

Data 100

Lecture 5: Data Cleaning & Exploratory Data Analysis

Slides by:
 Joseph E. Gonzalez, Deb Nolan, & Joe Hellerstein
jegonzal@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)

Today

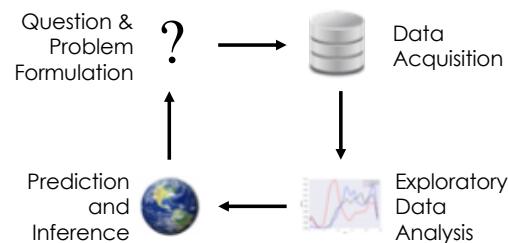


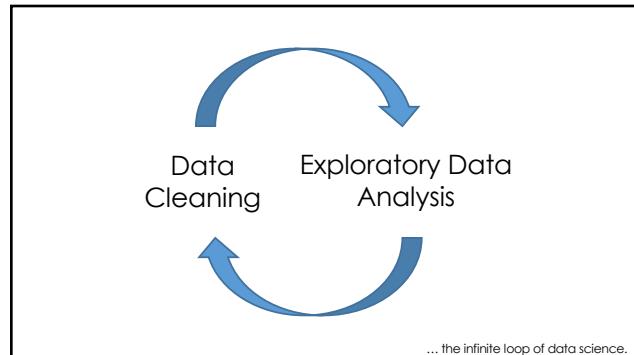
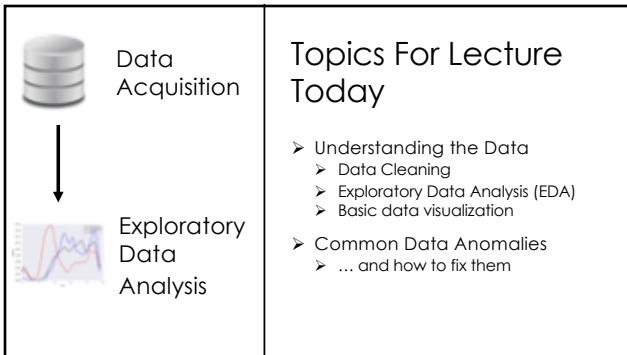
Congratulations!



You have **collected** or **been given** a box of data?

What do you do next?





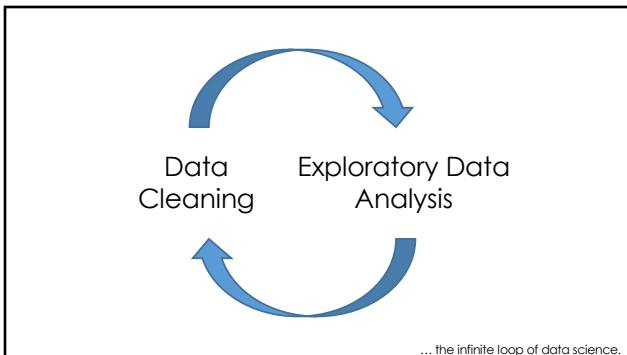
Data Cleaning

- 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...


Big Data
Borat
@BigDataBorat
Following

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

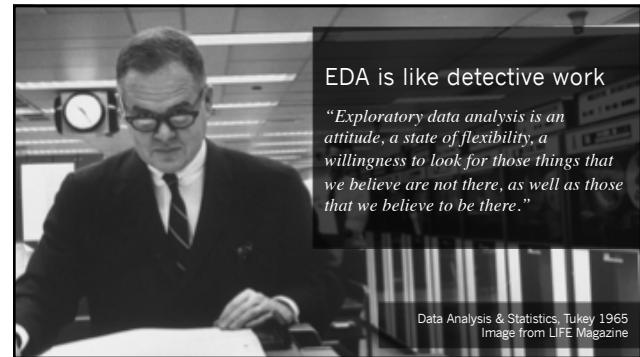
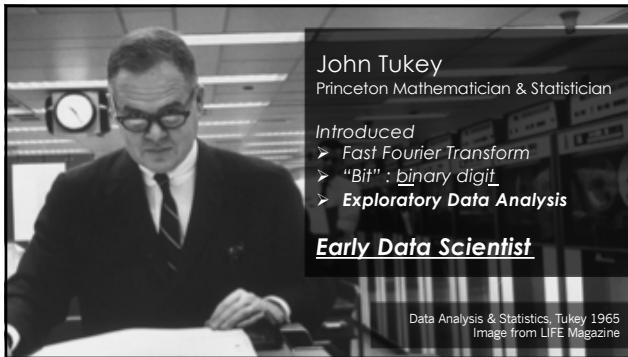


Exploratory Data Analysis (EDA)

"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



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"

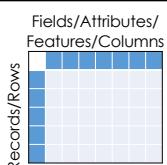
Rectangular Data

We prefer rectangular data for data analysis (why?)

- Regular structures are easy 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



How are these data files formatted?

The figure shows a terminal window with three sections side-by-side:

- TSV:** Tab-separated values. The data consists of several lines of tab-delimited text representing incidents.
- CSV:** Comma separated values. The data consists of several lines of comma-delimited text representing incidents.
- JSON:** The data is represented as a single JSON object with an array of incidents.

```

TSV
Tab separated values

CSV
Comma separated
values

JSON

```

Comma and Tab Separated Values Files

- Tabular data where
 - records are delimited by a newline: "\n", "\r\n"
 - Fields are delimited by ',' (comma) or '\t' (tab)
 - Very Common!
 - Issues?
 - Commas, tabs in records
 - Quoting
 - ...

JavaScript Object Notation (JSON)

```
1 { "a":  
2   "b":  
3     { "field1": "value1",  
4       "field2": ["list", "or", "values"],  
5       "myField3": { "is_recursive": true, "a null value": null}  
6     }  
7 }
```

- Widely used file format for nested data
 - Natural maps to python dictionaries (many tools for loading)
 - Strict formatting "quoting" addresses some issues in CSV/TSV
 - Issues
 - Each record can have different fields
 - Nesting means records can contain records → complicated

XML (another kind of nested data)

```
<catalog>
  <plant type='a'>
    <common>Bloodroot</common>
    <botanical>Sanguinaria canadensis</botanical>
    <zone>4</zone>
    <light>Mostly Shady</light>
    <price>2.44</price>
    <availability>03/15/2006</availability>
    <description>
      <color>white</color>
      <petals>true</petals>
    </description>
    <indoor>true</indoor>
  </plant>
</catalog>
```

We will study XML later in the class

Log data

Is this a csv file? tsv?
JSON/XML?

```
169.237.46.168 - - [26/Jan/2014:10:47:58 -0800] "GET /stat141/Winter04 HTTP/1.1" 301 328 "http://anson.ucdavis.edu/courses/" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0; .NET CLR 1.1.4322)"
```

169.237.6.168 - - [8/Jan/2014:10:47:58 -0800] "GET /stat141/Winter04/ HTTP/1.1" 200 2595 "<http://anson.ucavis.edu/courses/>" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0; .NET CLR 1.1.4322)"

Data can be **split across files** and **reference other data**.

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 ..

Merging/joining data across tables

Joining two tables

X

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			

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

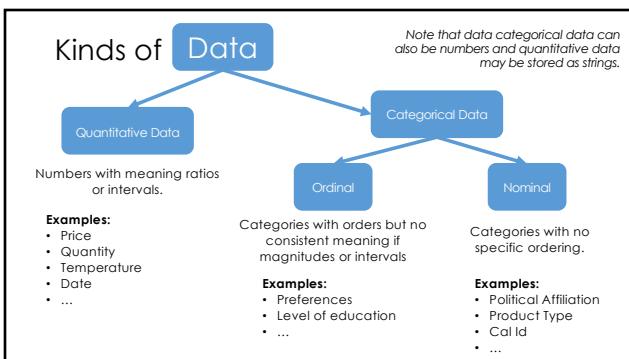


Pandas Merge Function

Demo

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

Manipulating Granularity: Group By

Key Data	
A	3
B	1
C	4
A	1
B	5
C	9
A	2
B	6
C	5

Manipulating Granularity: Group By

Key Data	
A	3
B	1
C	4
A	1
B	5
C	9
A	2
B	6
C	5

Manipulating Granularity: Group By

Key Data	
A	3
B	1
C	4
A	1
B	5
C	9
A	2
B	6
C	5

Manipulating Granularity: Group By

Key Data	
A	3
B	1
C	4
A	1
B	5
C	9
A	2
B	6
C	5

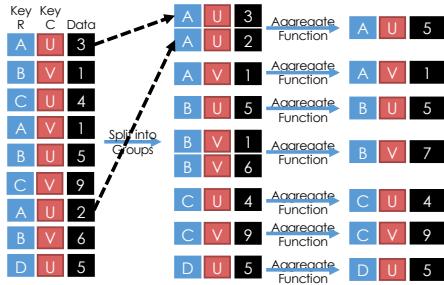
Manipulating Granularity: Group By

Key Data	
A	3
B	1
C	4
A	1
B	5
C	9
A	2
B	6
C	5

Manipulating Granularity: Pivot

Key R	Key C	Data
A	U	3
B	V	1
C	U	4
A	V	1
B	U	5
C	V	9
A	U	2
B	V	6
D	U	5

Manipulating Granularity: Pivot



Manipulating Granularity: Pivot

Aggregate Function	A	U	5
Aggregate Function	A	V	1
Aggregate Function	B	U	5
Aggregate Function	B	V	7
Aggregate Function	C	U	4
Aggregate Function	C	V	9
Aggregate Function	D	U	5

Manipulating Granularity: Pivot

Aggregate Function

A	U	5
A	V	1
B	U	5
B	V	7
C	U	4
C	V	9
D	U	5

U V

A	5	1
B	5	7
C	4	9
D	5	

Need to address missing values



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 ...

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"

Temporality

- What is the meaning of the time and date fields?
 - When the "event" happened?
 - When the data was collected or was entered into the system?
- Time depends on where? (Time zones & daylight savings)
 - Learn to use **datetime** python library
- Multiple string representation (depends on region): 08/08/08?
- Are there strange null values?
 - January 1st 1970, January 1st 1900
 - Date the data was copied into a database (look for many matching timestamps)
- Is there periodicity? Diurnal patterns

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"

Faithfulness: Do I trust this data?

- Does my data contain unrealistic or "incorrect" values?
 - Examples?
 - Dates in the future for events in the past
 - Locations that don't exist
 - Negative counts
 - Misspellings of names
 - Large outliers
- Does my data violate obvious dependencies?
 - E.g., age and birthday don't match
- Was the data entered by hand?
 - Spelling errors, fields shifted ...
 - Did the form require fields or provide default values?
- Are there obvious signs of curb stoning (data falsification)?
 - Repeated names, fake looking email addresses, repeated use of uncommon names or fields.

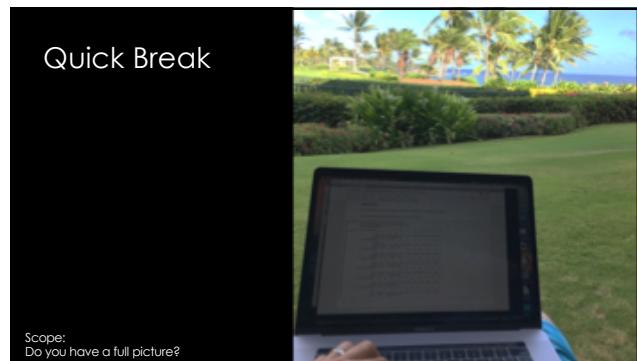
Signs that your data may not be faithful

- Missing Values/Default values: [0, -1, 999, 12345, NaN, Null, 1970, 1900, ... others?]
 - **Soln 1:** Drop records with missing values → implications on your sample!
 - **Soln 2:** Impute missing values → Bias your conclusions
- Time Zone Inconsistencies
 - **Soln 1:** convert to a common timezone (e.g., UTC)
 - **Soln 2:** convert to the timezone of the location – useful in modeling behavior.
- Duplicated Records or Fields
 - **Soln:** identify and eliminate (use primary key) → implications on sample?
- Spelling Errors
 - **Soln:** Apply corrections or drop records not in a dictionary → implications on sample?
- Units not specified or consistent
 - **Soln:** Infer units, check values are in reasonable ranges for data
- Truncated data (early excel limits: 65536 Rows, 255 Columns)
 - **Soln:** be aware of consequences in analysis → how did truncation affect sample?
- Others...

Quick Break



Quick Break



Berkeley Police Data Demo

Berkeley Police Public Datasets

- **Question:** For this analysis we will not begin with a detailed question but instead a rough goal of understanding Police activity.
- **Examine Two Data Sets:**
 - Call data
 - Stop data
- Today we will work through the basic process of data loading, some preliminary cleaning, and exploratory data analysis.

Call Data Description

Data pulled from Public Safety Server using data created for Berkeley's Crime View Community page. Displays **incidents reported** for the last 180 days along with **time, date, day of week** and **block level location information**.

The dataset reflects crimes as they have been reported to the BPD based on preliminary information **supplied by the reporting parties**. Preliminary crime classifications may change based on follow-up investigations. **Not all calls for police service are included (e.g.**

Animal Bite). The information provided on this site is intended for use by the community to enhance their awareness of crimes occurring in their neighborhoods and the entire City. **The data should not be used for in-depth crime analysis** as the initial information is subject to change.

Stops Data Description

This data was extracted from the Department's Public Safety Server and covers the **data beginning January 26, 2015**. On January 26, 2015 the department began collecting data pursuant to General Order B-4 (issued December 31, 2014). Under that order, **officers were required to provide certain data after making all vehicle detentions** (including bicycles) and pedestrian detentions (up to five persons). This data **set lists stops by police** in the categories of traffic, suspicious vehicle, pedestrian and bicycle stops. Incident number, date and time, location and disposition codes are also listed in this data.

Address **data has been changed from a specific address**, where applicable, and listed as the block where the incident occurred. Disposition codes were entered by officers who made the stop. These codes included the person(s) race, gender, age (range), reason for the stop, enforcement action taken, and whether or not a search was conducted.

Visualizing Univariate Relationships

➤ Quantitative Data

- Histograms, Box Plots, Rug Plots, Smoothed Interpolations (KDE – Kernel Density Estimators)
- Look for spread, shape, modes, outliers, unreasonable values ...

➤ Nominal & Ordinal Data

- Bar plots (sorted by frequency or ordinal dimension)
- Look for skew, frequent and rare categories, or invalid categories
- Consider grouping categories and repeating analysis

Histograms, Rug Plots, and KDE Interpolation

Describes distribution of data – relative prevalence of values

➤ Histogram

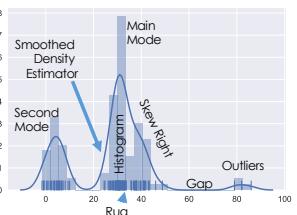
- relative frequency of values in bins (ranges)
- Tradeoff of bin sizes

➤ Rug Plot

- Shows the actual data locations

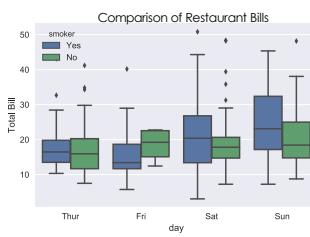
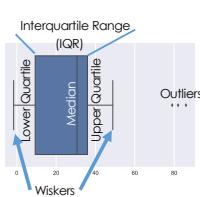
➤ Smoothed density estimator

- Tradeoff of "bandwidth" parameter (more on this later)



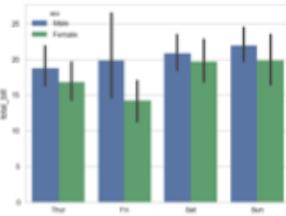
Box Charts

- Useful for summarizing distributions and comparing multiple distributions



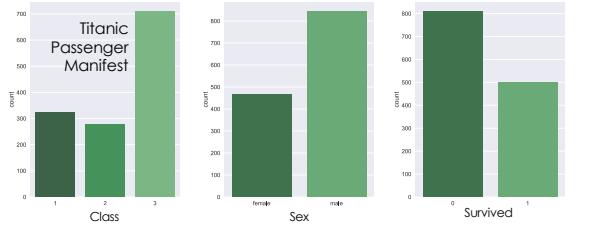
Bar Charts

- Used to compare nominal and ordinal data.
- Consider sorting by category or frequency



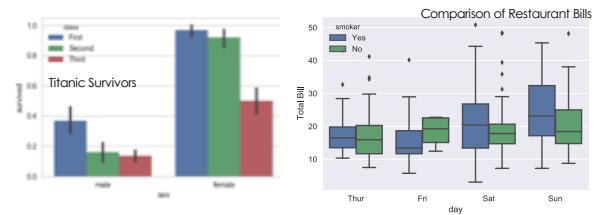
Bar Charts

- Used to compare nominal and ordinal data.
- Consider sorting by category or frequency



Visualizing Multivariate Relationships

- Conditioning on a range of values (e.g., ages in groups) and construct side by side box-plots or bar charts



Visualizing Multivariate Relationships

- Scatter Plots: try plotting variables against each other
 - Try to linearize relationships (e.g., logs, exponents, square-roots)
 - More on transformations when we return to visualizations
- Conditioning on a range of values (e.g., ages in groups) and construct side by side box-plots or bar charts

Caution about EDA

With enough data, if you look hard enough you will find something "**interesting**"

Important to differentiate **inferential conclusions** about world from **exploratory analysis of data**

