completa e analisa resultados

9.1 - Com resultado da aplicação do modelo na base de dados completa, cria novas colunas para cada parte da NCM

In [49]:

```
i = 0
for key, model in models.items():
    name = key + '_resul'
    print(name)
    if key == 'cap':
        df_itens[name] = encoders['capitulo'].encoder.inverse_transform(encoders['capitulo']
    elif key == 'pos':
        df_itens[name] = encoders['posicao'].encoder.inverse_transform(encoders['posicao'].
    elif key == 'subpos':
        df_itens[name] = encoders['subposicao'].encoder.inverse_transform(encoders['subposi
    elif key == 'item':
        df_itens[name] = encoders['item'].encoder.inverse_transform(encoders['item'].binari
    elif key == 'subitem':
        df_itens[name] = encoders['subitem'].encoder.inverse_transform(encoders['subitem'].
    i += 1
```

```
cap_resul
pos_resul
subpos_resul
item_resul
subitem_resul
```

In [50]:

```
df_itens.head()
```

Out[50]:

	descricao_limpa	capitulo	posicao	subposicao	item	subitem	cap_resul	pos_resul	subpo
0	mascara facial hidratante em embalagem de 25ml...	33	04	99	1	0	33	04	
1	diocetil ftalato em flexi bag d c 20 toneladas ...	29	17	32	0	0	29	17	
2	sola de calcado de borracha n 39 a ref xl belf...	64	06	20	0	0	64	06	
3	sola de calcado de borracha n 40 a ref xl lyon mx	64	06	20	0	0	64	06	
4	sola de calcado de borracha n 39 a ref xl lyon mx	64	06	20	0	0	64	06	

9.2 - Recria os campos de NCM e cria uma NCM_result com as colunas resultado da aplicação do modelo

In [51]:

```
df_itens['ncm'] = df_itens['capitulo'] + df_itens['posicao'] + df_itens['subposicao'] + df_
```

In [52]:

```
df_itens['ncm_resul'] = df_itens['cap_resul'] + df_itens['pos_resul'] + df_itens['subpos_re
```

In [105]:

```
df_itens.head()
```

Out[105]:

	descricao_limpa	capitulo	posicao	subposicao	item	subitem	cap_resul	pos_resul	subpo
0	mascara facial hidratante em embalagem de 25ml...	33	04	99	1	0	33	04	
1	diocetil ftalato em flexi bag d c 20 toneladas ...	29	17	32	0	0	29	17	
2	sola de calcado de borracha n 39 a ref xl belf...	64	06	20	0	0	64	06	
3	sola de calcado de borracha n 40 a ref xl lyon mx	64	06	20	0	0	64	06	
4	sola de calcado de borracha n 39 a ref xl lyon mx	64	06	20	0	0	64	06	

9.3 - Cria dataframe erro com os valores de NCM_resultado errados

In [106]:

```
df_erros = df_itens[df_itens['ncm'] != df_itens.ncm_resul]
```

In [107]:

```
print(f'Tamanho do dataset: {len(df_itens)} registros')
print(f'Quantidade de erros: {len(df_erros)}, o que representa {(len(df_erros)/len(df_itens))}
```

Tamanho do dataset: 26117 registros

Quantidade de erros: 6041, o que representa 23.13%

In [108]:

`df_erros.head()`

Out[108]:

	descricao_limpa	capitulo	posicao	subposicao	item	subitem	cap_resul	pos_resul	sub
9	tambor de metal d c 25kg de pasta de pigmento ...	32	19	90	3	0	32	19	
105	chaveiro imitacao da marca kipling	95	03	00	3	1	95	03	
124	caixa para relógio imitacao da marca diesel	42	02	99	0	0	42	02	
126	caixa para relógio imitacao da marca emporio a...	42	02	99	0	0	42	02	
127	certificado de garantia imitacao da marca rolex	49	11	10	1	0	49	11	

9.4 - Analisa os erros em cada parte da NCM

In [109]:

```

capitulos_err = list(df_erros[df_erros['capitulo'] != df_erros.cap_resul]['cap_resul'])

posicoes_err = list(df_erros[(df_erros['capitulo'] == df_erros.cap_resul) &
                             (df_erros['posicao'] != df_erros.pos_resul)]['pos_resul'])

subposicoes_err = list(df_erros[(df_erros['capitulo'] == df_erros.cap_resul) &
                                 (df_erros['posicao'] == df_erros.pos_resul) &
                                 (df_erros['subposicao'] != df_erros.subpos_resul)]['subpos_

itens_err = list(df_erros[(df_erros['capitulo'] == df_erros.cap_resul) &
                           (df_erros['posicao'] == df_erros.pos_resul) &
                           (df_erros['subposicao'] == df_erros.subpos_resul) &
                           (df_erros['item'] != df_erros.item_resul)]['item_resul'])

subitens_err = list(df_erros[(df_erros['capitulo'] == df_erros.cap_resul) &
                              (df_erros['posicao'] == df_erros.pos_resul) &
                              (df_erros['subposicao'] == df_erros.subpos_resul) &
                              (df_erros['item'] == df_erros.item_resul) &
                              (df_erros['subitem'] != df_erros.subitem_resul)]['subitem_resu

```

In [110]:

```

print(len(capitulos_err), len(posicoes_err), len(subposicoes_err), len(itens_err), len(subi
145 524 2867 1553 952

```

9.5 - Erros em capítulo - detalha quantidade de capítulos errados e quantos erros por capítulo

In [111]:

```
total_err = {}
for capitulo in capitulos_err:
    if total_err.get(capitulo):
        total_err[capitulo] += 1
    else:
        total_err[capitulo] = 1

print(f'Total de capítulos errados: {len(total_err.keys())}')
print(f'Total de erros em capítulos: {len(capitulos_err)} erros\n')

for k, v in total_err.items():
    total_value = len(df_itens[df_itens['capitulo'] == str(k).zfill(2)])
    print(f'Capítulo com erro: {k} => {v} erros em {total_value} = {(v/total_value)*100}:%.2f
```

Total de capítulos errados: 20

Total de erros em capítulos: 145 erros

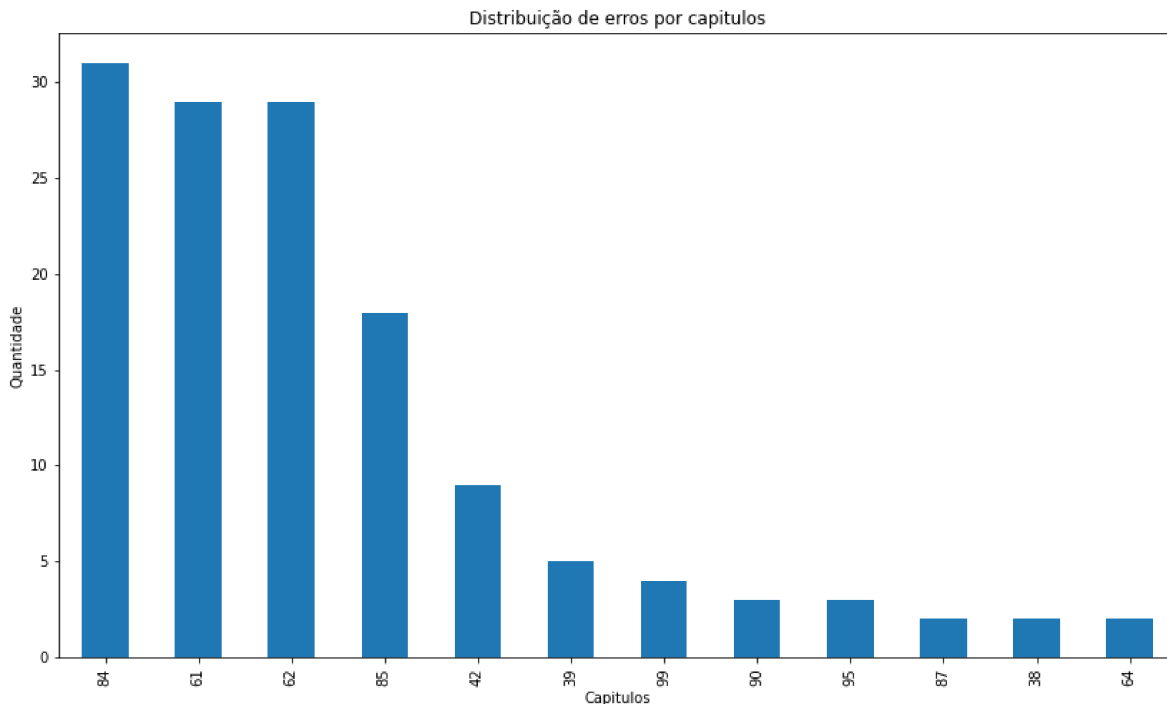
```
Capítulo com erro: 42 => 9 erros em 357 = 2.52%
Capítulo com erro: 62 => 29 erros em 1133 = 2.56%
Capítulo com erro: 39 => 5 erros em 596 = 0.84%
Capítulo com erro: 61 => 29 erros em 2176 = 1.33%
Capítulo com erro: 24 => 1 erros em 309 = 0.32%
Capítulo com erro: 94 => 1 erros em 100 = 1.00%
Capítulo com erro: 64 => 2 erros em 136 = 1.47%
Capítulo com erro: 49 => 1 erros em 40 = 2.50%
Capítulo com erro: 11 => 1 erros em 30 = 3.33%
Capítulo com erro: 68 => 1 erros em 77 = 1.30%
Capítulo com erro: 99 => 4 erros em 4 = 100.00%
Capítulo com erro: 87 => 2 erros em 371 = 0.54%
Capítulo com erro: 71 => 1 erros em 135 = 0.74%
Capítulo com erro: 84 => 31 erros em 3645 = 0.85%
Capítulo com erro: 85 => 18 erros em 5183 = 0.35%
Capítulo com erro: 95 => 3 erros em 1077 = 0.28%
Capítulo com erro: 90 => 3 erros em 737 = 0.41%
Capítulo com erro: 56 => 1 erros em 99 = 1.01%
Capítulo com erro: 38 => 2 erros em 310 = 0.65%
Capítulo com erro: 70 => 1 erros em 112 = 0.89%
```

In [128]:

```
# Cria gráfico de barras
df_temp = pd.DataFrame()
df_temp = df_erros[df_erros['capitulo'] != df_erros.cap_resul]
values = df_temp['cap_resul'].value_counts()
threshold = 1 # define limite inferior para exibição no gráfico (exibir 10 primeiros )
mask = values > threshold
values = values.loc[mask] # pega os valores que devem ser exibidos

# informações do gráfico
ax = values.plot.bar(figsize=(14,8), title="Distribuição de erros por capitulos")
ax.set_xlabel("Capitulos")
ax.set_ylabel("Quantidade")
print(f"Quantidade de capitulos errados: {len(df_temp['cap_resul'].value_counts())}")
```

Quantidade de capítulos errados: 20



9.6 - Erros em posição - detalha quantidade de posições erradass e quantos erros por posição

In [86]:

```

total_err = {}
for posicao in posicoes_err:
    if total_err.get(posicao):
        total_err[posicao] += 1
    else:
        total_err[posicao] = 1

print(f'Total de posições erradas: {len(total_err.keys())}')
print(f'Total de erros em posiçõss: {len(posicoes_err)} erros\n')

for k, v in total_err.items():
    total_value = len(df_itens[df_itens['posicao'] == str(k).zfill(2)])
    print(f'Posição com erro: {k} => {v} erro(s) em {total_value} = {(v/total_value)*100):

```

Total de posições erradas: 33

Total de erros em posiçõss: 524 erros

Posição com erro: 02 => 23 erro(s) em 1730 = 1.33%

Posição com erro: 01 => 5 erro(s) em 662 = 0.76%

Posição com erro: 17 => 20 erro(s) em 2294 = 0.87%

Posição com erro: 82 => 9 erro(s) em 1258 = 0.72%

Posição com erro: 07 => 11 erro(s) em 905 = 1.22%

Posição com erro: 09 => 22 erro(s) em 458 = 4.80%

Posição com erro: 06 => 22 erro(s) em 1262 = 1.74%

Posição com erro: 04 => 30 erro(s) em 2631 = 1.14%

Posição com erro: 05 => 9 erro(s) em 820 = 1.10%

Posição com erro: 14 => 58 erro(s) em 1518 = 3.82%

Posição com erro: 10 => 29 erro(s) em 584 = 4.97%

Posição com erro: 03 => 19 erro(s) em 2393 = 0.79%

Posição com erro: 18 => 3 erro(s) em 857 = 0.35%

Posição com erro: 26 => 5 erro(s) em 299 = 1.67%

Posição com erro: 08 => 7 erro(s) em 744 = 0.94%

Posição com erro: 15 => 9 erro(s) em 374 = 2.41%

Posição com erro: 11 => 5 erro(s) em 412 = 1.21%

Posição com erro: 13 => 4 erro(s) em 275 = 1.45%

Posição com erro: 12 => 4 erro(s) em 334 = 1.20%

Posição com erro: 19 => 6 erro(s) em 227 = 2.64%

Posição com erro: 16 => 2 erro(s) em 242 = 0.83%

Posição com erro: 27 => 1 erro(s) em 152 = 0.66%

Posição com erro: 25 => 9 erro(s) em 178 = 5.06%

Posição com erro: 62 => 1 erro(s) em 21 = 4.76%

Posição com erro: 51 => 5 erro(s) em 17 = 29.41%

Posição com erro: 77 => 3 erro(s) em 20 = 15.00%

Posição com erro: 46 => 2 erro(s) em 18 = 11.11%

Posição com erro: 71 => 155 erro(s) em 879 = 17.63%

Posição com erro: 73 => 4 erro(s) em 224 = 1.79%

Posição com erro: 43 => 3 erro(s) em 331 = 0.91%

Posição com erro: 28 => 37 erro(s) em 447 = 8.28%

Posição com erro: 44 => 1 erro(s) em 85 = 1.18%

Posição com erro: 24 => 1 erro(s) em 235 = 0.43%

In [133]:

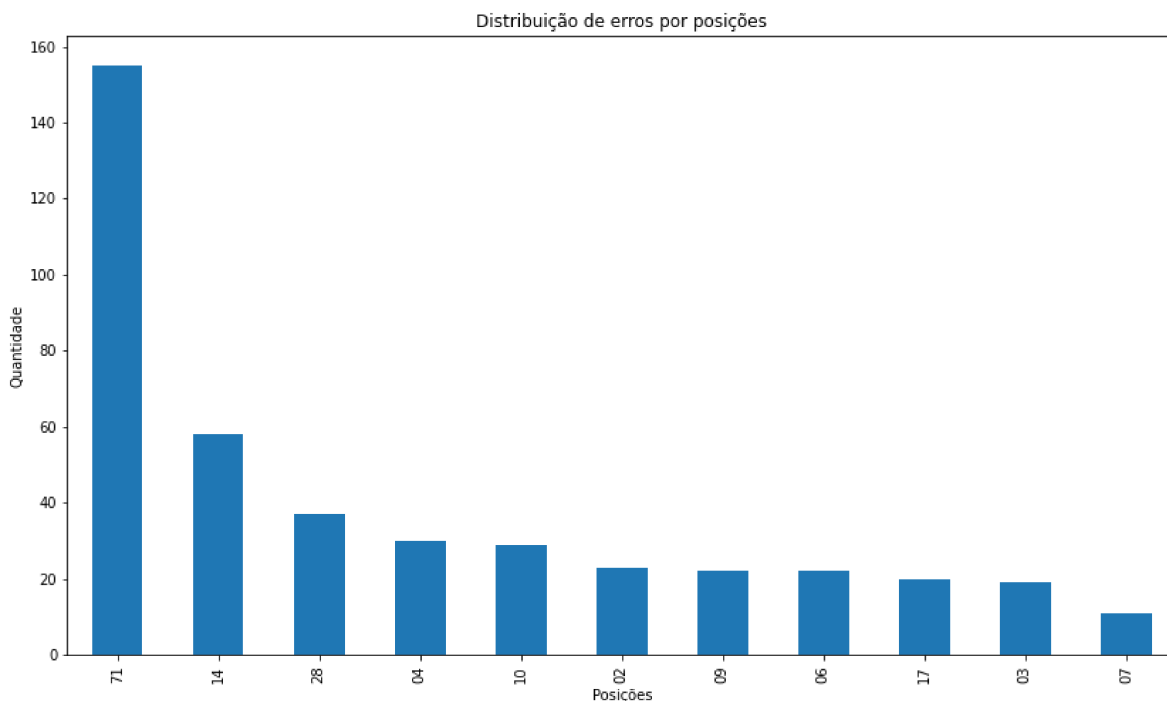
```

# Cria gráfico de barras
df_temp = pd.DataFrame()
df_temp = df_erros[(df_erros['capitulo'] == df_erros.cap_resul) &
                    (df_erros['posicao'] != df_erros.pos_resul)]
df_temp.head()
values = df_temp['pos_resul'].value_counts()
threshold = 10 # define limite inferior para exibição no gráfico (exibir 10 primeiros )
mask = values > threshold
values = values.loc[mask] # pega os valores que devem ser exibidos

# informações do gráfico
ax = values.plot.bar(figsize=(14,8), title="Distribuição de erros por posições")
ax.set_xlabel("Posições")
ax.set_ylabel("Quantidade")
print(f"Quantidade de posições erradas: {len(df_temp['cap_resul'].value_counts())}")

```

Quantidade de posições errados: 51



9.7 - Erros em subposição, item e subitem - calcula o erro por categoria e a quantidade de erros em cada categoria

In [87]:

```

total_err = {}
for subposicao in subposicoes_err:
    if total_err.get(subposicao):
        total_err[subposicao] += 1
    else:
        total_err[subposicao] = 1

print(f'Total de subposições erradas: {len(total_err.keys())}')
print(f'Total de erros em subposições: {len(subposicoes_err)} erros\n')

for k, v in total_err.items():
    total_value = len(df_itens[df_itens['subposicao'] == str(k).zfill(2)])
    print(f'Subposição com erro: {k} => {v} erro(s) em {total_value} = {(v/total_value)*100}')

```

Total de subposições erradas: 47

Total de erros em subposições: 2867 erros

```

Subposição com erro: 50 => 27 erro(s) em 1001 = 2.70%
Subposição com erro: 32 => 77 erro(s) em 257 = 29.96%
Subposição com erro: 11 => 83 erro(s) em 507 = 16.37%
Subposição com erro: 30 => 146 erro(s) em 2627 = 5.56%
Subposição com erro: 10 => 57 erro(s) em 2517 = 2.26%
Subposição com erro: 39 => 137 erro(s) em 393 = 34.86%
Subposição com erro: 21 => 81 erro(s) em 974 = 8.32%
Subposição com erro: 63 => 34 erro(s) em 209 = 16.27%
Subposição com erro: 40 => 90 erro(s) em 824 = 10.92%
Subposição com erro: 92 => 21 erro(s) em 276 = 7.61%
Subposição com erro: 43 => 82 erro(s) em 167 = 49.10%
Subposição com erro: 52 => 11 erro(s) em 50 = 22.00%
Subposição com erro: 90 => 343 erro(s) em 2599 = 13.20%
Subposição com erro: 99 => 238 erro(s) em 901 = 26.42%
Subposição com erro: 14 => 16 erro(s) em 64 = 25.00%
Subposição com erro: 42 => 61 erro(s) em 156 = 39.10%
Subposição com erro: 22 => 34 erro(s) em 191 = 17.80%
Subposição com erro: 12 => 62 erro(s) em 1280 = 4.84%
Subposição com erro: 53 => 21 erro(s) em 51 = 41.18%
Subposição com erro: 20 => 81 erro(s) em 1580 = 5.13%
Subposição com erro: 29 => 157 erro(s) em 628 = 25.00%
Subposição com erro: 19 => 220 erro(s) em 794 = 27.71%
Subposição com erro: 13 => 66 erro(s) em 97 = 68.04%
Subposição com erro: 49 => 125 erro(s) em 224 = 55.80%
Subposição com erro: 44 => 17 erro(s) em 63 = 26.98%
Subposição com erro: 89 => 72 erro(s) em 134 = 53.73%
Subposição com erro: 59 => 53 erro(s) em 139 = 38.13%
Subposição com erro: 23 => 23 erro(s) em 106 = 21.70%
Subposição com erro: 69 => 59 erro(s) em 113 = 52.21%
Subposição com erro: 62 => 38 erro(s) em 900 = 4.22%
Subposição com erro: 31 => 62 erro(s) em 244 = 25.41%
Subposição com erro: 51 => 40 erro(s) em 442 = 9.05%
Subposição com erro: 83 => 5 erro(s) em 13 = 38.46%
Subposição com erro: 94 => 7 erro(s) em 51 = 13.73%
Subposição com erro: 41 => 15 erro(s) em 737 = 2.04%
Subposição com erro: 81 => 8 erro(s) em 44 = 18.18%
Subposição com erro: 71 => 15 erro(s) em 391 = 3.84%
Subposição com erro: 91 => 36 erro(s) em 312 = 11.54%
Subposição com erro: 60 => 13 erro(s) em 246 = 5.28%
Subposição com erro: 00 => 13 erro(s) em 2007 = 0.65%
Subposição com erro: 70 => 15 erro(s) em 627 = 2.39%

```

Subposição com erro: 79 => 14 erro(s) em 50 = 28.00%
Subposição com erro: 77 => 1 erro(s) em 16 = 6.25%
Subposição com erro: 93 => 34 erro(s) em 86 = 39.53%
Subposição com erro: 80 => 34 erro(s) em 357 = 9.52%
Subposição com erro: 33 => 22 erro(s) em 141 = 15.60%
Subposição com erro: 61 => 1 erro(s) em 77 = 1.30%

In [88]:

```
total_err = {}
for item in itens_err:
    if total_err.get(item):
        total_err[item] += 1
    else:
        total_err[item] = 1

print(f'Total de itens erradas: {len(total_err.keys())}')
print(f'Total de erros em itens: {len(itens_err)} erros\n')

for k, v in total_err.items():
    total_value = len(df_itens[df_itens['item'] == str(k)])
    print(f'Itens com erro: {k} => {v} erro(s) em {total_value} = {(v/total_value)*100):.2%
```

Total de itens erradas: 10

Total de erros em itens: 1553 erros

Itens com erro: 1 => 267 erro(s) em 4628 = 5.77%
Itens com erro: 9 => 764 erro(s) em 5794 = 13.19%
Itens com erro: 2 => 254 erro(s) em 1646 = 15.43%
Itens com erro: 0 => 27 erro(s) em 10732 = 0.25%
Itens com erro: 3 => 137 erro(s) em 1603 = 8.55%
Itens com erro: 8 => 18 erro(s) em 135 = 13.33%
Itens com erro: 4 => 37 erro(s) em 842 = 4.39%
Itens com erro: 6 => 25 erro(s) em 110 = 22.73%
Itens com erro: 7 => 21 erro(s) em 291 = 7.22%
Itens com erro: 5 => 3 erro(s) em 336 = 0.89%

In [89]:

```
total_err = {}
for subitem in subitens_err:
    if total_err.get(subitem):
        total_err[subitem] += 1
    else:
        total_err[subitem] = 1

print(f'Total de subitens erradas: {len(total_err.keys())}')
print(f'Total de erros em subitens: {len(subitens_err)} erros\n')

for k, v in total_err.items():
    total_value = len(df_itens[df_itens['subitem'] == str(k)])
    print(f'Subitens com erro: {k} => {v} erro(s) em {total_value} = {(v/total_value)*100}
```

Total de subitens erradas: 10

Total de erros em subitens: 952 erros

Subitens com erro: 0 => 152 erro(s) em 19438 = 0.78%

Subitens com erro: 9 => 287 erro(s) em 2602 = 11.03%

Subitens com erro: 7 => 51 erro(s) em 258 = 19.77%

Subitens com erro: 2 => 86 erro(s) em 910 = 9.45%

Subitens com erro: 1 => 253 erro(s) em 2109 = 12.00%

Subitens com erro: 8 => 9 erro(s) em 50 = 18.00%

Subitens com erro: 4 => 33 erro(s) em 191 = 17.28%

Subitens com erro: 3 => 48 erro(s) em 343 = 13.99%

Subitens com erro: 5 => 32 erro(s) em 140 = 22.86%

Subitens com erro: 6 => 1 erro(s) em 76 = 1.32%

In [90]:

```

for i, row in enumerate(df_erros.iloc[:,0]):
    if df_erros.iloc[i,1] != df_erros.iloc[i,6]:
        print(f'errou capítulo {df_erros.iloc[i,1]} - ncm: {df_erros.iloc[i, 11]} - ncm_resu
    elif df_erros.iloc[i,2] != df_erros.iloc[i,7]:
        print(f'errou posicao {df_erros.iloc[i,2]} - ncm: {df_erros.iloc[i, 11]} - ncm_resu
    elif df_erros.iloc[i,3] != df_erros.iloc[i,8]:
        print(f'errou subposicao {df_erros.iloc[i,3]} - ncm: {df_erros.iloc[i, 11]} - ncm_r
    elif df_erros.iloc[i,4] != df_erros.iloc[i,9]:
        print(f'errou item {df_erros.iloc[i,4]} - ncm: {df_erros.iloc[i, 11]} - ncm_resul=
    elif df_erros.iloc[i,5] != df_erros.iloc[i,10]:
        print(f'errou subitem {df_erros.iloc[i,5]} - ncm: {df_erros.iloc[i, 11]} - ncm_resu

```

```

errou item 3 - ncm: 32199030 - ncm_resul= 32199010
errou subitem 1 - ncm: 95030031 - ncm_resul= 95030030
errou item 0 - ncm: 42029900 - ncm_resul= 42029990
errou item 0 - ncm: 42029900 - ncm_resul= 42029910
errou item 1 - ncm: 49111010 - ncm_resul= 49111090
errou item 9 - ncm: 39209990 - ncm_resul= 39209910
errou capítulo 91 - ncm: 91021110 - ncm_resul= 42021110
errou item 0 - ncm: 39172900 - ncm_resul= 39172920
errou subitem 6 - ncm: 22041096 - ncm_resul= 22041090
errou subposicao 00 - ncm: 95030021 - ncm_resul= 95035021
errou subitem 1 - ncm: 95030021 - ncm_resul= 95030020
errou item 9 - ncm: 84186999 - ncm_resul= 84186919
errou subitem 7 - ncm: 95030097 - ncm_resul= 95030099
errou posicao 26 - ncm: 39269090 - ncm_resul= 39023290
errou subposicao 90 - ncm: 39269090 - ncm_resul= 39263290
errou posicao 04 - ncm: 64041100 - ncm_resul= 64021900
errou subposicao 12 - ncm: 91021290 - ncm_resul= 91021190
errou item 9 - ncm: 91021190 - ncm_resul= 91021110
errou capítulo 61 - ncm: 61059000 - ncm_resul= 62053000

```

In []: