

Character Encodings

Avoid UnicodeDecodeErrors when loading CSV files.

Tutorial Data



Learn Tutorial

Data Cleaning

In this notebook, we're going to be working with different character encodings.

Let's get started!

Get our environment set up

The first thing we'll need to do is load in the libraries we'll be using. Not our dataset, though: we'll get to it later!

```
# modules we'll use  
import pandas as pd  
import numpy as np  
  
# helpful character encoding module  
import chardet  
  
# set seed for reproducibility  
np.random.seed(0)
```

What are encodings?

Character encodings are specific sets of rules for mapping from raw binary byte strings (that look like this: 0110100001101001) to characters that make up human-readable text (like "hi"). There are many different encodings, and if you tried to read in text with a different encoding than the one it was originally written in, you ended up with scrambled text called "mojibake" (said like mo-gee-bah-kay). Here's an example of mojibake:

æ-ſå—å€-ã??

You might also end up with a "unknown" characters. There are what gets printed when there's no mapping between a particular byte and a character in the encoding you're using to read your byte string in and they look like this:

????????

Character encoding mismatches are less common today than they used to be, but it's definitely still a problem. There are lots of different character encodings, but the main one you need to know is UTF-8.

UTF-8 is **the** standard text encoding. All Python code is in UTF-8 and, ideally, all your data should be as well. It's when things aren't in UTF-8 that you run into trouble.

It was pretty hard to deal with encodings in Python 2, but thankfully in Python 3 it's a lot simpler. (Kaggle Notebooks only use Python 3.) There are two main data types you'll encounter when working with text in Python 3. One is the string, which is what text is by default.

```
# start with a string  
before = "This is the euro symbol: €"  
  
# check to see what datatype it is  
type(before)
```

```
str
```

The other data is the `bytes` data type, which is a sequence of integers. You can convert a string into bytes by specifying which encoding it's in:

```
# encode it to a different encoding, replacing characters that raise errors  
after = before.encode("utf-8", errors="replace")  
  
# check the type  
type(after)
```

```
bytes
```

If you look at a bytes object, you'll see that it has a b in front of it, and then maybe some text after. That's because bytes are printed out as if they were characters encoded in ASCII. (ASCII is an older character encoding that doesn't really work for writing any language other than English.) Here you can see that our euro symbol has been replaced with some mojibake that looks like `"\xe2\x82\xac"` when it's printed as if it were an ASCII string.

```
# take a look at what the bytes look like
after
```

```
:
b'This is the euro symbol: \xe2\x82\xac'
```

When we convert our bytes back to a string with the correct encoding, we can see that our text is all there correctly, which is great! :)

```
# convert it back to utf-8
print(after.decode("utf-8"))
```

```
This is the euro symbol: €
```

However, when we try to use a different encoding to map our bytes into a string, we get an error. This is because the encoding we're trying to use doesn't know what to do with the bytes we're trying to pass it. You need to tell Python the encoding that the byte string is actually supposed to be in.

You can think of different encodings as different ways of recording music. You can record the same music on a CD, cassette tape or 8-track. While the music may sound more-or-less the same, you need to use the right equipment to play the music from each recording format. The correct decoder is like a cassette player or a CD player. If you try to play a cassette in a CD player, it just won't work.

:

```
# try to decode our bytes with the ascii encoding  
print(after.decode("ascii"))
```

We can also run into trouble if we try to use the wrong encoding to map from a string to bytes. Like I said earlier, strings are UTF-8 by default in Python 3, so if we try to treat them like they were in another encoding we'll create problems.

For example, if we try to convert a string to bytes for ASCII using `encode()`, we can ask for the bytes to be what they would be if the text was in ASCII. Since our text isn't in ASCII, though, there will be some characters it can't handle. We can automatically replace the characters that ASCII can't handle. If we do that, however, any characters not in ASCII will just be replaced with the unknown character. Then, when we convert the bytes back to a string, the character will be replaced with the unknown character. The dangerous part about this is that there's not way to tell which character it *should* have been. That means we may have just made our data unusable!

```
# start with a string
before = "This is the euro symbol: €"

# encode it to a different encoding, replacing characters that raise errors
after = before.encode("ascii", errors = "replace")

# convert it back to utf-8
print(after.decode("ascii"))

# We've lost the original underlying byte string! It's been
# replaced with the underlying byte string for the unknown character :(
```

This is the euro symbol: ?

Reading in files with encoding problems

Most files you'll encounter will probably be encoded with UTF-8. This is what Python expects by default, so most of the time you won't run into problems. However, sometimes you'll get an error like this:

```
# try to read in a file not in UTF-8  
kickstarter_2016 = pd.read_csv("../input/kickstarter-projects/ks-projects-2016  
12.csv")
```

Notice that we get the same `UnicodeDecodeError` we got when we tried to decode UTF-8 bytes as if they were ASCII! This tells us that this file isn't actually UTF-8. We don't know what encoding it actually *is* though. One way to figure it out is to try and test a bunch of different character encodings and see if any of them work. A better way, though, is to use the `chardet` module to try and automatically guess what the right encoding is. It's not 100% guaranteed to be right, but it's usually faster than just trying to guess.

I'm going to just look at the first ten thousand bytes of this file. This is usually enough for a good guess about what the encoding is and is much faster than trying to look at the whole file. (Especially with a large file this can be very slow.) Another reason to just look at the first part of the file is that we can see by looking at the error message that the first problem is the 11th character. So we probably only need to look at the first little bit of the file to figure out what's going on.

:

```
# look at the first ten thousand bytes to guess the character encoding  
with open("../input/kickstarter-projects/ks-projects-201801.csv", 'rb') as raw  
data:  
    result = chardet.detect(rawdata.read(10000))  
  
# check what the character encoding might be  
print(result)
```

```
{'encoding': 'Windows-1252', 'confidence': 0.73, 'language': ''}
```

So chardet is 73% confidence that the right encoding is "Windows-1252". Let's see if that's correct:

```
:  
# read in the file with the encoding detected by chardet  
kickstarter_2016 = pd.read_csv("../input/kickstarter-projects/ks-projects-2016  
12.csv", encoding='Windows-1252')  
  
# look at the first few lines  
kickstarter_2016.head()
```

```
/opt/conda/lib/python3.7/site-packages/IPython/core/interactiveshell.py:344  
4: DtypeWarning: Columns (13,14,15) have mixed types.Specify dtype option o  
n import or set low_memory=False.  
    exec(code_obj, self.user_global_ns, self.user_ns)
```

Yep, looks like `chardet` was right! The file reads in with no problem (although we do get a warning about datatypes) and when we look at the first few rows it seems to be fine.

What if the encoding `chardet` guesses isn't right? Since `chardet` is basically just a fancy guesser, sometimes it will guess the wrong encoding. One thing you can try is looking at more or less of the file and seeing if you get a different result and then try that.

Saving your files with UTF-8 encoding

Finally, once you've gone through all the trouble of getting your file into UTF-8, you'll probably want to keep it that way. The easiest way to do that is to save your files with UTF-8 encoding. The good news is, since UTF-8 is the standard encoding in Python, when you save a file it will be saved as UTF-8 by default:

:

```
# save our file (will be saved as UTF-8 by default!)
kickstarter_2016.to_csv("ks-projects-201801-utf8.csv")
```


