

## INTRO TO DATA WRANGLING

- GATHERING DATA
- BEST PRACTICE: DOWNLOAD FILES VIA SCRIPT
- SCALABLE & REPRODUCIBLE

## - ASSESSING DATA

### - DATA QUALITY

- MISSING DATA
- INVALID DATA
- INACCURATE DATA
- INCONSISTENT DATA

FACTORS THAT AFFECT QUALITY

### - TIDINESS

- STRUCTURAL FORM OF DATASET

### - TYPES OF ASSESSMENT

- VISUAL: LOOK AT DATA IN TABLE
- PROGRAMMATIC: LOOK AT RANDAS .info() & EQUIVALENT FUNCTIONS

## - CLEANING DATA

### - PROGRAMMATIC DATA CLEANING

1. DEFINE
2. CODE
3. TEST

## - FINISH w/ REASSESS & ITERATE

## - WRANGLING vs. EDA vs. ETL

- EXPLORATORY DATA ANALYSIS (EDA): SIMPLE VISUALIZATION & ANALYSIS TO MAXIMIZE POTENTIAL OF LATER ANALYSIS
- EXTRACT TRANSFORM - LOAD PROCESS (ETL): COPYING DATA FROM ONE OR MORE SOURCES INTO SINGLE DATASET

## GATHERING DATA

### - FLAT FILES

- ENTRIES OCCUPY ONE LINE w/ SPECIFIC DELIMITERS

### - WEB SCRAPING

#### - HTML ELEMENTS

- HEADING ELEMENTS: USED FOR SECTION HEADING  
`<h1> THIS IS A HEADING. </h1>`

- PARAGRAPH ELEMENTS: USED FOR STANDARD BLOCKS OF TEXT  
`<p> THIS IS A PARAGRAPH. </p>`

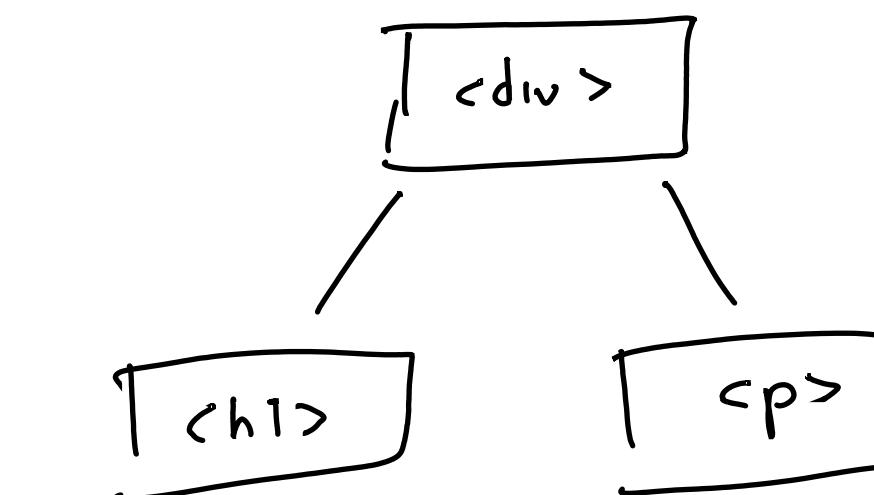
- SPAN ELEMENTS: USED TO GROUP TEXT w/IN BLOCK OF TEXT

`<p> THIS BLOCK HAS A <span> ELEMENT </span></p>`

- IMAGE ELEMENTS: USED TO EMBED IMAGES

``

### - TREES



`<div>  
 <h1> HEADING </h1>  
 <p> BODY OF TEXT </p>  
</div>`

## - HTML w/ PYTHON

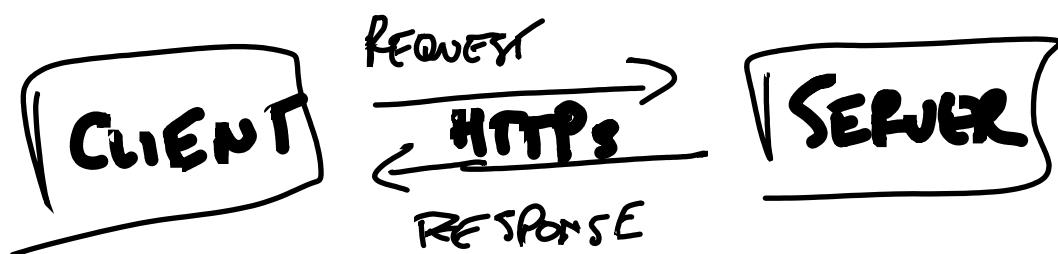
```
import requests  
response = requests.get(url)  
with open('title.html', mode='wb') as file:  
    file.write(response.content)
```

SAVE  
TO  
FILE

```
from bs4 import BeautifulSoup  
soup = BeautifulSoup(response.content, 'xml')  
soup.find('TITLE')
```

PARSE  
FILE

- DOWNLOADING FILES FROM THIS INTERNET
- HTTP IS A REQUEST / RESPONSE PROTOCOL



### - TEXT FILE STRUCTURE

- ENCODING: SCHEME FOR CONVERTING CHARACTER SET BITS TO NUMBERS AND LETTERS
- CHARACTER SETS: COLLECTION OF CHARACTERS AVAILABLE TO USE W/IN A SYSTEM

### - APPLICATION PROGRAMMING INTERFACE (API)

- JSON FILE STRUCTURE
- GREAT FOR ACCESSING COMPLICATED DATA ARCHITECTURE
- INTERPRETED AS DICTIONARIES IN PYTHON

### - STORING DATA

- RELATIONAL DATABASE STRUCTURE
- FETCHED W/ STRUCTURED QUERY LANGUAGE (SQL)
- DATABASE TABLES
  - DIVIDES INTO ROWS & COLUMNS
  - DESCRIPTIVE NAMES FOR COLUMNS, ALL SAME DATATYPE

`SELECT *  
FROM demo.orders` ] ~ READ ALL FROM DEMO ORDERS  
INTO EXPANDABLE TABLE

### - OTHER FILE FORMATS

- FLAT FILES
- HTML FILES
- JSON FILES
- TXT FILES
- EXCEL
- PICKLE
- HDFS
- SAS \*
- STATA \*

### - DATA GATHERING CAN BE ITERATIVE

## ■ ACCESSING DATA

- TWO FACTORS TO CONSIDER
- DATA QUALITY ISSUES: MISSING, DUPLICATE, OR INCORRECT
- LACK OF TIDINESS: STRUCTURAL ISSUES THAT MAY SLOW OR INHIBIT CLEANING, ANALYSIS, OR VISUALIZATION
- CASE STUDY: DIABETES DRUG PHASE II TRIAL

### - DIRTY VS MESSY DATA

- DIRTY DATA: ISSUES W/ CONTENT (i.e., LOW QUALITY)
- MESSY DATA: ISSUES W/ STRUCTURE (i.e., UNTIDY)

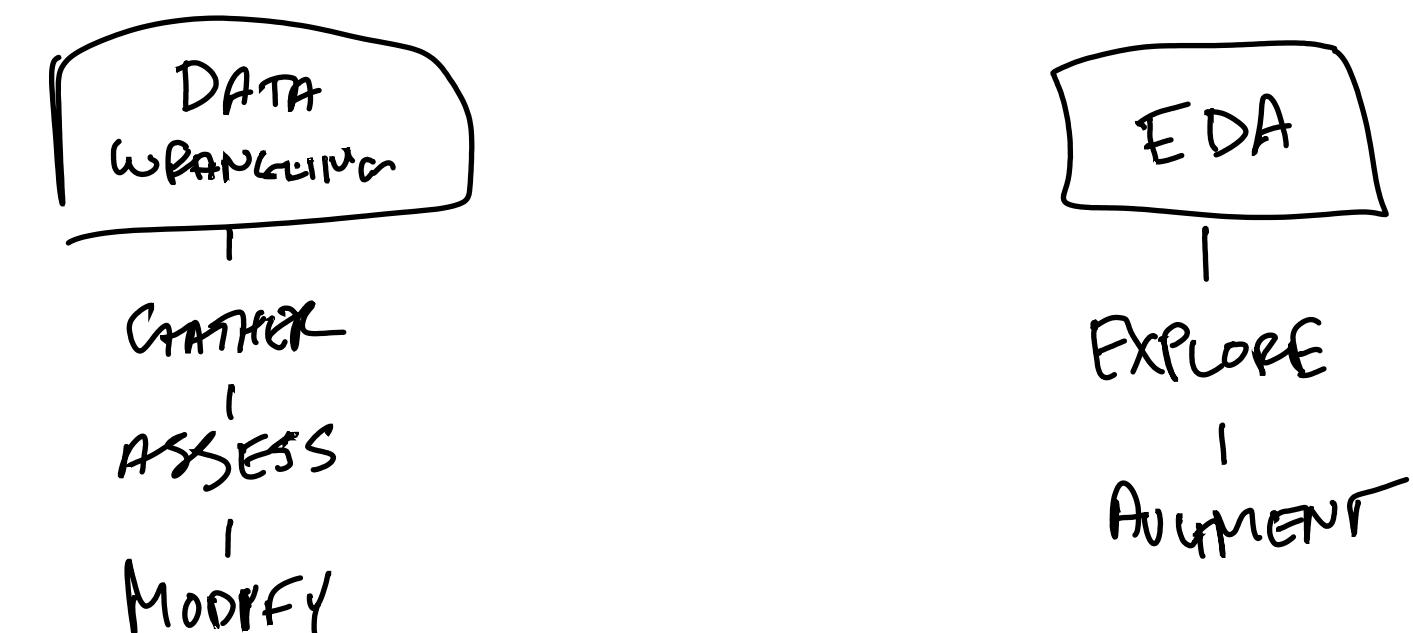
### - TYPES VS. STEPS

- TYPES
  - VISUAL
  - PROGRAMMATIC
- STEPS
  - DETECTING AN ISSUE
  - DOCUMENTING AN ISSUE

### - VISUAL ASSESSMENT

- DIRECTED: MANUALLY LOOK THROUGH EACH ROW
- NON-DIRECTED: LET'S READ THROUGH TABLE

### - ASSESSING VS EXPLORING



## - DATA QUALITY DIMENSIONS

1. COMPLETENESS
  2. VALIDITY
  3. ACCURACY
  4. CONSISTENCY
- ↑  
DIMENSIONS

## - PROGRAMMATIC ASSESSMENT

- VALUE COUNTS
- SAMPLE

## - TIDINESS: VISUAL ASSESSMENT

- EACH VARIABLE FORMS A COLUMN
- EACH OBSERVATION FORMS A ROW
- EACH TYPE OF OBSERVATION FORMS A TABLE

## ■ CLEANING DATA

### - THREE STEPS

1. DEFINE
2. CODE
3. TEST

### - DATA CLEANING PROCESS

- ① MAKE COPY
- ② ADDRESS MISSING DATA
- ③ CLEAN FOR TIDINESS
- ④ CLEAN FOR QUALITY

## ■ PROJECT OVERVIEW

- GOAL: WRANGLE WE RATE DOGS TWITTER DATA
- DATA: ENHANCED TWITTER ARCHIVE
  - TWEET TEXT
  - RATING
  - DOG NAME
  - RETWEET & FAV. COUNTS
  - IMG PREDICTED BREED
- STEP 1: GATHERING DATA