

How can I one-hot encode in Python?

Asked 7 years, 9 months ago Modified 1 year, 4 months ago Viewed 565k times

253

Highly active question

You have enough reputation to answer or unprotect this question.

I have a machine learning classification problem with 80% categorical variables. Must I use one-hot encoding if I want to use some classifier for the classification? Can I pass the data to a classifier without the encoding?

I am trying to do the following for feature selection:

- I read the train file:

```
num_rows_to_read = 10000
train_small = pd.read_csv("../dataset/train.csv", nrows=num_rows_to_read)
```

- I change the type of the categorical features to 'category':

```
non_categorical_features = ['orig_destination_distance',
                             'srch_haltz_cnt',
                             'srch_children_cnt',
                             'srch_rm_cnt',
                             'cnt']

for categorical_feature in list(train_small.columns):
    if categorical_feature not in non_categorical_features:
        train_small[categorical_feature] =
        train_small[categorical_feature].astype('category')
```

- I use one-hot encoding:

```
train_small_with_dummies = pd.get_dummies(train_small, sparse=True)
```

The problem is that the 3rd part often get stuck, although I am using a strong machine.

Thus, without the one-hot encoding I can't do any feature selection, for determining the importance of the features.

What do you recommend?

python pandas machine-learning one-hot-encoding Edit tags

Share Edit Follow Close Flag Unprotect

edited Aug 31, 2020 at 14:54

yafu
87.2k 12 90 143

asked May 18, 2016 at 7:26

avicohen
2,947 6 16 16

24

Answers

Sorted by
Reset to default

Date modified (newest first)

A simple example using vectorize in numpy and apply example in pandas:

```
import numpy as np
a = np.array(['male','female','female','male'])

def onehot_function(x):
    return 1.0 if (x=="male") else 0.0

onehot_a = np.vectorize(onehot_function)(a)

print(onehot_a)
# [1., 0., 0., 1.]

import pandas as pd
s = pd.Series(['male','female','female','male'])
onehot_s = s.apply(onehot_function)

print(onehot_s)
# 0    1.0
# 1    0.0
# 2    0.0
# 3    1.0
# dtype: float64
```

Share Edit Follow Flag

answered Sep 27, 2022 at 8:08

turbo
777 6 18

Lets assume out of 10 variables, you have 3 categorical variables in your data frame named as cname1, cname2 and cname3. Then following code will automatically create one-hot encoded variable in the new dataframe.

```
import category_encoders as ce
encoder_var=ce.OneHotEncoder(cols=
['cname1','cname2','cname3']).handle_unknown='return_nan',return_df=True,use_cat_names=True
new_df = encoder_var.fit_transform(old_df)
```

Share Edit Follow Flag

answered Oct 19, 2021 at 11:00

DSBLR
283 5 10

139

Much easier to use Pandas for basic one-hot encoding. If you're looking for more options you can use `scikit-learn`.

For basic one-hot encoding with Pandas you pass your data frame into the `get_dummies` function.

For example, if I have a dataframe called `imdb_movies`.

| | imdbRating | Awards | Language | Plot | Genre | Rated |
|---|------------|---|----------|---|----------------------|-------|
| 0 | 4.9 | NaN | English | A young girl and her coach overcome adversity ... | Drama, Family, Sport | PG |
| 1 | 8.2 | Won 4 Oscars. Another 33 wins & 67 nominations. | English | After John Nash, a brilliant but asocial mathe... | Biography, Drama | PG-13 |
| 2 | 5.7 | 1 win & 1 nomination. | English | A beautiful dancer balances on the razor's edg... | Comedy, Drama, Music | R |
| 3 | 7.8 | 1 nomination. | English | In a world rapidly being torn asunder by viole... | Drama | NaN |
| 4 | 7.6 | 11 wins & 2 nominations. | English | After his wife dies of cancer, an overworked e... | Drama, Fantasy | PG |

...and I want to one-hot encode the Rated column, I do this:

```
pd.get_dummies(imdb_movies, Rated)
```

| | APPROVED | G | NOT RATED | PASSED | PG | PG-13 | R | TV-14 | TV-G | TV-MA | TV-PG | TV-Y | TV-Y7 | TV-Y7-FV | UNRATED |
|---|----------|---|-----------|--------|----|-------|---|-------|------|-------|-------|------|-------|----------|---------|
| 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 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

This returns a new `dataframe` with a column for every "level" of rating that exists, along with either a 1 or 0 specifying the presence of that rating for a given observation.

Usually, we want this to be part of the original `dataframe`. In this case, we attach our new dummy coded frame onto the original frame using "column-binding".

We can column-bind by using Pandas `concat` function:

```
rated_dummies = pd.get_dummies(imdb_movies, Rated)
pd.concat([imdb_movies, rated_dummies], axis=1)
```

Usage:

```
encode_and_bind(imdb_movies, 'Rated')
```

Result:

| | APPROVED | G | NOT RATED | PASSED | PG | PG-13 | R | TV-14 | TV-G | TV-MA | TV-PG | TV-Y | TV-Y7 | TV-Y7-FV | UNRATED |
|-----|----------|---|-----------|--------|----|-------|---|-------|------|-------|-------|------|-------|----------|---------|
| 4.9 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8.2 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5.7 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Can we now run an analysis on our full `dataframe`.

SIMPLE UTILITY FUNCTION

I would recommend making yourself a utility function to do this quickly:

```
def encode_and_bind(original_dataframe, feature_to_encode):
    dummies = pd.get_dummies(original_dataframe[[feature_to_encode]])
    res = pd.concat([original_dataframe, dummies], axis=1)
    return(res)
```

Usage:

```
encode_and_bind(imdb_movies, 'Rated')
```

Result:

| | APPROVED | G | NOT RATED | PASSED | PG | PG-13 | R | TV-14 | TV-G | TV-MA | TV-PG | TV-Y | TV-Y7 | TV-Y7-FV | UNRATED |
|-----|----------|---|-----------|--------|----|-------|---|-------|------|-------|-------|------|-------|----------|---------|
| 4.9 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8.2 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5.7 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Also, as per @pmlbu comment, if you would like the function to remove the original `feature_to_encode` then use this version:

```
def encode_and_bind(original_dataframe, feature_to_encode):
    dummies = pd.get_dummies(original_dataframe[[feature_to_encode]])
    res = pd.concat([original_dataframe, dummies], axis=1)
    res = res.drop([feature_to_encode], axis=1)
    return(res)
```

You can encode multiple features at the same time as follows:

```
features_to_encode = ['feature_1', 'feature_2', 'feature_3',
                      'feature_4']
for feature in features_to_encode:
    res = encode_and_bind(train_set, feature)
```

Also, as per @pmlbu comment, if you would like the function to remove the original `feature_to_encode` then use this version:

```
def encode_and_bind(original_dataframe, feature_to_encode):
    dummies = pd.get_dummies(original_dataframe[[feature_to_encode]])
    res = pd.concat([original_dataframe, dummies], axis=1)
    res = res.drop([feature_to_encode], axis=1)
    return(res)
```

You can encode multiple features at the same time as follows:

```
features_to_encode = ['feature_1', 'feature_2', 'feature_3',
                      'feature_4']
for feature in features_to_encode:
    res = encode_and_bind(train_set, feature)
```

Share Edit Follow Flag

edited Oct 5, 2020 at 13:51

Harry B
2,874 1 24 44

answered Oct 22, 2018 at 18:07

Cybernetic
12,306 16 96 136

Approach 1: You can use pandas' `pd.get_dummies`.

Example 1:

```
import pandas as pd
s = pd.Series([list('abca')])
pd.get_dummies(s)
Out[1]:
   a  b  c
0  1.0  0.0  0.0
1  0.0  1.0  0.0
2  0.0  0.0  1.0
3  1.0  0.0  0.0
```

Example 2:

The following will transform a given column into one-hot. Use prefix to have multiple dummies.

```
import pandas as pd

df = pd.DataFrame({
    'A': ['a', 'b', 'a'],
    'B': ['b', 'a', 'c']
})

df
Out[1]:
   A  B
0  a  b
1  b  a
2  a  c

# Get one-hot encoding of columns B
one_hot = pd.get_dummies(df, B)
# Drop column B as it is one-encoded
df = df.drop('B', axis=1)
# Join the encoded df
df = df.join(one_hot)
df
Out[1]:
   A  a  b  c
0  a  0  1  0
1  b  1  0  0
2  a  0  0  1
```

Approach 2: Use Scikit-learn

Using a `OneHotEncoder` has the advantage of being able to `fit` on some training data and then `transform` on some other data using the same instance. We also have `handle_unknown` to further control what the encoder does with `unseen` data.

Given a dataset with three features and four samples, we let the encoder find the maximum value per feature and transform the data to a binary one-hot encoding.

```
>>> from sklearn.preprocessing import OneHotEncoder
>>> enc = OneHotEncoder()
>>> enc.fit([[0, 0, 3], [1, 1, 0], [0, 2, 1], [1, 0, 2]])
OneHotEncoder(categorical_features='all', dtype='float64',
             handle_unknown='error', n_values='auto', sparse=True)
>>> enc.n_values_
array([2, 3, 4])
>>> enc.feature_indices_
array([0, 2, 5, 9], dtype=int32)
>>> enc.transform([[0, 1, 1]], toarray())
array([[ 1.,  0.,  0.,  1.,  0.,  1.,  0.,  0.]])
```

Here is the link for this example: <http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>

Share Edit Follow Flag

edited Aug 31, 2020 at 15:00

answered Sep 2, 2016 at 7:55

yafu 87.2k 12 90 143

Sayali Sonawane 12.4k 5 46 47

30 setting `drop_first=True` with `get_dummies` removes the need to drop the original column separately - [OverflowingTheGlass Feb 28, 2018 at 15:14](#)

1 In example 2, is there a way to join the new columns to the dataframe without using join? I'm dealing with a really big dataset and get MemoryError when I try to do that - [JD May 31, 2018 at 12:27](#)

31 @OverflowingTheGlass- drop-first: True does not remove the original column. It drops the first level of the categorical feature so that you end up with k-1 columns instead of k columns, k being the cardinality of the categorical feature. - [Garima Jain Feb 28, 2019 at 13:50](#)

3 the df.join() does not work here, it creates more rows... do not know why though. - [Chenxi Zeng Oct 6, 2019 at 21:10](#)

2 df.join() creates more rows for me, so I used pd.concat([data, cat_encoded], axis=1) to join the encoded columns with the original dataset - [Ajay Bhary Dec 7, 2020 at 13:17](#)

Try this:

```
2  [!p] install category encoders
import category_encoders as ce

categorical_columns = [...] the list of names of the columns you want to one-hot-encode
...
encoder = ce.OneHotEncoder(cols=categorical_columns, use_cat_names=True)
df_train_encoded = encoder.fit_transform(df_train_small)

df_encoded.head()

The resulting dataframe df_train_encoded is the same as the original, but the categorical features are now replaced with their one-hot-encoded versions.

More information on category_encoders here.
```

Share Edit Follow Flag

answered Mar 5, 2020 at 10:03

Andres Araldo 1,382 14 20

Short Answer

1 Here is a function to do one-hot-encoding without using numpy, pandas, or other packages. It takes a list of integers, booleans, or strings (and perhaps other types too).

```
import typing

def one_hot_encode(items: list) -> typing.List[list]:
    results = []
    # find the unique items (we want to unique items b/c duplicate items will have the same encoding)
    unique_items = list(set(items))
    # sort the unique items
    sorted_items = sorted(unique_items)
    # find how long the list of each item should be
    max_index = len(unique_items)

    for item in items:
        # create a list of zeros the appropriate length
        one_hot_encoded_result = [0 for i in range(0, max_index)]
        # find the index of the item
        one_hot_index = sorted_items.index(item)
        # change the zero at the index from the previous line to a one
        one_hot_encoded_result[one_hot_index] = 1
        # add the result
        results.append(one_hot_encoded_result)

    return results
```

Example:

```
one_hot_encode([2, 1, 1, 2, 5, 3])

# [[0, 1, 0, 0],
#  [1, 0, 0, 0],
#  [1, 0, 0, 0],
#  [0, 1, 0, 0],
#  [0, 0, 1, 1],
#  [0, 0, 1, 0]]

one_hot_encode([True, False, True])

# [[0, 1], [1, 0], [0, 1]]

one_hot_encode(['a', 'b', 'c', 'a', 'a'])

# [[1, 0, 0, 0], [0, 1, 0, 0], [0, 0, 1, 0], [1, 0, 0, 0], [0, 0, 0, 1]]
```

Long(er) Answer

I know there are already a lot of answers to this question, but I noticed two things. First, most of the answers use packages like numpy and/or pandas. And this is a good thing. If you are writing production code, you should probably be using robust, fast algorithms like those provided in the numpy/pandas packages. But, for the sake of education, I think someone should provide an answer which has a transparent algorithm and not just an implementation of someone else's algorithm. Second, I noticed that many of the answers do not provide a robust implementation of one-hot encoding because they do not meet one of the requirements below. Below are some of the requirements (as I see them) for a useful, accurate, and robust one-hot encoding function:

A one-hot encoding function must:

- handle list of various types (e.g. integers, strings, floats, etc.) as input
- handle an input list with duplicates
- return a list of lists corresponding (in the same order as) to the inputs
- return a list of lists where each list is as short as possible

I tested many of the answers to this question and most of them fail on one of the requirements above.

Share Edit Follow Flag

answered Feb 5, 2020 at 2:04

Lloyd 1,382 20 25

Expanding @Martin Thoma's answer

```
def one_hot_encode(y):
    """Convert an iterable of indices to one-hot encoded labels."""
    y = y.flatten() # Sometimes not flattened vector is passed e.g. [(18,1) in these cases
    # the function ends up creating a tensor e.g. [(18, 2, 1)]. flatten removes this issue
    nb_classes = len(np.unique(y)) # get the number of unique classes
    standardised_labels = dict(zip(np.arange(nb_classes), np.arange(nb_classes))) # get the class labels as a dictionary
    # which then is standardised. E.g. imagine class labels are (4,7,9) if a vector of y containing 4,7 and 9 is
    # directly passed then np.eye(nb_classes)[4] or 7,9 throws an out of index error.
    # standardised_labels fixes this issue by returning a dictionary.
    # standardised_labels = {4:0, 7:1, 9:2}. The values of the dictionary are mapped to keys in y array.
    # standardised_labels also removes the error that is raised if the labels are floats. E.g. 1.0: element
    # cannot be called by integer index e.g. y[1.0] - throws an index error.
    targets = np.vectorize(standardised_labels.get)(y) # map the dictionary values to array.
    return np.eye(nb_classes)[targets]
```

Share Edit Follow Flag

answered Dec 29, 2019 at 12:36

mcsgiaradic 71 8

You can do the following as well. Note for the below you don't have to use `pd.concat`.

```
import pandas as pd
# Initialise data of lists.
data = {'Color': ['Red', 'Yellow', 'Red', 'Yellow'], 'length': [20.1, 21.1, 19.1, 18.1],
        'Group': [1, 2, 1, 2]}

# Create DataFrame
df = pd.DataFrame(data)

for c in df.select_dtypes(include='object').columns:
    print(c)
    df[c] = pd.Categorical(df[c])
df_transformed = pd.get_dummies(df)
df_transformed
```

You can also change explicit columns to categorical. For example, here I am changing the `Color` and `Group`

```
import pandas as pd
# Initialise data of lists.
data = {'Color': ['Red', 'Yellow', 'Red', 'Yellow'], 'length': [20.1, 21.1, 19.1, 18.1],
        'Group': [1, 2, 1, 2]}

# Create DataFrame
df = pd.DataFrame(data)
columns_to_change = list(df.select_dtypes(include='object').columns)
for c in columns_to_change:
    print(c)
    df[c] = pd.Categorical(df[c])
df_transformed = pd.get_dummies(df)
df_transformed
```

Share Edit Follow Flag

edited Aug 30, 2019 at 4:23

answered Aug 30, 2019 at 4:18

Aashmit 4,487 2 37 39

It can and it should be easy as :

```
class OneHotEncoder:
    def __init__(self, optionKeys):
        length=len(optionKeys)
        self.__dict__={optionKeys[i]:0 if i!=j else 1 for i in range(length)} for j in range(length)

Usage:

ohe=OneHotEncoder(["A","B","C","D"])
print(ohe.A)
print(ohe.B)
```

Share Edit Follow Flag

answered May 24, 2019 at 5:07

