# 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

---

## Last Lecture

➤ Started discussing exploratory data analysis

➤ **Structure** -- *the "shape" of a data file (how is it organized)*



---

## Last Lecture

➤ Started discussing exploratory data analysis

➤ **Structure** -- *the "shape" of a data file (how is it organized)*

➤ ***Granularity*** *-- how fine/coarse is each datum*



Fine Grained — Coarse Grained

---

## **Group By** – manipulating granularity



---

## **Pivot** – A kind of Group By Operation



Need to address missing values

---

## Last Lecture

➤ Started discussing exploratory data analysis

➤ **Structure** -- *the "shape" of a data file (how is it organized)*

➤ ***Granularity*** *-- how fine/coarse is each datum*

➤ ***Scope*** *-- how (in)complete is the data*



Need to Filter — Population — Data

Need to Filter — Data — Population — Need more data

## Last Lecture

➢ Started discussing exploratory data analysis

➢ **Structure** -- *the "shape" of a data file (how is it organized)*

➢ *Granularity -- how fine/coarse is each datum*

➢ *Scope -- how (in)complete is the data*

➢ *Temporality -- how is the data situated in time*

## Temporality

➢ Data changes → When was the data collected!

➢ What is the meaning of a the time and date fields?
  ➢ When the "event" **happened**?
  ➢ When the data was **collected** or was **entered** into the system?
  ➢ Date the data was copied into a database (look for many matching timestamps)

➢ Time depends on where! (Time zones & daylight savings)
  ➢ Learn to use **datetime** python library
  ➢ Multiple string representation (depends on region): 07/08/09?

➢ Are there strange null values?
  ➢ January 1st 1970, January 1st 1900

➢ Is there periodicity? Diurnal patterns

## Unix Time / POSIX Time

➢ Time **measured in seconds** since January 1st 1970
  ➢ Minus leap seconds …

➢ Unix time follows Coordinated Universal Time (UTC)
  ➢ International time standard
  ➢ Measured at 0 degrees latitude
    ➢ Similar to Greenwich Mean Time (GMT)
  ➢ No daylight savings
  ➢ Time codes

➢ Time Zones:
  ➢ San Francisco (UTC-8) without daylight savings

https://en.wikipedia.org/wiki/Coordinated_Universal_Time

## 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"*
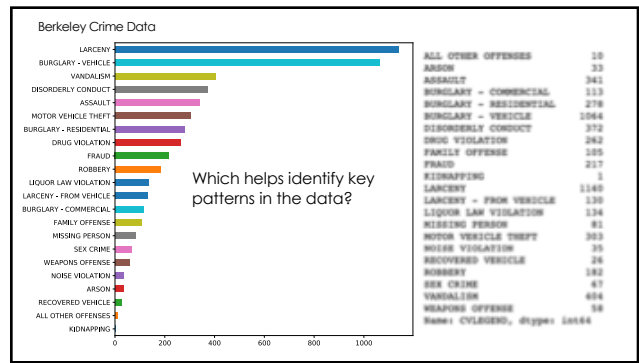
## Faithfulness: *Do I trust this data?*

➢ 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
  - **Solns:** 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…

## How do you do EDA?

- Examine data and meta-data:
  - What is the date, size, organization, and structure of the data?
- Examine each field/attribute/dimension individually
- Examine pairs of related dimensions
  - Stratifying earlier analysis: break down grades by major …
- Along the way:
  - Visualize/summarize the data
  - Validate assumptions about data and collection process
  - Identify and address anomalies
  - Apply data transformations and corrections
  - ***Record everything you do! (why?)***

# Visualization and EDA

Berkeley Crime Data

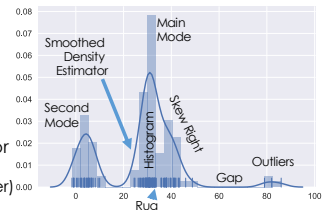Which helps identify key patterns in the data?

## 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 oridinal dimension)
  - Look for skew, frequent and rare categories, or invalid categories
  - Consider grouping categories and repeating analysis

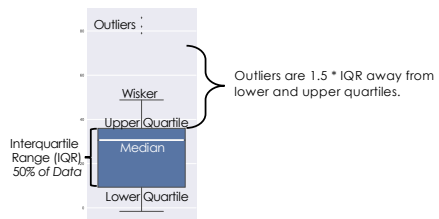## Histograms, Rug Plots, and KDE Interpolation

Describes distribution of data – relative prevalence of values

- Histogram
  - relative frequency of values
  - Tradeoff of bin sizes
- Rug Plot
  - Shows the actual data locations
- Smoothed density estimator
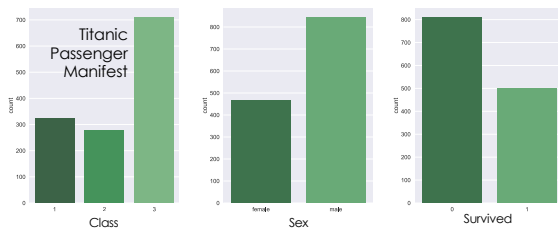  - Tradeoff of "bandwidth" parameter (more on this later)

## Box Charts

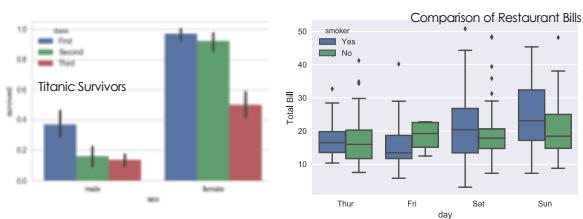➤ Useful for summarizing distributions and comparing multiple distributions



Outliers are 1.5 * IQR away from lower and upper quartiles.

## Bar Charts

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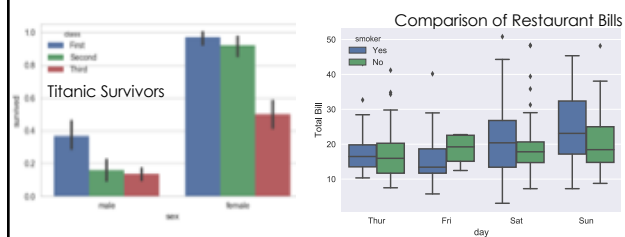  ➤ Consider sorting by category or frequency



## Visualizing Multivariate Relationships

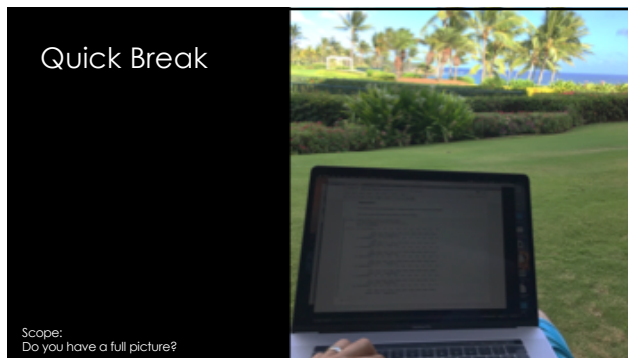➤ Conditioning on a range of values (e.g., ages in groups) and construct side by side box-plots or bar charts



## http://bit.ly/ds100-sp18-eda





Quick Break



Quick Break

Scope:
Do you have a full picture?

4

# 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.

## 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**



U.C. BERKELEY STATISTICS

If you torture the data enough, it may confess