

## Boring but important disclaimers:

- ▶ If you are not getting this from the GitHub repository or the associated Canvas page (e.g. CourseHero, Chegg etc.), you are probably getting the substandard version of these slides Don't pay money for those, because you can get the most updated version for free at

<https://github.com/julianmak/academic-notes>

The repository principally contains the compiled products rather than the source for size reasons.

- ▶ Associated Python code (as Jupyter notebooks mostly) will be held on the same repository. The source data however might be big, so I am going to be naughty and possibly just refer you to where you might get the data if that is the case (e.g. JRA-55 data). I know I should make properly reproducible binders etc., but I didn't...
- ▶ I do not claim the compiled products and/or code are completely mistake free (e.g. I know I don't write Pythonic code). Use the material however you like, but use it at your own risk.
- ▶ As said on the repository, I have tried to honestly use content that is self made, open source or explicitly open for fair use, and citations should be there. If however you are the copyright holder and you want the material taken down, please flag up the issue accordingly and I will happily try and swap out the relevant material.

# OCES 3301 : basic Data Analysis in ocean sciences

## Session 2: basic manipulations and statistics

# Outline

(Just overview here; for actual content see Jupyter notebooks)

- ▶ basic stats with basic example
- ▶ El Nino 3.4 SST data
  - demonstration of data
  - overview
  - some plotting + exercises

# Basic stats

Suppose I have some data samples as the following:

$$x_i = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

- ▶ **sample size**  $N$ , the number of samples
- ▶ **range**, largest minus smallest of sample  
→ crude measure of spread

**averages** (but actually three of these):

1. **mode**, most frequent occurrence
2. **median**, rank these, and find the middle one
3. **mean**, THE average

## Basic stats

$$x_i = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

- ▶ **lower/upper (25/75 percent) quartile**, rank data, value at which 25/75 percent of data lie below
- ▶ **inter-quartile range**, the different between upper and lower quartile  
→ measures spread

# Basic stats

Summary as a **box-and-whisker plot**

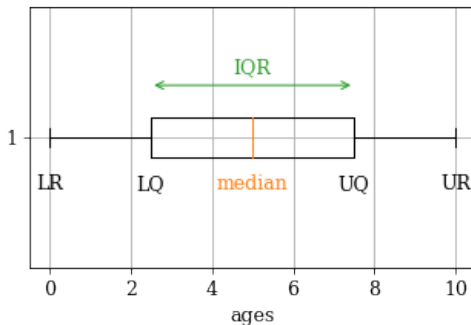


Figure: Nobel prize winning box plot.

## Basic stats

$$x_i = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

► **mean**, THE average

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_N}{N} = \frac{1}{N} \sum_{i=1}^N x_i$$

→ sum up, divide by number going into sum

## Basic stats

$$x_i = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

► (unadjusted) **variance**

$$s^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_N - \bar{x})^2}{N} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2,$$

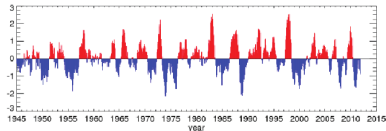
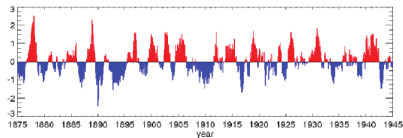
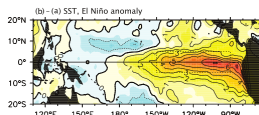
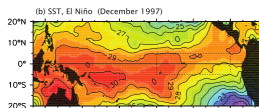
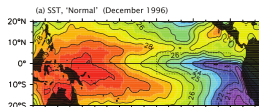
→ take mean off sample, square each result, sum, divide by number going into sum

→ square-root of variance is the **standard deviation (s.t.d.)**



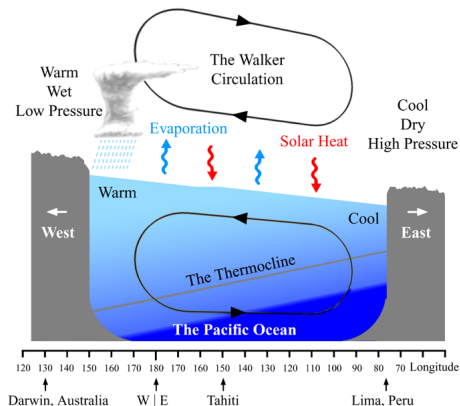
# El-Niño (see also ENVS 3004; probably see also OCES 4001)

- ▶ “the little boy”, known to fisherman in South America for a long time
- ▶ generally starts around Christmas time
- ▶ warming in Eastern equatorial Pacific ocean
- ▶ signal in SST in modern day  
→ proxy data from corals



# Southern Oscillation

- ▶ discovered by Gilbert Walker (1868–1958)
  - correlation with monsoon and thus famine and drought in India
  - Companion of the Order of the Star of India in 1911
- ▶ winds change directions periodically
  - Walker circulation changes (E-W, part of N-S Hadley circulation)



(Ocean heat source moving affects atmospheric circulation)

# “Normal” + El-Niño event

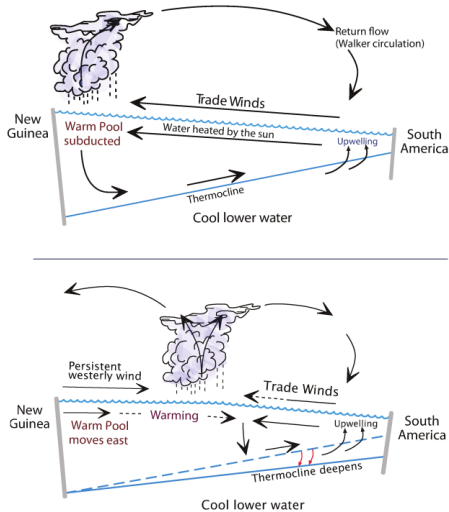


Figure: Schematic of ENSO, from Vallis (2019).

## El-Niño 3.4 region

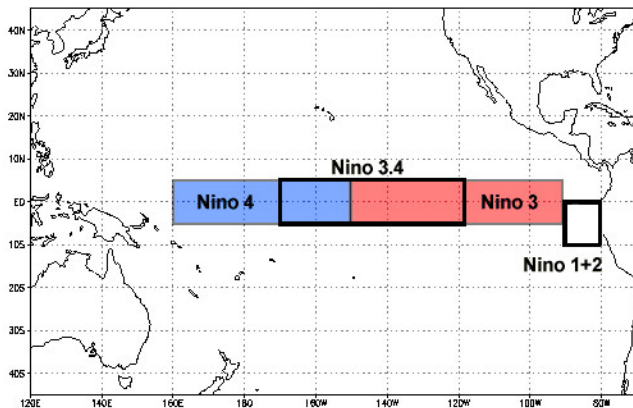


Figure: Pre-defined regions related to El-Niño indices. Picture probably (?) from NOAA.

# El-Niño 3.4 SST

Example of time-series data (see S07 and S08 also)

1	1948	1949	1950	1951	1952	1953	1954	1955	1956	1957	1958	1959	1960	1961	1962	1963	1964	1965	1966	1967	1968	1969
2	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99
3	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99	-99.99
4	24.55	25.06	25.87	26.28	26.18	26.46	26.29	25.88	25.74	25.69	25.47	25.29	25.15	25.01	24.87	24.73	24.59	24.45	24.31	24.17	24.03	23.89
5	25.24	25.71	26.90	27.58	27.92	27.73	27.60	27.02	27.23	27.20	27.25	26.91	26.77	26.63	26.49	26.35	26.21	26.07	25.93	25.79	25.65	25.51
6	26.67	26.74	27.17	27.80	27.79	27.18	26.53	26.30	26.36	26.26	25.92	26.21	26.07	25.93	25.79	25.65	25.51	25.37	25.23	25.09	24.95	24.81
7	26.74	27.00	27.57	28.04	28.28	28.12	27.43	26.94	27.01	26.87	26.88	27.00	26.86	26.72	26.58	26.44	26.30	26.16	26.02	25.88	25.74	25.60
8	26.98	27.03	26.90	26.64	27.12	26.80	26.11	25.43	25.12	25.23	25.57	25.26	25.12	24.98	24.84	24.70	24.56	24.42	24.28	24.14	24.00	23.86
9	25.61	25.81	26.22	26.60	26.66	26.55	26.15	25.51	25.28	24.41	24.25	24.57	24.43	24.29	24.15	24.01	23.87	23.73	23.59	23.45	23.31	23.17
10	25.34	25.76	26.46	26.85	27.13	26.81	26.23	25.68	25.73	25.75	25.56	25.71	25.62	25.48	25.34	25.20	25.06	24.92	24.78	24.64	24.50	24.36
11	26.04	26.54	27.46	28.23	28.55	28.36	28.17	27.69	27.44	27.42	27.62	27.90	27.76	27.62	27.48	27.34	27.20	27.06	26.92	26.78	26.64	26.50
12	28.33	28.24	28.27	28.27	28.31	27.99	27.32	26.85	26.40	26.45	26.75	26.62	26.48	26.34	26.20	26.06	25.92	25.78	25.64	25.50	25.36	25.22
13	27.07	27.18	27.47	27.88	27.70	27.37	26.44	26.09	25.92	26.24	26.04	26.18	26.04	25.90	25.76	25.62	25.48	25.34	25.20	25.06	24.92	24.78
14	26.27	26.29	26.98	27.49	27.68	27.24	26.88	26.70	26.44	26.22	26.26	26.22	26.08	25.94	25.80	25.66	25.52	25.38	25.24	25.10	24.96	24.82
15	26.23	26.56	26.94	27.36	27.75	27.67	26.89	26.19	25.78	25.71	26.07	25.97	25.83	25.69	25.55	25.41	25.27	25.13	24.99	24.85	24.71	24.57
16	25.96	26.19	26.80	27.13	27.05	27.08	26.76	26.33	25.94	25.97	25.75	25.67	25.53	25.39	25.25	25.11	24.97	24.83	24.69	24.55	24.41	24.27
17	25.77	26.22	27.18	27.78	27.63	27.62	27.78	27.48	27.40	27.36	27.47	27.62	27.48	27.34	27.20	27.06	26.92	26.78	26.64	26.50	26.36	26.22
18	27.34	27.13	27.02	26.95	26.82	26.59	26.33	25.60	25.32	25.37	25.26	25.23	25.09	24.95	24.81	24.67	24.53	24.39	24.25	24.11	23.97	23.83
19	25.66	26.19	26.94	27.38	27.99	28.09	27.90	27.97	28.01	28.17	28.12	27.96	27.82	27.68	27.54	27.40	27.26	27.12	26.98	26.84	26.70	26.56
20	27.67	27.55	28.21	28.16	27.55	27.64	27.33	26.48	26.27	26.22	26.23	26.03	25.89	25.75	25.61	25.47	25.33	25.19	25.05	24.91	24.77	24.63
21	25.88	26.11	26.50	26.74	27.35	27.47	26.97	26.44	25.86	25.97	26.08	25.95	25.81	25.67	25.53	25.39	25.25	25.11	24.97	24.83	24.69	24.55
22	25.69	25.68	26.33	27.10	27.19	27.88	27.58	27.01	26.72	26.75	27.20	27.27	27.13	26.99	26.85	26.71	26.57	26.43	26.29	26.15	26.01	25.87
23	27.50	27.86	27.82	28.13	28.29	27.69	27.08	27.02	27.15	27.34	27.10	26.98	26.84	26.70	26.56	26.42	26.28	26.14	26.00	25.86	25.72	25.58

Figure: Sample content of elnino34.sst.data.

# Iris data

Example of categorical + numerical data (multivariate), in `iris.csv`

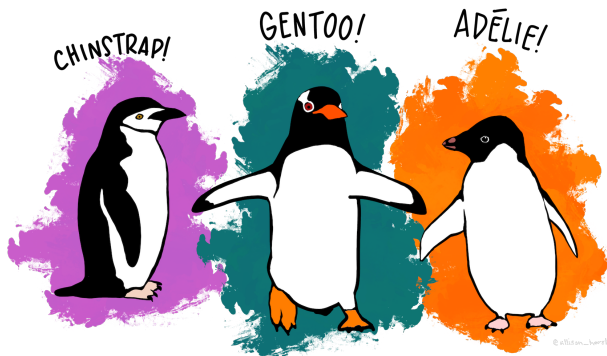
- ▶ from Ronald Fisher's dataset (will see him time and again in this course)
- ▶ the “hello world” of statistics and machine learning



**Figure:** *Iris setosa*, *versicolor*, and *virginica*, in the iris dataset. Often used in machine learning (e.g. clustering analysis) and useful for demonstrating statistical concepts. Pictures from Wikipedia.

# Palmer Penguin data

Fisher was known to be a proponent of eugenics (whether he was a racist as such is ongoing debate), so if you want an alternative to iris data, there is the penguin data in `penguins.csv`



**Figure:** Artwork from Palmer Penguins dataset, by Allison Horst. See <https://allisonhorst.github.io/palmerpenguins/articles/intro.html>.

# Jupyter notebook

Go to 02 Jupyter notebook to play around with these datasets in Python