# Lazily Looking at pandas DataFrames

Eric_A_Moore@yahoo.com

8/21/18

# "The Problem"

- See what's in a DataFrame, toward identifying interesting relationships among variables
  - Concisely
  - Clearly
  - Don't overwhelm
  - Don't lie
- Code in order to learn: pandas, Altair, Vega
- Initial application: how much do athlete's bodies predict their position? (baseball, football, Olympics)

# Key Ideas

Phases of analyses, series of views

Simple schema – permute over

Simple invocation

Embody visual judgment

Two modes:
- **Explore**: Design for screen first (over paper) – sequence, scroll, page, gesture
- **Explain**: Legible static image, highlight the "tell"

# Phases of Analyses

1. ***Data***: Quality, Quantity, Qualify.
   - Formats; gaps, causes for gaps; need for transformation.
2. ***Variables***: what is potentially informative – similar, different, distinct?
   - See Variables – distribution of values
   - See Variable – Variable relationships
3. ***Model***: what is predictive?; what are useful abstractions?
   - See relationships in context of a Target Variable
   - Measures on Measures – toward statistical inference
   - "Fit" or abstract variable distributions and relationships
   - "Focus" a predictive model to most essential form

# Simple Schema

Statistics

- Compare effects of a process across subpopulations
- Effects are measured, estimated given samples of a population
- *Category* – divides population
- *Measure* -- describes population

Data

- Symbolic -- *Categorical, Nominal, Ordinal, Key*
- Numeric -- *Quantitative, Temporal*
- *Set* - collection of records ~ Sample ~ Population
- Field ~ Variable ~ Feature

# Embody visual judgment – describe, compare

| Subject | Chart | Comments |
|---------|-------|----------|
| Categorical 1D | Horizontal Histogram | Group by category. |
| Scalar 1D | Vertical Histogram | Binned X.<br>Fit later w/ parametric Distribution |
| Category vs Category | Crosstab Table | Counts |
| Measure vs Measure | Scatter Plot | Fit later w/ Loess or regression. |
| Measure over Category | Hozo Bar Chart | Measure, conditioned on each category. |
| Measures on Measure over Category | Hozo Box Chart | Using mean, std deviation vs median and quartiles.<br>"Separation" in 1D, with one measure. |
| Category over Measures | Scatter Plot with Category | Separation in 2D, with two joint measures.<br>Only compare 2 or 3 categories at once. |
| Category over Category differences | Crosstab Table with residual density | Highlight non uniformity in joint distribution |
| Set of Correlations | Hozo Bar Chart (Correlations) | Highlight strong or weak correlations |

# Simple Schema

```
view_phases = [ 'Data', 'Variables', 'Model' ]

view_schema = {
    'Data' : [] ,
    'Variables' : [ '1D', '2D' ],
    'Model' : [ 'CategoryCategory', 'CategoryMeasure', 'Correlations' ] }

schema = {
    'Collection' : 'Baseball',
    'Measures' : [ 'Age', 'BMI', 'Height', 'Weight' ],
    'Categories' : ['Position' , 'Team'],
    'Target' : [ 'Position' ]
}
```
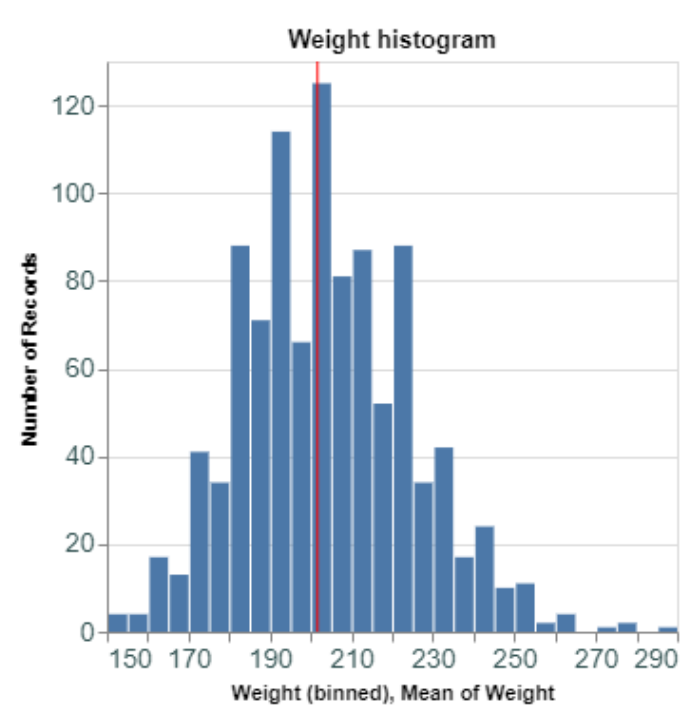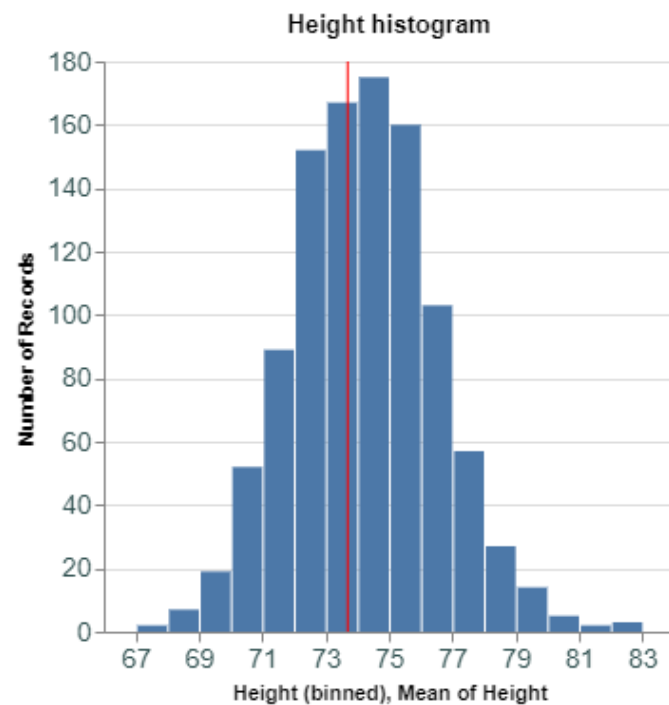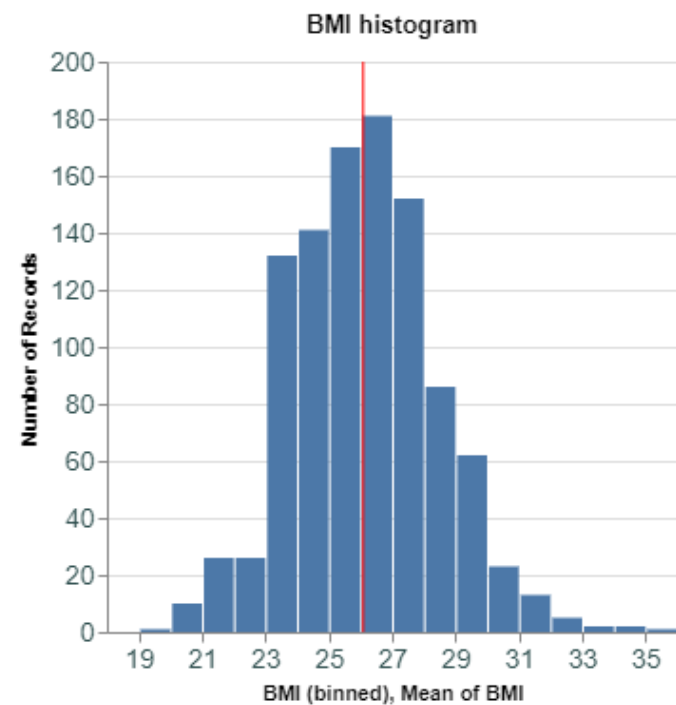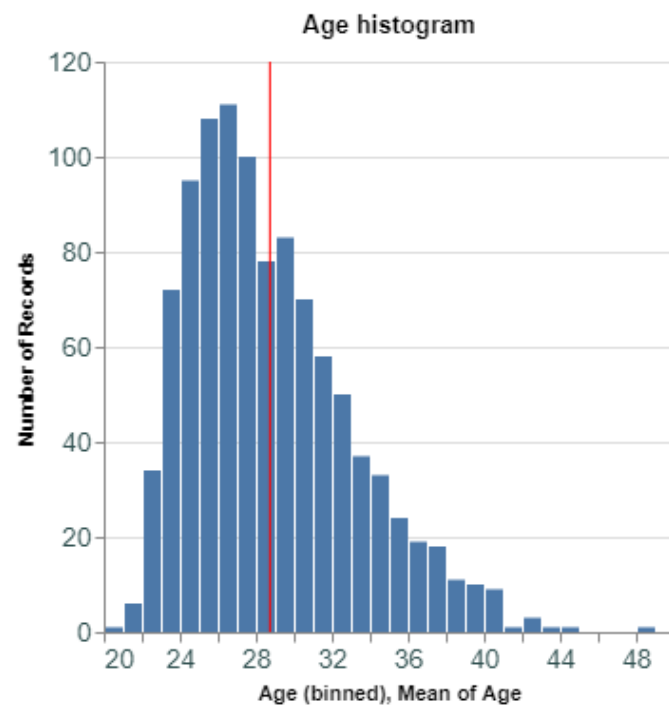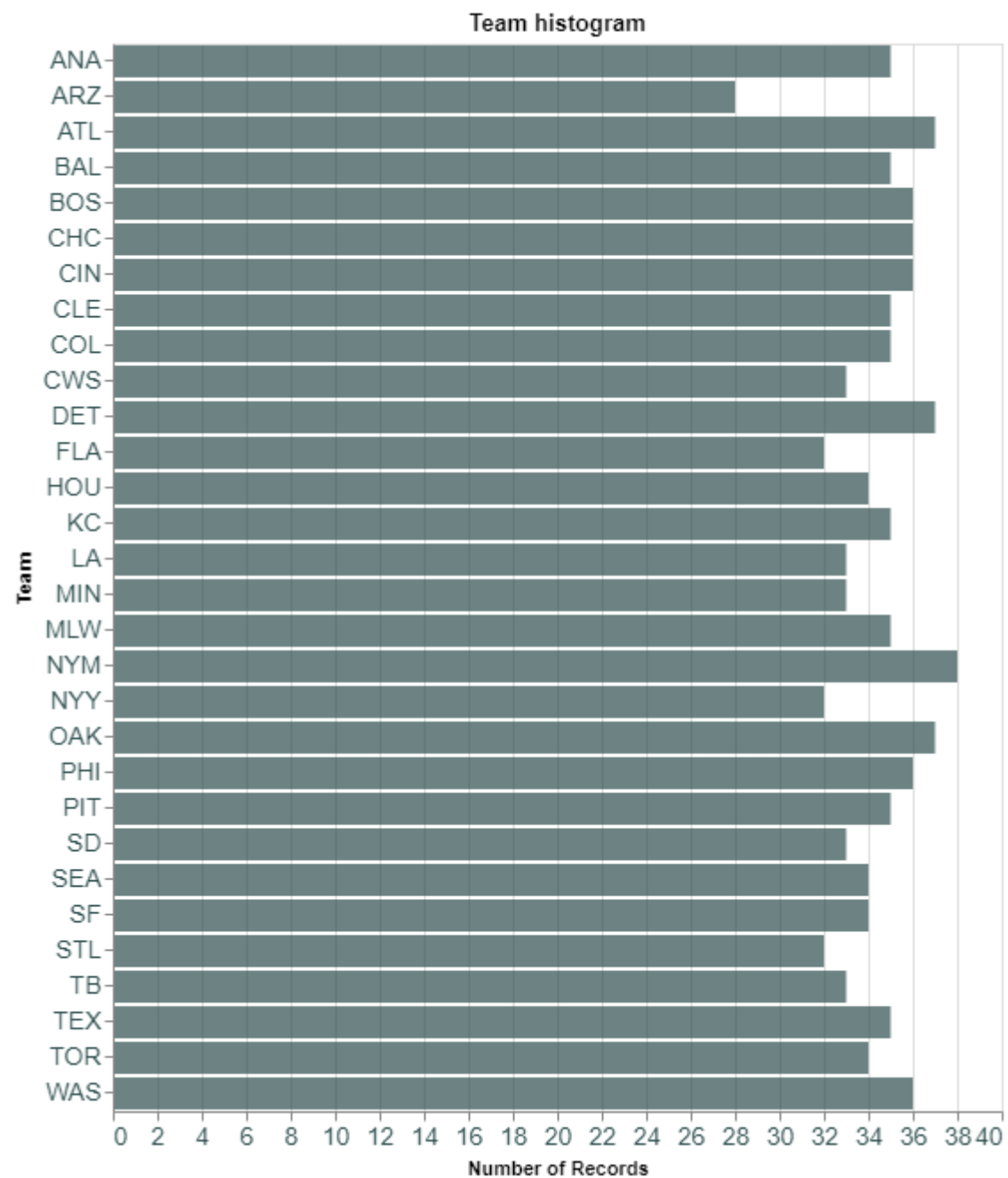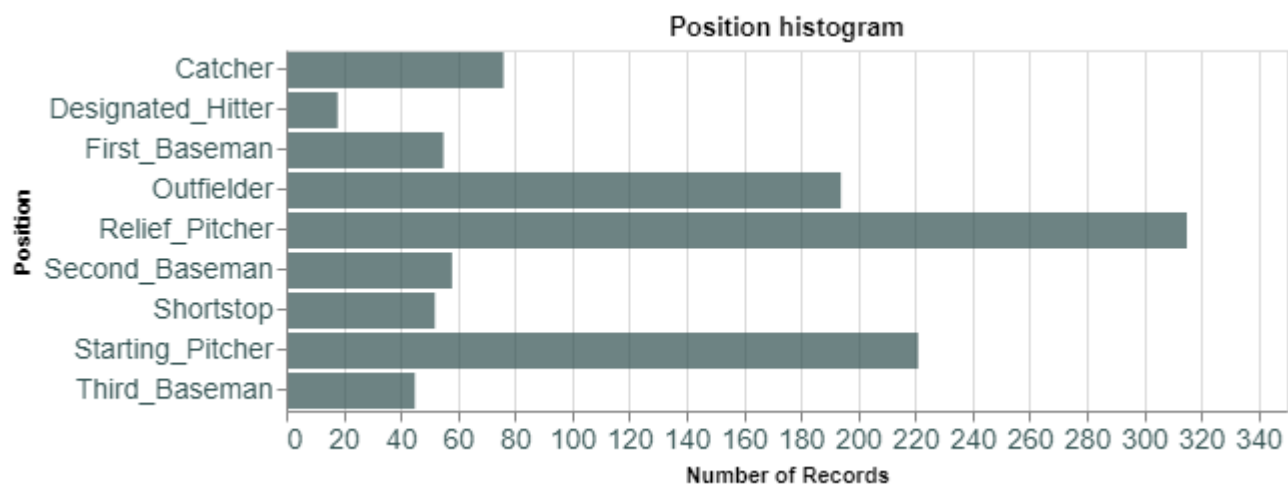
# Simple invocation

```python
schema['Sport'] = 'Baseball'
df_baseball = load_baseball('./baseball2016.csv')
df_baseball_summary = make_summary_df(df_baseball, schema)
view_all(df_baseball, df_baseball_summary, schema, phases, view_schema)
```
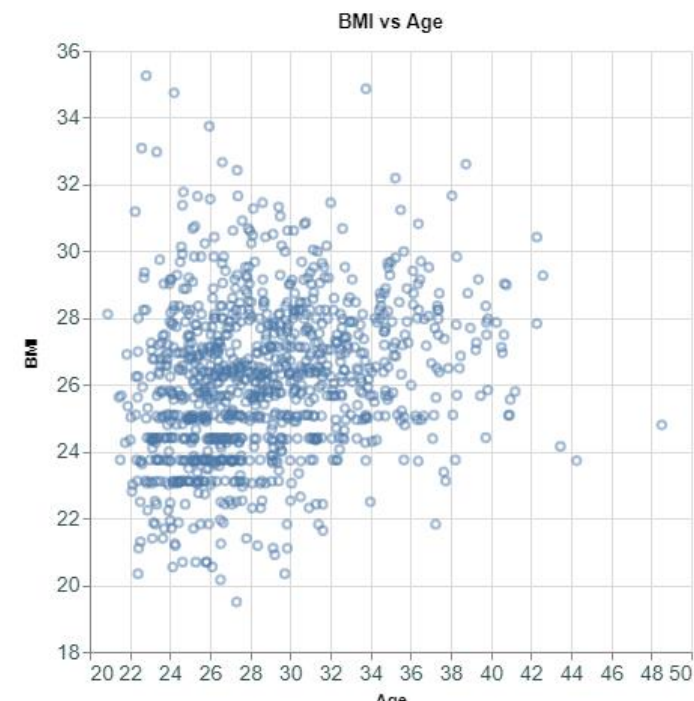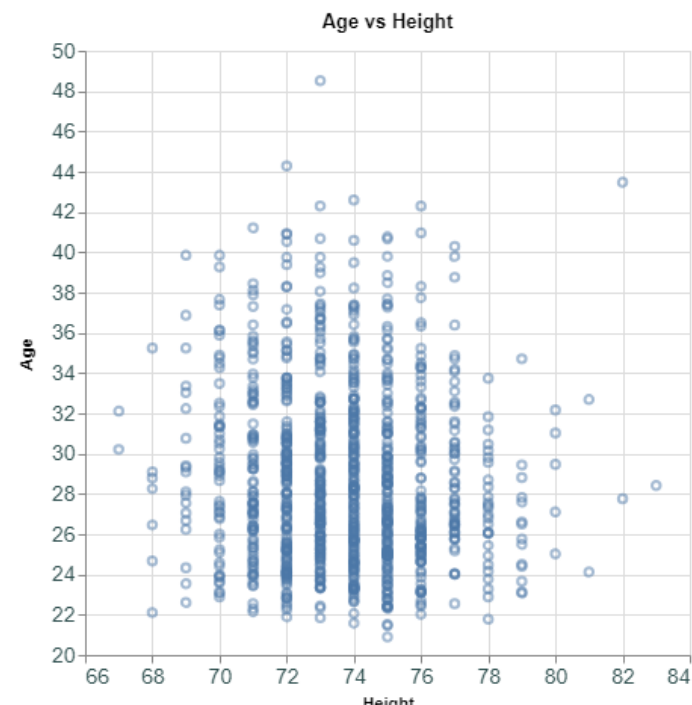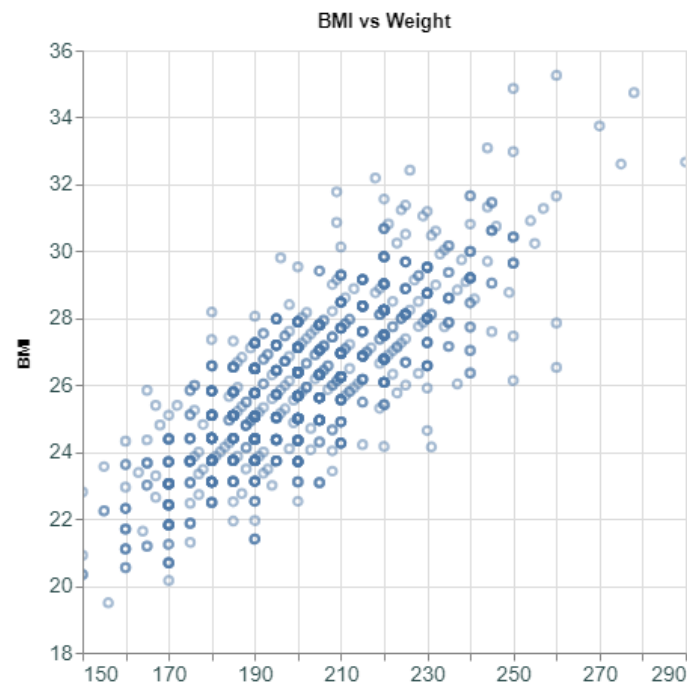
```python
def load_baseball(file):
    df = pd.read_csv(file)
    # format df.
    baseball_columns_map = {
        'Position': 'Position',
        'Height(inches)': 'Height',
        'Weight(pounds)': 'Weight',
        'Age' : 'Age'
    }

    df = df.drop('Name', axis=1)
    df.rename(columns=baseball_columns_map,inplace=True)

    df['BMI'] = calculate_BMI(df['Height'], df['Weight'])
    return df
```

| Age histogram | BMI histogram |
| Height histogram | Weight histogram |

Position histogram

Team histogram

## Position x Team

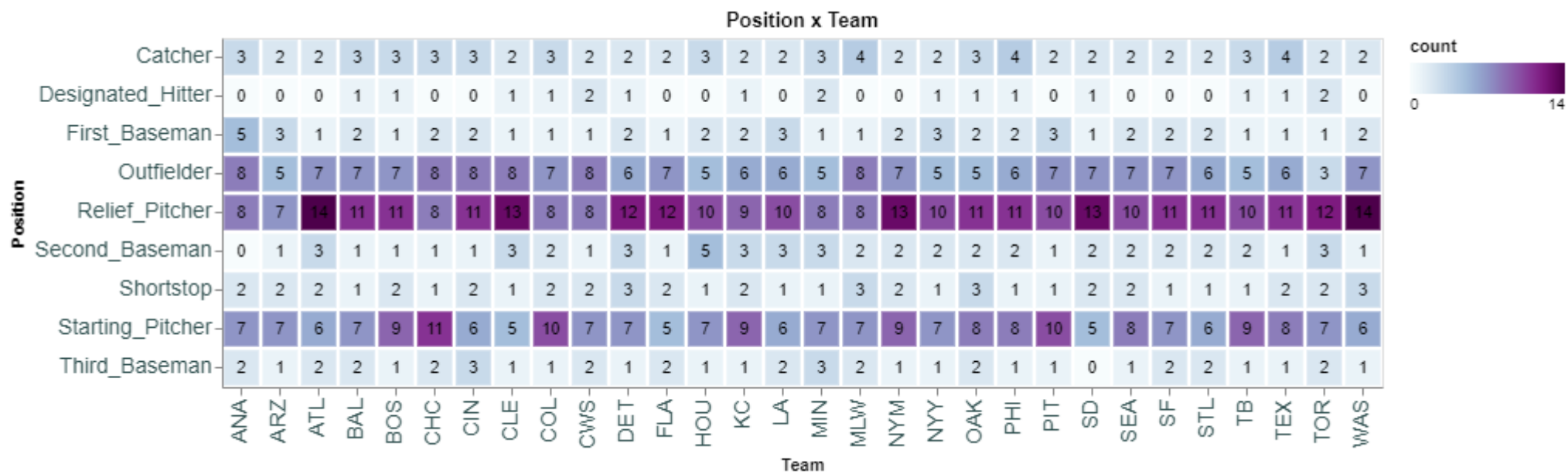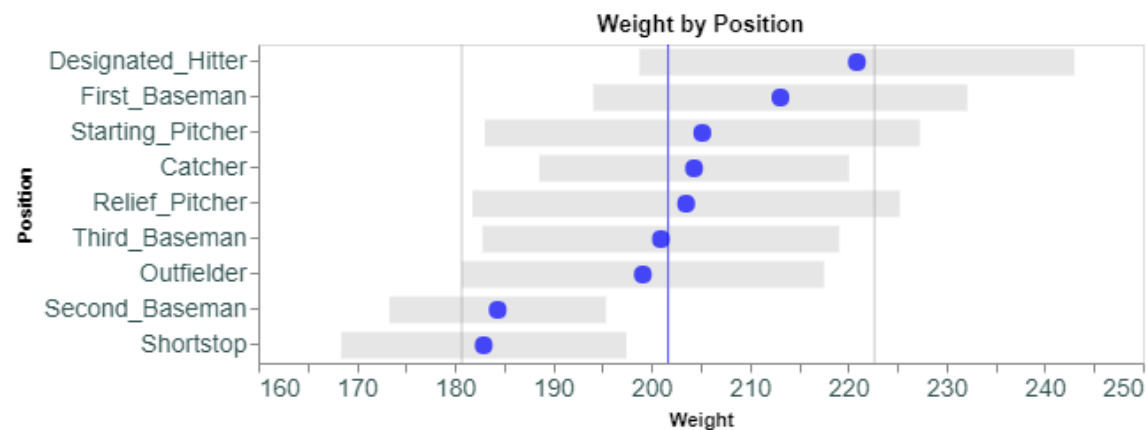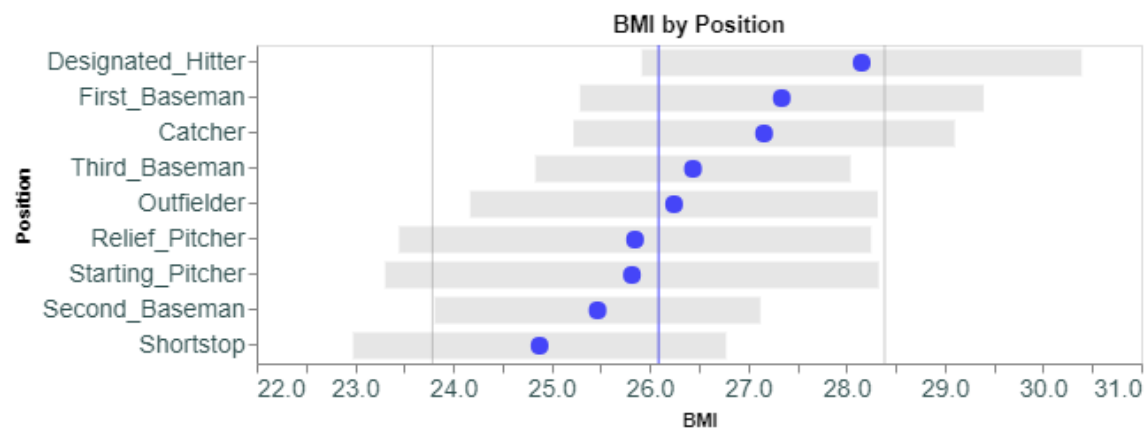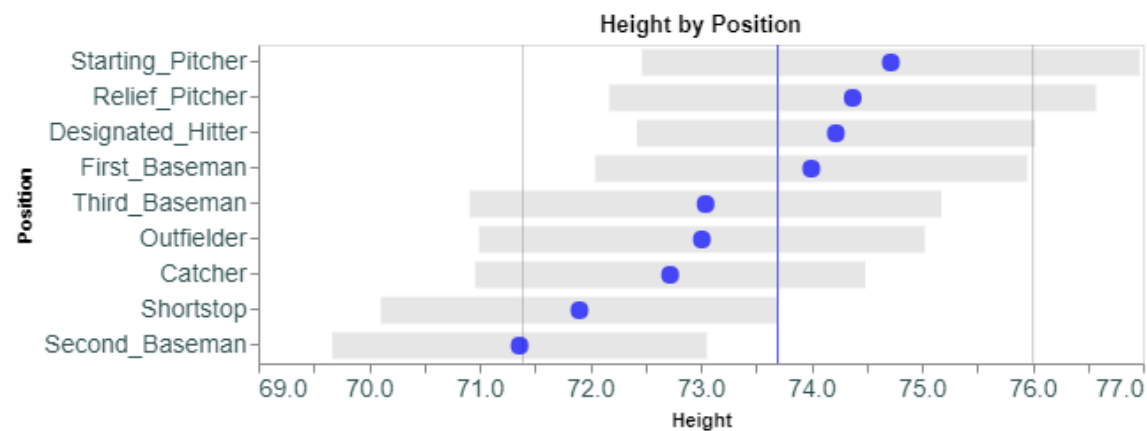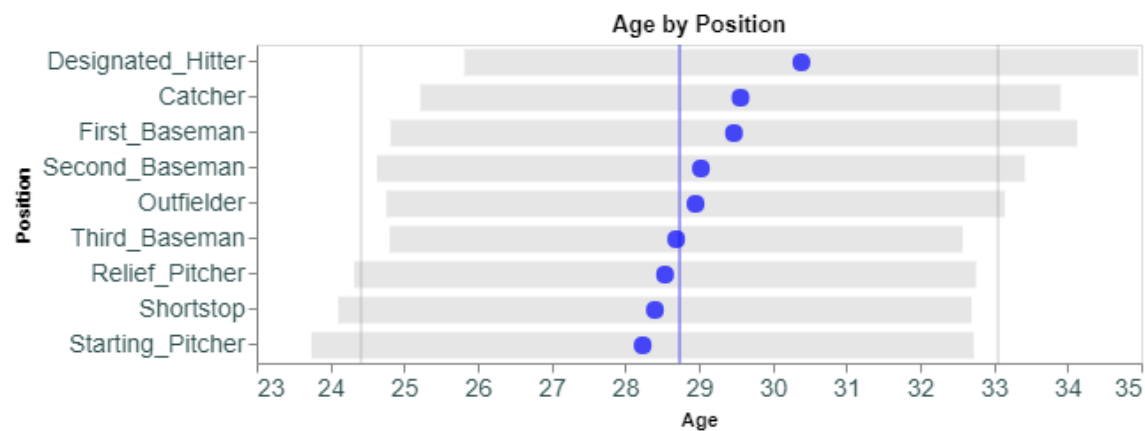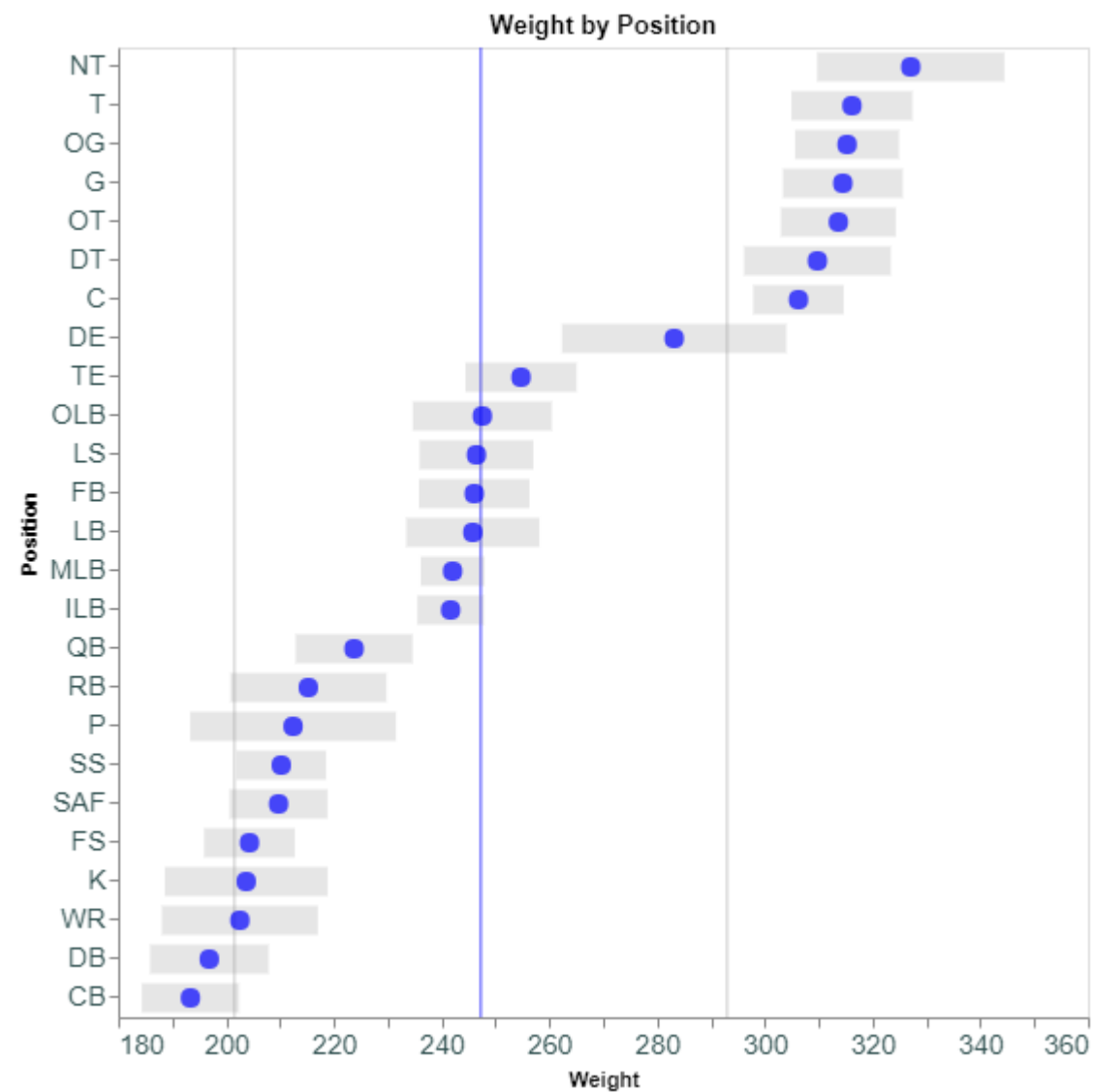| Position | ANA | ARZ | ATL | BAL | BOS | CHC | CIN | CLE | COL | CWS | DET | FLA | HOU | KC | LA | MIN | MLW | NYM | NYY | OAK | PHI | PIT | SD | SEA | SF | STL | TB | TEX | TOR | WAS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Catcher | 3 | 2 | 2 | 3 | 3 | 3 | 3 | 2 | 3 | 2 | 2 | 2 | 3 | 2 | 2 | 3 | 4 | 2 | 2 | 3 | 4 | 2 | 2 | 2 | 2 | 2 | 3 | 4 | 2 | 2 |
| Designated_Hitter | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 2 | 1 | 0 | 0 | 1 | 0 | 2 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 2 | 0 |
| First_Baseman | 5 | 3 | 1 | 2 | 1 | 2 | 2 | 1 | 1 | 1 | 2 | 1 | 2 | 2 | 3 | 1 | 1 | 2 | 3 | 2 | 2 | 3 | 1 | 2 | 2 | 2 | 1 | 1 | 1 | 2 |
| Outfielder | 8 | 5 | 7 | 7 | 7 | 8 | 8 | 8 | 7 | 8 | 6 | 7 | 5 | 6 | 6 | 5 | 8 | 7 | 5 | 5 | 6 | 7 | 7 | 7 | 7 | 6 | 5 | 6 | 3 | 7 |
| Relief_Pitcher | 8 | 7 | 14 | 11 | 11 | 8 | 11 | 13 | 8 | 8 | 12 | 12 | 10 | 9 | 10 | 8 | 8 | 13 | 10 | 11 | 11 | 10 | 13 | 10 | 11 | 11 | 10 | 11 | 12 | 14 |
| Second_Baseman | 0 | 1 | 3 | 1 | 1 | 1 | 1 | 3 | 2 | 1 | 3 | 1 | 5 | 3 | 3 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 3 | 1 |
| Shortstop | 2 | 2 | 2 | 1 | 2 | 1 | 2 | 1 | 2 | 2 | 3 | 2 | 1 | 2 | 1 | 1 | 3 | 2 | 1 | 3 | 1 | 1 | 2 | 2 | 1 | 1 | 1 | 2 | 2 | 3 |
| Starting_Pitcher | 7 | 7 | 6 | 7 | 9 | 11 | 6 | 5 | 10 | 7 | 7 | 5 | 7 | 9 | 6 | 7 | 7 | 9 | 7 | 8 | 8 | 10 | 5 | 8 | 7 | 6 | 9 | 8 | 7 | 6 |
| Third_Baseman | 2 | 1 | 2 | 2 | 1 | 2 | 3 | 1 | 1 | 2 | 1 | 2 | 1 | 1 | 2 | 3 | 2 | 1 | 1 | 2 | 1 | 1 | 0 | 1 | 2 | 2 | 1 | 1 | 2 | 1 |

count: 0 — 14

# Baseball…
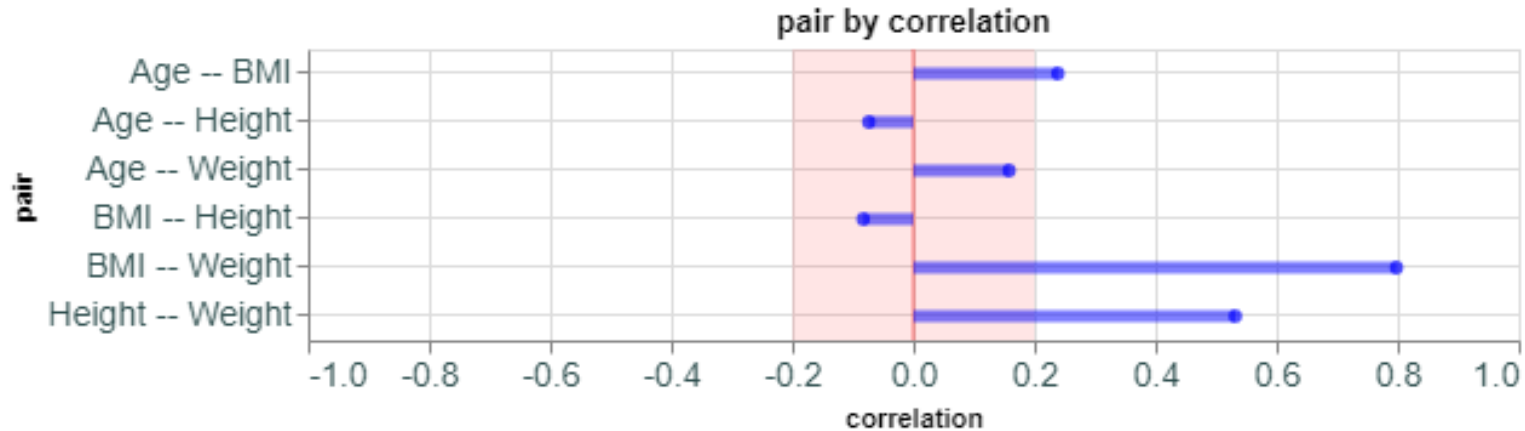
# Football …

# Swimmers vs Runners, Gagnon, et al (2018)

# Correlations ... not enough by themselves



pair by correlation

```
def calculate_BMI(h, w):
    # inches => meters, pounds => kg
    w = w / 2.205
    h = h * 0.0254
    bmi = w / (h * h)
    return bmi
```

# Charts in Altair

(DataFrame + field) + (axis + encoding + mark) = chart
- Tidy data preferred
- Composition of chart by layers
- Peered charts – cross select
- Data transformation – binning, sort, grouping
- Stack of Altair, Vega Lite, Vega. Backed by visionaries … UW IDL.

Challenges (for me)
- Composition of scenes (working outside Notebooks)
- Navigation events handled outside the chart (how to "click through") **
- Resolution of conflicting attributes in scope hierarchy – axes, domains
- Should know how Vega Lite, Vega work

# References

altair-viz.github.io

github.com/jakevdp/altair-examples

vallandingham.me/altair_intro.html

U Washington Interactive Data Lab https://idl.cs.washington.edu/
Making Data Visual  A Practical Guide to Using Visualization for Insight
By Miriah Meyer, Danyel Fisher

How Body Type May Determine Runners' and Swimmers' Destinies –
NYTimes 08-14-18

# Development Observations

- Libraries … as collective works, and as works in progress
  - Many Channels: docs, git issues (open & closed), project google group/slack, project contributor blogs, user blogs,  contributor presentations
  - Different ways to do the same thing … bushy interfaces

- Library as a lever, vs utility
  - Great to copy from Examples gallery
  - "Why did that break?" Value in knowing what it really does, and how.

- VS Code
  - Code folding, Vim plugin support, a lot of "push".

# DataViz design issues

What am I looking at?

→ clearly display clear labels for fields

Is this "a thing"? (is there an 'effect' here?)

→ visualize statistical measures too

Scale # fields

→ groups, hierarchies in schema

→ rank and trim display via measures over measures

Scale # records

→ measure, aggregate, bin *before* handing off to display substrate

# Analyses – humans in the loop

To Judge –
   to test a hypothesis,
   to ask and answer a question in some medium

Operations toward judgment:
   Describe, Compare, Abstract, Infer.

Visual "medium":
   Person, through DataViz.
   Look → Judge

Computational "medium":
   System, through Algorithms.
   Compute → Look → Judge.