

A wireframe figure of a person in a running pose, composed of a grid of points and lines, set against a dark green background with a grid pattern and glowing circular light trails.

Sports Analytics Second Interactive



Sports Analytics
Winter 2023

Topics for Today (20+20+20 minutes)

- 1. Pandas Intro**
- 2. Case Study - Data Prep for Analysis**

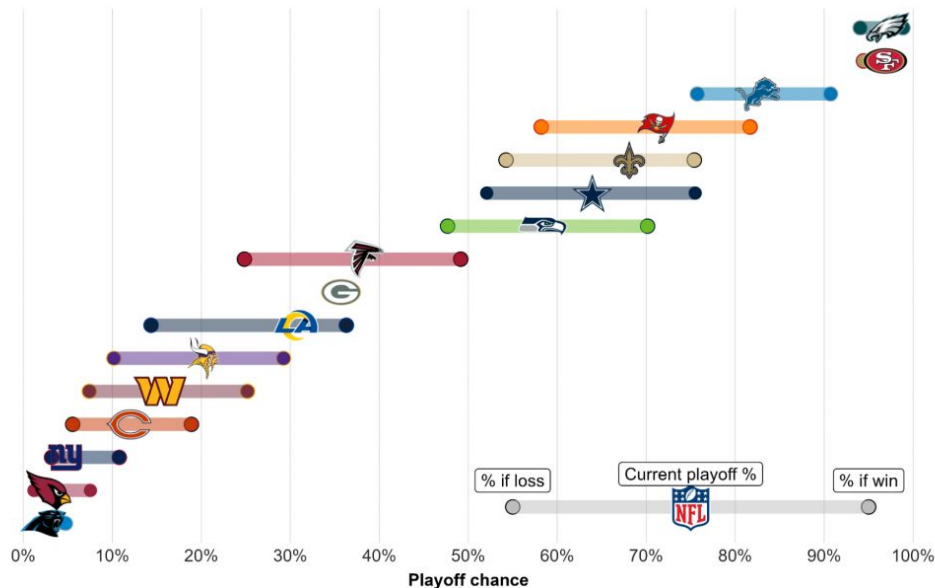
Warm up!

Just think. Don't shout out the answer

Playoff chances - back in Week 6

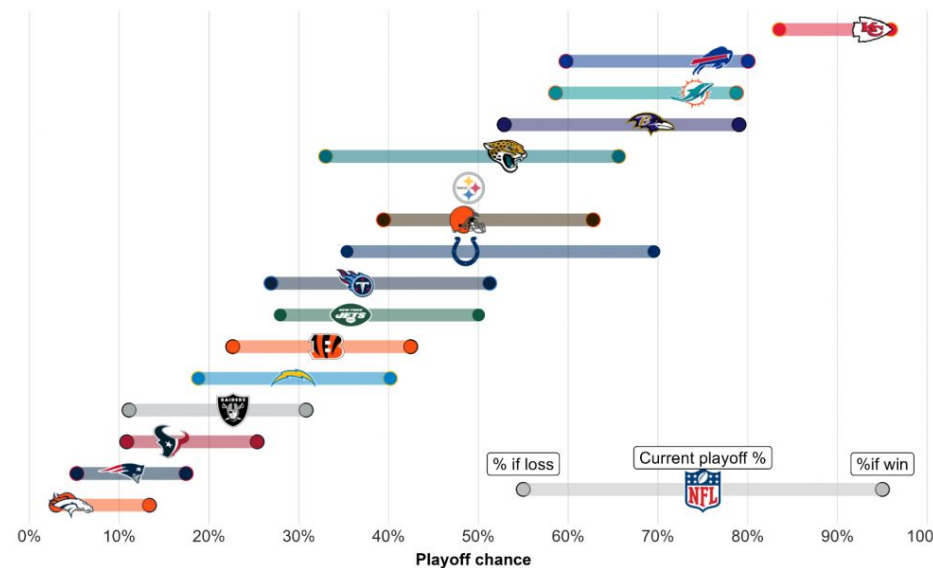
NFC playoff chances and leverage

2023 NFL season, going into week 6



AFC playoff chances and leverage

2023 NFL season, going into week 6





The Sports Analytics “Journey”

1. Ask some (what-if) questions
2. Frame as a “Story”
3. Gather Data
4. Descriptive
5. Visualization
6. Comparative Analysis
7. Predictive
8. Simulations
9. Communicate the Story /Presentation



If you do not know how to ask the right question,
you discover nothing.

(W. Edwards Deming)



Pythagorean Expectation



Sports Analytics
Winter 2023

$$\text{Win \%} = \frac{(\text{Runs Scored})^2}{(\text{Runs Scored})^2 + (\text{Runs Against})^2}$$

Pythagorean Expectation in Baseball

A really simple way to predict Win Percentages!

What's the big deal with PE - Pythagorean Expectation?

$$Win \% = \frac{Runs\ Scored^2}{Runs\ Scored^2 + Runs\ Conceded^2}$$

- Mid-season, future outlook
 - “Deserve to win” vs Actual wins
 - Reversion to the Mean - powerful mathematical concept
 - Applies to Baseball, Basketball, NFL
 - But not as such to Ice Hockey. Why would that be?!
 - Can play around with Exponents
-
- Go over to Colab
 - Idea: For any team, calculate “actual win percentage” and calculate “PE” and plot the two to see if the Expectation is being met

Pandas Basics



Someone dumps a Dataset on us.

What could we *possibly* look for?

What is the very first thing
we could look for?

Size of Dataset

We can try and
“Describe” our dataset

Basic Stats (Measures)

Sum

Minimum / Maximum

Average

Median / Mode

Duplicates

Quantiles



What if there a few columns
that are text? (Not numerical)

What can you do with text
columns?

Text columns could be ways
to group “stats” into useful
categories

| Name | Country | Operating System | Sales in Q2 |
|---------|---------|------------------|-------------|
| Andy | USA | Mac OS | 35,000 |
| Beverly | Canada | Windows | 21,000 |
| Charlie | Canada | Windows | 17,000 |
| Donna | Mexico | Windows | 28,500 |
| Eliza | USA | Mac OS | 77,400 |
| Farooq | Mexico | Windows | 9,800 |
| George | USA | Mac OS | 39,200 |
| Harry | Mexico | Mac OS | 27,200 |

What is Pandas?



pandas is a Python package providing fast and flexible data structures designed to make working with data both easy and intuitive.

It aims to be the fundamental high-level building block for doing practical, real world data analysis in Python.

Pandas has 2 fundamental Structures

Series and Data Frame

- Series = 1-dimensional
- Data Frame = 2 (or more) Dimensions

Pandas Series

- How many columns are in this Pandas Series?

| | |
|---|-----|
| 1 | 'A' |
| 2 | 'B' |
| 3 | 'C' |
| 4 | 'D' |
| 5 | 'E' |

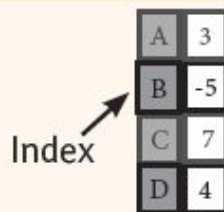
| index | | values |
|-------|---|--------|
| A | → | 5 |
| B | → | 6 |
| C | → | 12 |
| D | → | -5 |
| E | → | 6.7 |

Pandas Series is very similar to a List

Pandas Data Structures

Series

A one-dimensional labeled array capable of holding any data type



| | |
|---|----|
| A | 3 |
| B | -5 |
| C | 7 |
| D | 4 |

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

DataFrame

Columns



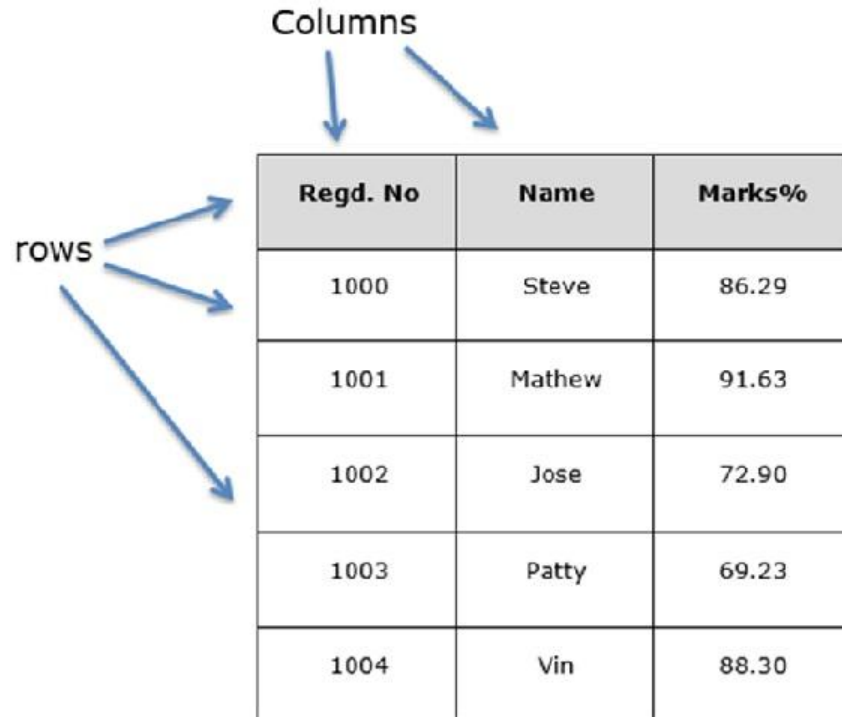
| | Country | Capital | Population |
|---|---------|-----------|------------|
| 1 | Belgium | Brussels | 11190846 |
| 2 | India | New Delhi | 1303171035 |
| 3 | Brazil | Brasília | 207847528 |

A two-dimensional labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],  
            'Capital': ['Brussels', 'New Delhi', 'Brasília'],  
            'Population': [11190846, 1303171035, 207847528]}
```

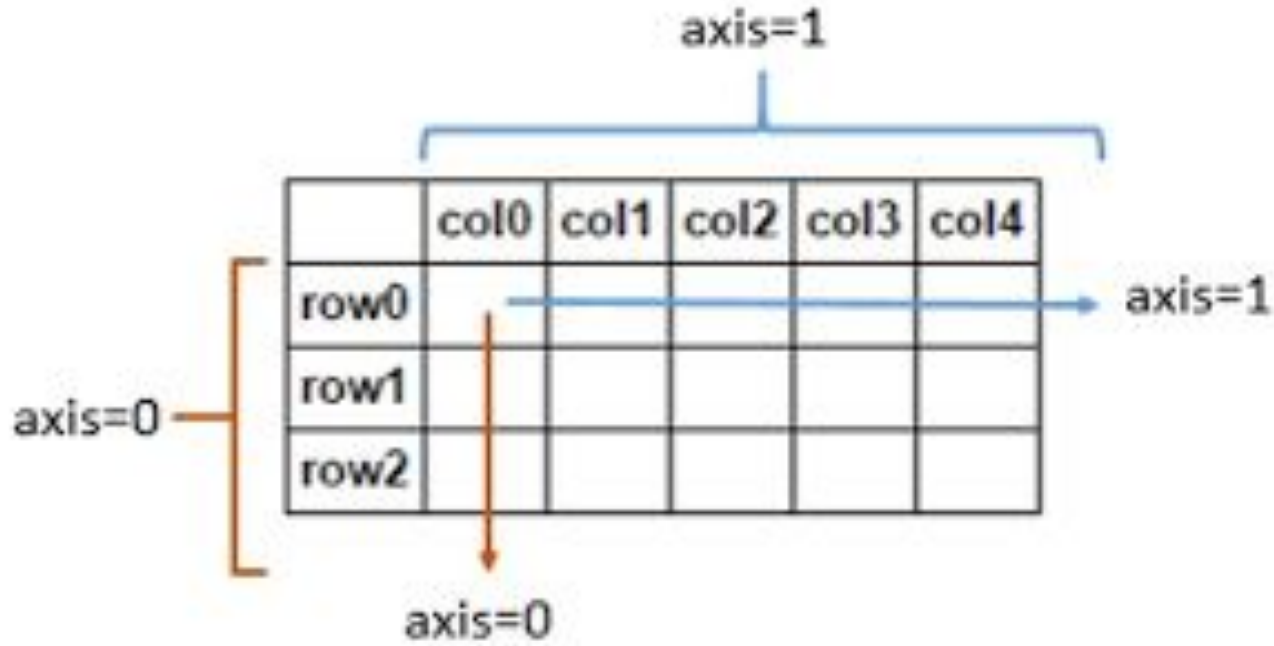
```
>>> df = pd.DataFrame(data,  
                       columns=['Country', 'Capital', 'Population'])
```

Pandas Data Frame



The diagram illustrates a Pandas Data Frame as a table. The word "Columns" is positioned above the table with two blue arrows pointing to the "Regd. No" and "Name" headers. The word "rows" is positioned to the left of the table with three blue arrows pointing to the first three rows of data.

| Regd. No | Name | Marks% |
|----------|--------|--------|
| 1000 | Steve | 86.29 |
| 1001 | Mathew | 91.63 |
| 1002 | Jose | 72.90 |
| 1003 | Patty | 69.23 |
| 1004 | Vin | 88.30 |



For simplicity, you can think of
INDEX = ROW NAMES (ROW LABELS)
COLUMNS = HEADER (COLUMN NAMES)

Pandas Dataframe

Shape = (# rows, # columns)

In [37]: df1

Out[37]:

| | GEOID | State | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 |
|---|-----------|------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 0 | 04000US01 | Alabama | 37150 | 37952 | 42212 | 44476 | 39980 | 40933 | 42590 | 43464 | 41381 |
| 1 | 04000US02 | Alaska | 55891 | 56418 | 62993 | 63989 | 61604 | 57848 | 57431 | 63648 | 61137 |
| 2 | 04000US04 | Arizona | 45245 | 46657 | 47215 | 46914 | 45739 | 46896 | 48621 | 47044 | 50602 |
| 3 | 04000US05 | Arkansas | 36658 | 37057 | 40795 | 39586 | 36538 | 38587 | 41302 | 39018 | 39919 |
| 4 | 04000US06 | California | 51755 | 55319 | 55734 | 57014 | 56134 | 54283 | 53367 | 57020 | 57528 |

How many columns are there?
How Many Rows?

In this data frame, what are we doing in each column?

| FinancialYear | 2014/2015 | 2015/2016 | 2016/2017 | 2017/2018 |
|---------------|-----------|-----------|-----------|-----------|
| Month | | | | |
| April | 8.7% | 6.3% | 6.3% | 23.6% |
| May | 7.3% | 6.8% | 6.0% | 29.1% |
| June | 4.7% | 10.4% | 6.3% | 24.0% |
| July | 5.9% | 10.2% | 5.3% | 23.3% |
| August | 9.1% | 7.9% | 9.9% | 0.0% |
| September | 8.9% | 6.1% | 12.1% | 0.0% |
| October | 9.2% | 9.9% | 7.8% | 0.0% |
| November | 9.7% | 8.3% | 8.9% | 0.0% |
| December | 9.2% | 10.9% | 6.6% | 0.0% |
| January | 8.7% | 8.0% | 10.4% | 0.0% |
| February | 9.3% | 7.0% | 10.8% | 0.0% |
| March | 9.2% | 8.1% | 9.7% | 0.0% |

What can you do with Columns?

- Number of Unique values
- Counts (Frequency of each unique value)
- Sort By Column (Ascending, Descending)
- Count Number of Missing Values
- Calculate Means for each column
- Add two columns

Some Advanced Things we can do...

- SELECT only rows that follow a particular condition
 - Very powerful!
- Remove Duplicate Values
- Apply any function to any column



**You don't have to remember
all these commands.**

**I use them a lot, but I don't
remember all the
commands!**

I just look them up!

What can you do with A DATA FRAME?

df = data frame

1. # Rows and columns
2. Examine (Head and Tail) of the df
3. Describe() each of the columns
 1. Statistics about the column
4. Statistics [mean, count, max, min, median]
5. Combine (concatenate) two data frames
6. Merge two data frames
7. Filter, Sort and GroupBy

Row index
(df.index)



Series of data



Series of data



Series of data



Series of data



Series of data



Series of data



Series of data



Column index (df.columns)

An Important Idea: Boolean Masks

- In Python & Pandas, **True** has the value of 1
- **False** has a value of 0

[True, True, False, True].sum() is 3

temps = [10, 7, 3, 9]

We can apply a “condition” to any column.

The result will be a True/False column.

temps > 7 will give [True, False, False, True]

Understanding Pandas conditional filtering

- Boolean Mask

Temp > 9

| |
|----|
| 12 |
| 8 |
| 9 |
| 7 |
| 6 |
| 10 |

| |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |

| |
|----|
| 12 |
| |
| |
| |
| |
| 10 |

| |
|----|
| 12 |
| 10 |

- True means keep the row
- False means “I don’t want the row”

```
temps = [10, 7, 3, 9]
```

```
condition = temps > 7
```

```
temps > 7 will give [True, False, False, True]
```

```
temps[condition]
```

```
temps[temps > 7]
```

True means Keep it, False means leave it out

Data Frame + Boolean Mask

Apply the mask

| | | | | |
|--|--|--|--|--|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

Create the mask

| |
|---|
| |
| T |
| T |
| F |
| T |
| F |
| T |
| T |

Get the result. Store it with
A new name

Keep only the rows
that match condition.
Mask (black out) all
other rows

End up with fewer
rows

Data Frame + Boolean Mask

df + condition = Boolean
Column

df

+

| T |
|---|
| T |
| F |
| T |
| F |
| T |
| T |

Boolean
Column

=

= Smaller_df

Grouping Operations

Sub-group Insights

Getting insights at each subgroup level is a core part of exploratory sports data analysis

To understand how various sub-categories are performing

1. Win Percentages by Region | Conference | League
2. “Production” by Player
3. High performing Venues or Stadiums for a League
4. Chronic reasons for failures. (Group by type of unforced errors)



| | species | sepal_length | sepal_width | petal_length | petal_width |
|-----|------------|--------------|-------------|--------------|-------------|
| 0 | setosa | 5.1 | 3.5 | 1.4 | 0.2 |
| 1 | setosa | 4.9 | 3.0 | 1.4 | 0.2 |
| 2 | setosa | 4.7 | 3.2 | 1.3 | 0.2 |
| 3 | setosa | 4.6 | 3.1 | 1.5 | 0.2 |
| 4 | setosa | 5.0 | 3.6 | 1.4 | 0.2 |
| 50 | versicolor | 7.0 | 3.2 | 4.7 | 1.4 |
| 51 | versicolor | 6.4 | 3.2 | 4.5 | 1.5 |
| 52 | versicolor | 6.9 | 3.1 | 4.9 | 1.5 |
| 53 | versicolor | 5.5 | 2.3 | 4.0 | 1.3 |
| 54 | versicolor | 6.5 | 2.8 | 4.6 | 1.5 |
| 100 | virginica | 6.3 | 3.3 | 6.0 | 2.5 |
| 101 | virginica | 5.8 | 2.7 | 5.1 | 1.9 |
| 102 | virginica | 7.1 | 3.0 | 5.9 | 2.1 |
| 103 | virginica | 6.3 | 2.9 | 5.6 | 1.8 |
| 104 | virginica | 6.5 | 3.0 | 5.8 | 2.2 |

| | species | sepal_length | sepal_width | petal_length | petal_width |
|--|------------|--------------|-------------|--------------|-------------|
| | setosa | 24.3 | 16.4 | 7.0 | 1.0 |
| | versicolor | 32.3 | 14.6 | 22.7 | 7.2 |
| | virginica | 32.0 | 14.9 | 28.4 | 10.5 |

Sort into groups based on “Grouping Column(s)”

For each group, take some OTHER Column (numeric)

Apply some statistic function. (SUM, Count, Average etc.)

Combine and present the results

The background features a wireframe illustration of a person in a running pose, overlaid on a grid of glowing green data points and lines. The overall color scheme is dark teal and green.

Data Prep Case Study



Sports Analytics
Winter 2023

Real Life Case Stude: Soccer Players Dataset - Physical attributes, Financial, Clubs

| | Name | Age | Nationality | Club | Loaned From | Value | Wage | Height | Weight |
|-------|-------------------|-----|-------------|---------------------|-------------|---------|-------|--------|--------|
| 0 | L. Messi | 31 | Argentina | FC Barcelona | NaN | €110.5M | €565K | 5'7 | 159lbs |
| 1 | Cristiano Ronaldo | 33 | Portugal | Juventus | NaN | €77M | €405K | 6'2 | 183lbs |
| 2 | Neymar Jr | 26 | Brazil | Paris Saint-Germain | NaN | €118.5M | €290K | 5'9 | 150lbs |
| 3 | De Gea | 27 | Spain | Manchester United | NaN | €72M | €260K | 6'4 | 168lbs |
| 4 | K. De Bruyne | 27 | Belgium | Manchester City | NaN | €102M | €355K | 5'11 | 154lbs |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 17316 | R. Akbari | 18 | Australia | Melbourne Victory | NaN | €140K | €1K | 5'11 | 150lbs |
| 17317 | A. Reyes Ambuila | 18 | Colombia | Atlético Nacional | NaN | €130K | €1K | 5'10 | 165lbs |
| 17318 | I. Sadiq | 18 | Ghana | FC Nordsjælland | NaN | €140K | €1K | 5'9 | 154lbs |
| 17319 | C. Merrie | 19 | England | Wigan Athletic | NaN | €140K | €3K | 5'11 | 150lbs |
| 17320 | L. Morrison | 19 | Scotland | Sligo Rovers | NaN | €130K | €1K | NaN | NaN |

| | Loaned From | Club |
|-------|-----------------|-----------------------------|
| 12269 | NaN | Independiente Medellín |
| 10834 | NaN | Envigado FC |
| 12811 | NaN | NaN |
| 5746 | NaN | SK Slavia Praha |
| 10514 | SV Darmstadt 98 | TSV 1860 München |
| 5747 | NaN | Gimnasia y Esgrima La Plata |
| 2345 | NaN | Stoke City |
| 11282 | NaN | Independiente Medellín |
| 6856 | NaN | Portimonense SC |
| 14757 | NaN | GIF Sundsvall |

| | Name | Value | Wage | Release Clause |
|-------|---------------|-------|------|----------------|
| 4287 | D. Andrade | €2M | €2K | 2361 |
| 12182 | Yang Xiaotian | €290K | €4K | 11472 |
| 4233 | J. Quiñones | €3.9M | €25K | 11370 |
| 5781 | Roldão Riso | €925K | €12K | 13780 |
| 16321 | B. Vera | €200K | €1K | 9962 |
| 3770 | M. Abeid | €3.2M | €10K | 6720 |
| 5912 | K. Mączyński | €700K | €5K | 957 |
| 10387 | N. Mezquida | €650K | €2K | 10490 |
| 6098 | M. Petković | €1.1M | €1K | 653 |
| 2066 | S. Sinclair | €6.5M | €48K | 4366 |
| | | | | €11./M |

| Name | Height | Weight |
|----------------|--------|--------|
| S. Deli | 6'4 | 185lbs |
| P. Taylor | 5'7 | 157lbs |
| S. Salinas | 5'10 | 150lbs |
| J. Kotzke | 6'0 | 165lbs |
| M. Bakker | 6'0 | 159lbs |
| R. Gall | 5'9 | 154lbs |
| Wilson Eduardo | 5'10 | 161lbs |
| Álvaro Traver | 6'2 | 168lbs |
| C. Wilson | 5'11 | 146lbs |
| Bruno Alves | 6'2 | 179lbs |



Data Clean up activities Needed

1. **Drop** Loaned From **Column**: Not relevant to the analysis.
2. If a player has no Club affiliation, we want to **Drop that row**
3. **Remove** (strip) the '€' symbol, and
4. There are 3 financial columns -- Wages Value and Release Clause. These are in M, or K formats –we want to **Convert them to proper integers**.
5. Convert Dollar values **strings to Integers**
6. If a Player's "value" is missing, we want to drop that player (**Drop Row**)
7. The Heights are given as a string 5' 11" to be **Converted to Integer** – 71 inches.
8. **Drop the rows** where Player Height is missing
9. Each player's weight has lbs attached to it. **Strip “lbs”, Convert into an integer**.

Let's Go to Colab
And do it step by step

Case Study: Interactive