Data Cleaning Basics: Takeaways 🖻

by Dataquest Labs, Inc. - All rights reserved © 2020

Syntax

READING A CSV IN WITH A SPECIFIC ENCODING

• Reading in a CSV file using Latin encoding:

```
laptops = pd.read_csv('laptops.csv', encoding='Latin-1')
```

• Reading in a CSV file using UTF-8:

```
laptops = pd.read_csv('laptops.csv', encoding='UTF-8')
```

• Reading in a CSV file using Windows-1251:

```
laptops = pd.read_csv('laptops.csv', encoding='Windows-1251')
```

MODIFYING COLUMNS IN A DATAFRAME

• Renaming An Existing Column:

```
laptops.rename(columns={'MANUfacturer' : 'manufacturer'}, inplace=True)
```

• Converting A String Column To Float:

```
laptops["screen_size"] = laptops["screen_size"].str.replace('"','').astype(float)
```

• Converting A String Column To Integer:

```
laptops["ram"] = laptops["ram"].str.replace('GB','')
laptops["ram"] = laptops["ram"].astype(int)
```

STRING COLUMN OPERATIONS

• Extracting Values From Strings:

FIXING VALUES

• Replacing Values Using A Mapping Dictionary:

```
mapping_dict = {
   'Android': 'Android',
   'Chrome OS': 'Chrome OS',
   'Linux': 'Linux',
   'Mac OS': 'macOS',
   'No OS': 'No OS',
   'Windows': 'Windows',
   'macOS': 'macOS'
}
laptops["os"] = laptops["os"].map(mapping_dict)
```

• Dropping Missing Values:

```
laptops_no_null_rows = laptops.dropna(axis=0)
```

EXPORTING CLEANED DATA

• Exporting Cleaned Data:

```
df.to_csv("laptops_cleaned.csv", index=False)
```

Concepts

- Computers, at their lowest levels, can only understand binary.
- Encodings are systems for representing all other values in binary so a computer can work with them.
- UTF-8 is the most common encoding and is very friendly to work with in Python 3.
- When converting text data to numeric data, we usually follow the following steps:
 - Explore the data in the column.
 - Identify patterns and special cases.
 - Remove non-digit characters.
 - Convert the column to a numeric dtype.
 - Rename column if required.

Resources

- Python Encodings
- Indexing and Selecting Data



Takeaways by Dataquest Labs, Inc. - All rights reserved $\ensuremath{\text{@}}$ 2020