


Data 100

Lecture 5: Data Cleaning & Exploratory Data Analysis

Slides by:
 Joseph E. Gonzalez, Deb Nolan, & Joe Hellerstein
jgonzoli@berkeley.edu
deborah.nolan@berkeley.edu
hellerstein@berkeley.edu




Pandas and Jupyter Notebooks

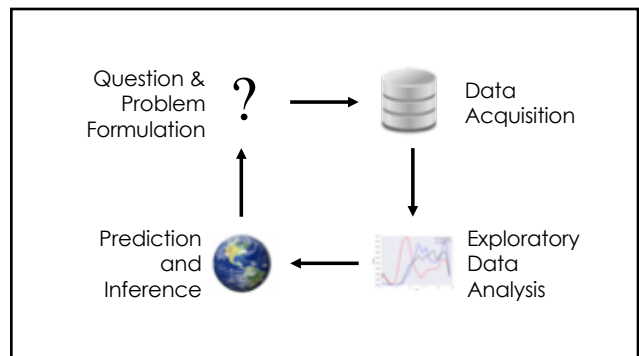
- Reviewed Jupyter Notebook Environment
- Introduced DataFrame concepts
 - **Series:** A named column of data with an index
 - **Indexes:** The mapping from keys to rows
 - **DataFrame:** collection of series with common index
- Dataframe access methods
 - **Filtering** on predicts and **slicing**
 - **df.loc:** location by index
 - **df.iloc:** location by integer address
 - **groupby** & **pivot** (we will review these again today)


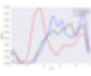


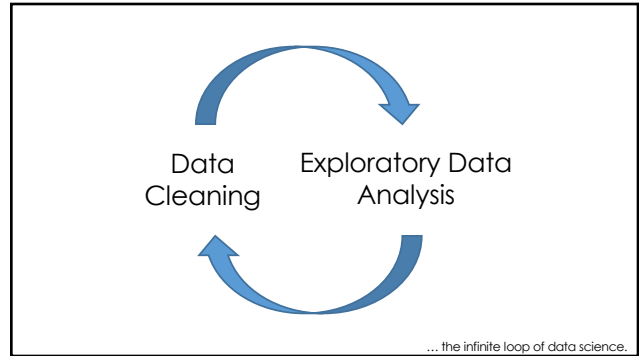
Congratulations!



You have **collected** or **been given** a box of data?
 What do you do next?




 <p>Data Acquisition</p> <p style="text-align: center;">↓</p>  <p>Exploratory Data Analysis</p>	<h2>Topics For Lecture Today</h2> <ul style="list-style-type: none"> ➤ Understanding the Data <ul style="list-style-type: none"> ➤ Data Cleaning ➤ Exploratory Data Analysis (EDA) ➤ Basic data visualization ➤ Common Data Anomalies <ul style="list-style-type: none"> ➤ ... and how to fix them
---	--

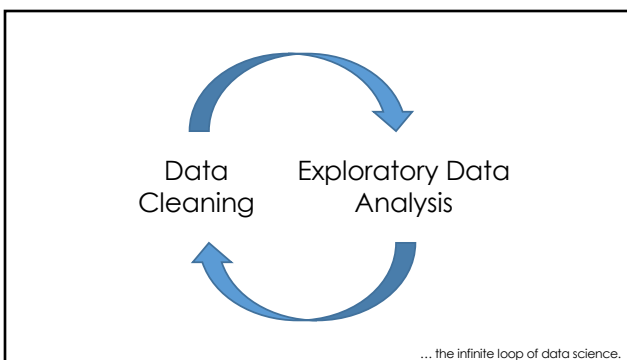


Data Cleaning

- The process of transforming raw data to facilitate subsequent analysis
- Data cleaning often addresses
 - structure / formatting
 - missing or corrupted values
 - unit conversion
 - encoding text as numbers
 - ...
- Sadly data cleaning is a big part of data science...

- Data cleaning often addresses
 - structure / formatting
 - missing or corrupted values
 - unit conversion
 - encoding text as numbers
 - ...
- Sadly data cleaning is a big part of data science...



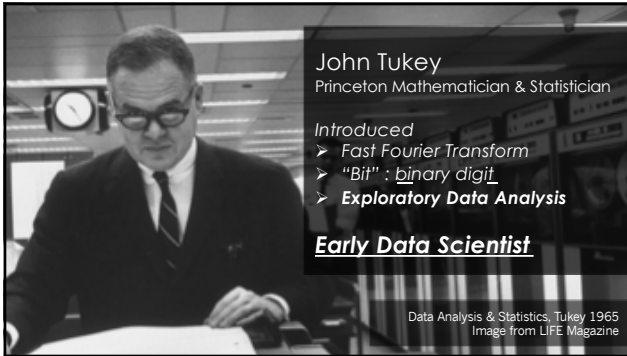


Exploratory Data Analysis (EDA)

"Getting to know the data"

The process of **transforming**, **visualizing**, and **summarizing** data to:

- Build/confirm understanding of the data and its provenance
- Identify and address potential issues in the data
- Inform the subsequent analysis
- discover *potential* hypothesis ... (be careful)
- **EDA is an open ended analysis**
 - Be willing to find something surprising



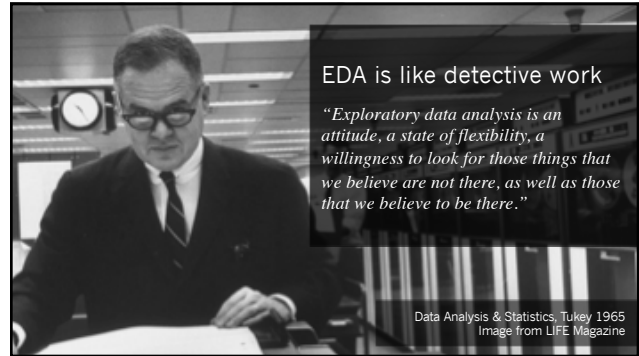
John Tukey
Princeton Mathematician & Statistician

Introduced

- Fast Fourier Transform
- "Bit" : binary digit
- **Exploratory Data Analysis**

Early Data Scientist

Data Analysis & Statistics, Tukey 1965
Image from LIFE Magazine



EDA is like detective work

"Exploratory data analysis is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those that we believe to be there."

Data Analysis & Statistics, Tukey 1965
Image from LIFE Magazine

What should we look for?

Key Data Properties to Consider in EDA

- **Structure** -- the "shape" of a data file
- **Granularity** -- how fine/coarse is each datum
- **Scope** -- how (in)complete is the data
- **Temporality** -- how is the data situated in time
- **Faithfulness** -- how well does the data capture "reality"

Key Data Properties to Consider in EDA

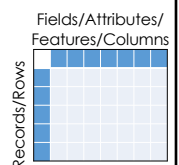
- **Structure** -- the "shape" of a data file
- **Granularity** -- how fine/coarse is each datum
- **Scope** -- how (in)complete is the data
- **Temporality** -- how is the data situated in time
- **Faithfulness** -- how well does the data capture "reality"

Rectangular Data

- We prefer rectangular data for data analysis (why?)
- Regular structures are easy to manipulate and analyze
 - A big part of data cleaning is about transforming data to be more rectangular

Two kinds of rectangular data: *Tables* and *Matrices*
(what are the differences?)

- 1. Tables** (a.k.a. data-frames in R/Python and relations in SQL)
 - Named columns with different types
 - Manipulated using data transformation languages (map, filter, group by, join, ...)
- 2. Matrices**
 - Numeric data of the same type
 - Manipulated using linear algebra



Structure: Keys

- Often data will reference other pieces of data
- **Primary key:** the column or set of columns in a table that determine the values of the remaining columns
 - Primary keys are unique
 - Examples: SSN, ProductIDs, ...
- **Foreign keys:** the column or sets of columns that reference primary keys in other tables.

Purchases.csv		
OrderNum	ProdID	Quantity
1	42	3
1	999	2
2	42	1

Orders.csv		
OrderNum	CustID	Date
1	171345	8/21/2017
2	281139	8/30/2017

Products.csv	
ProdID	Cost
42	3.14
999	2.72

Customers.csv	
CustID	Addr
171345	Harmon...
281139	Main ...



Joining two tables

OrderNum	ProdID	Name
1	42	Gum
2	999	NullFood
2	42	Towel

X

Orderid	Cust Name	Date
1	Joe	8/21/2017
2	Arthur	8/14/2017

Left "key"			Right "key"		
OrderNum	ProdID	Name	Orderid	Cust Name	Date
1	42	Gum	1	Joe	8/21/2017
1	42	Gum	2	Arthur	8/14/2017
2	999	NullFood	1	Joe	8/21/2017
2	999	NullFood	2	Arthur	8/14/2017
2	42	Towel	1	Joe	8/21/2017
2	42	Towel	2	Arthur	8/14/2017

Drop rows that don't match on the key

OrderNum	ProdID	Name	Orderid	Cust Name	Date
1	42	Gum	1	Joe	8/21/2017
2	999	NullFood	1	Joe	8/21/2017
2	999	NullFood	2	Arthur	8/14/2017
2	42	Towel	1	Joe	8/21/2017
2	42	Towel	2	Arthur	8/14/2017

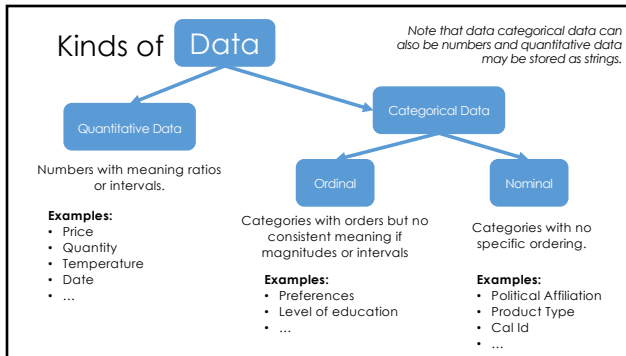
Drop rows that don't match on the key

OrderNum	ProdID	Name	Orderid	Cust Name	Date
1	42	Gum	1	Joe	8/21/2017
2	999	NullFood	2	Arthur	8/14/2017
2	42	Towel	2	Arthur	8/14/2017



Questions to ask about Structure

- Are the data in a standard format or encoding?
 - **Tabular data:** CSV, TSV, Excel, SQL
 - **Nested data:** JSON or XML
- Are the data organized in "records"?
 - No: Can we define records by parsing the data?
- Are the data nested? (records contained within records...)
 - Yes: Can we reasonably un-nest the data?
- Does the data reference other data?
 - Yes: can we join/merge the data
- What are the fields in each record?
 - How are they encoded? (e.g., strings, numbers, binary, dates ...)
 - What is the type of the data?



Quiz

<http://bit.ly/ds100-sp18-eda>

- Price in dollars of a product?
 - (A) Quantitative, (B) Ordinal, (C) Nominal
- Star Rating on Yelp?
 - (A) Quantitative, (B) Ordinal, (C) Nominal
- Date an item was sold?
 - (A) Quantitative, (B) Ordinal, (C) Nominal
- What is your Credit Card Number?
 - (A) Quantitative, (B) Ordinal, (C) Nominal

Key Data Properties to Consider in EDA

- **Structure** -- the "shape" of a data file
- **Granularity** -- how fine/coarse is each datum
- **Scope** -- how (in)complete is the data
- **Temporality** -- how is the data situated in time
- **Faithfulness** -- how well does the data capture "reality"

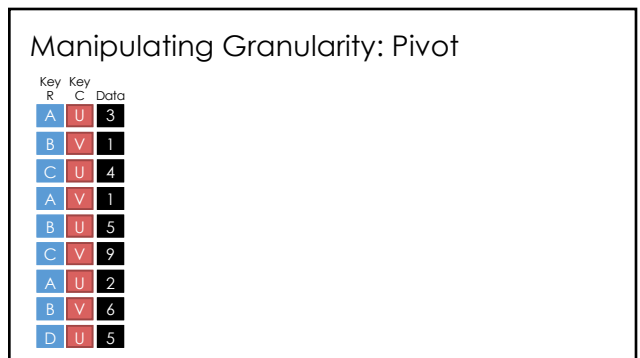
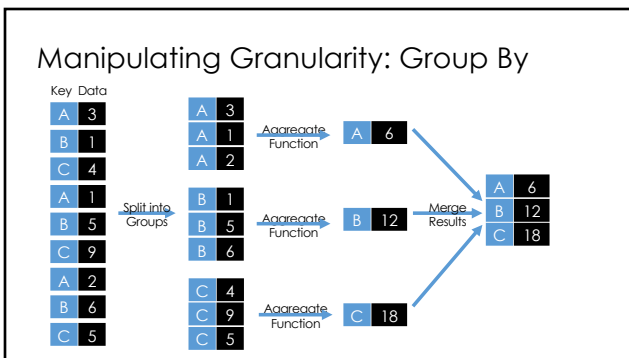
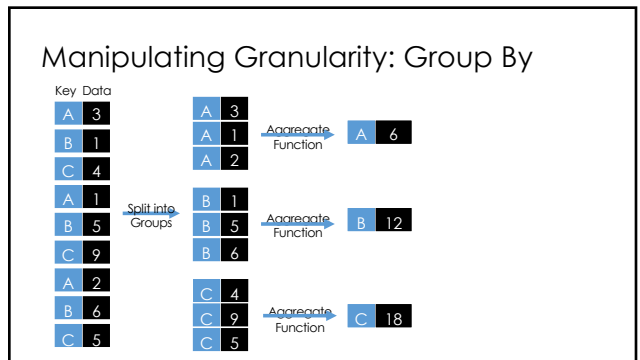
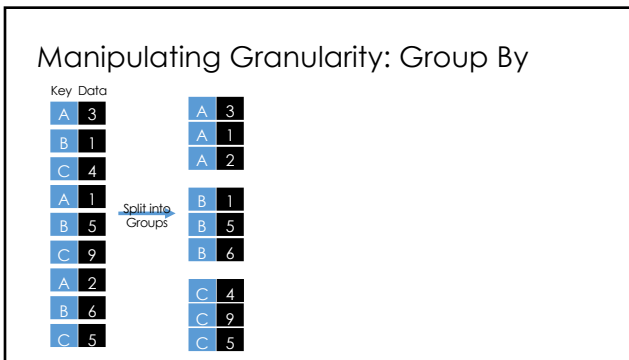
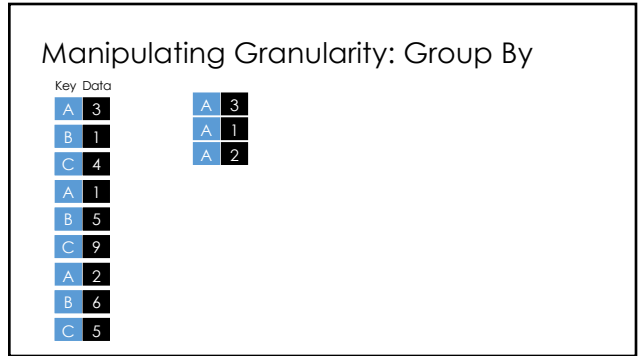
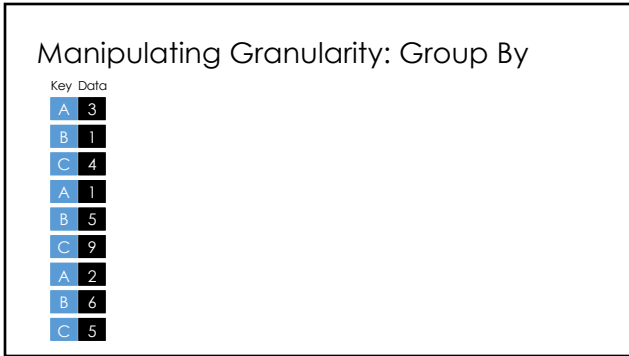
Key Data Properties to Consider in EDA

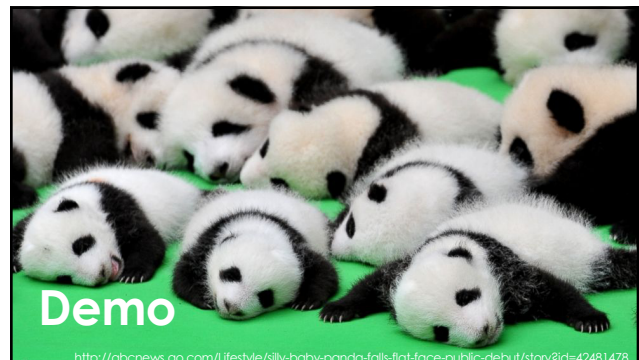
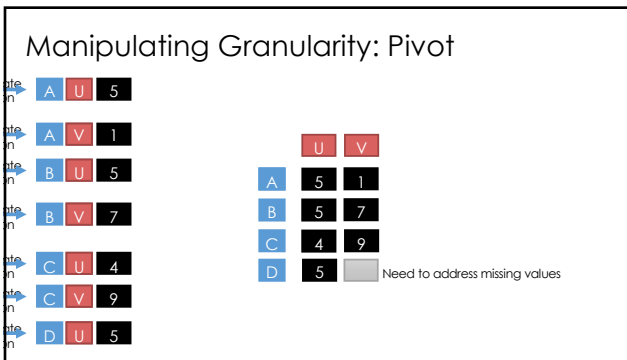
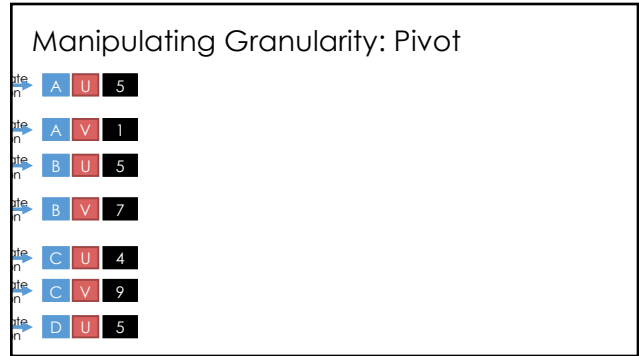
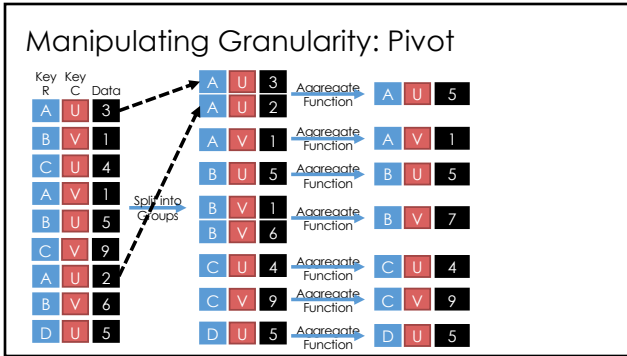
- **Structure** -- the "shape" of a data file
- **Granularity** -- how fine/coarse is each datum
- **Scope** -- how (in)complete is the data
- **Temporality** -- how is the data situated in time
- **Faithfulness** -- how well does the data capture "reality"

Granularity

- What does each record represent?
 - Examples: a purchase, a person, a group of users
- Do all records capture granularity at the same level?
 - Some data will include summaries as records
- If the data are coarse how was it aggregated?
 - Sampling, averaging, ...
- What kinds of aggregation is possible/desirable?
 - From individual people to demographic groups?
 - From individual events to totals across time or regions?
 - Hierarchies (city/county/state, second/minute/hour/days)
- Understanding and manipulating granularity can help reveal patterns.

Reviewing Group By and Pivot





- ### Key Data Properties to Consider in EDA
- > Structure -- the "shape" of a data file
 - > Granularity -- how fine/coarse is each datum
 - > Scope -- how (in)complete is the data
 - > Temporality -- how is the data situated in time
 - > Faithfulness -- how well does the data capture "reality"

- ### Key Data Properties to Consider in EDA
- > Structure -- the "shape" of a data file
 - > Granularity -- how fine/coarse is each datum
 - > Scope -- how (in)complete is the data
 - > Temporality -- how is the data situated in time
 - > Faithfulness -- how well does the data capture "reality"

Scope

- Does my data cover my area of interest?
 - **Example:** *I am interested in studying crime in California but I only have Berkeley crime data.*
- Is my data too expansive?
 - **Example:** *I am interested in student grades for DS100 but have student grades for all statistics classes.*
 - **Solution:** *Filtering → Implications on sample?*
 - *If the data is a sample I may have poor coverage after filtering ...*
- Does my data cover the right time frame?
 - *More on this in temporality ...*

To be continued ...

In the next lecture