

# Programming Principles in Python (CSCI 503/490)

---

Data

Dr. David Koop

# Quiz

# Quiz

---

1. What is a difference between a pandas Series and a polars Series?
  - (a) A pandas Series has an index; a polars Series does not
  - (b) A pandas Series has a name; a polars Series does not
  - (c) A pandas Series is typed; a polars Series is untyped
  - (d) A pandas Series is mutable; a polars Series is immutable

# Quiz

---

2. Which of the following is invalid code?

- (a) `numpy.array([[1.5, 2, 3], [4, 5, 6]], dtype='float')`
- (b) `numpy.array([[1.5, 2, 3], [4, 5, 6]])`
- (c) `numpy.array([[1, 2, 3], [4, 5]])`
- (d) `numpy.array([[1, 2, 3], [4, 5, 6]], dtype="float")`

# Quiz

---

3. Which of the following is **true** about numpy arrays and python lists?
- (a) Arrays are mutable; lists are not
  - (b) Lists are faster to access than arrays
  - (c) Arrays and lists both can have elements of different types
  - (d) Assigning to a slice (`[a[1:3] = [10,20]`) is a valid operation for both arrays and lists

# Quiz

---

4. Evaluate `pd.Series([1, 2, 3]) + pd.Series([3, 2, 1], [1, 2, 0])`.

(a) `pd.Series([4, 4, 4], [0, 1, 2])`

(b) `pd.Series([2, 4, 6], [0, 1, 2])`

(c) `pd.Series([2, 5, 5], [0, 1, 2])`

(d) There is an error

# Quiz

---

5. Given the array `arr = np.array([[1, 2], [3, 4], [5, 6]])`, what is

`arr[:, 1].shape`?

- (a) `(3, )`
- (b) `(2, 1)`
- (c) `(3, 1)`
- (d) `(1, 3)`

# pandas

---

- Contains high-level data structures and manipulation tools designed to make data analysis fast and easy in Python
- Originally built on top of NumPy
- Built with the following requirements:
  - Data structures with labeled axes (aligning data)
  - Support time series data
  - Do arithmetic operations that include metadata (labels)
  - Handle missing data
  - Add merge and relational operations

# polars

---

- Contains high-level data structures and manipulation tools designed to make data analysis "**lightning**" fast and easy in Python
  - Built using Apache Arrow
  - Written from scratch using Rust but with a Python API
  - Parallelized (uses multiple cores)
  - Intuitive API

# Series

---

- A one-dimensional data structure (with a type)
  - `s = pl.Series([1, 2, 3])`
  - `t = pd.Series([1, 2, 3])`
- May also have a name and dtype
  - `s = pl.Series('name', ['a', 'b', 'c'], dtype=pl.Float)`
  - `t = pd.Series([1, 2, 3], name='num', dtype='float')`
- In pandas, a series has an index
  - `ti = pd.Series([1, 2, 3], ['a', 'b', 'c'])` # index ['a', 'b', 'c']
  - `ti = pd.Series({'a': 1, 'b': 2, 'c': 3})` # same index
- Indexing: `s[0]`, `t[0]`, `ti['a']`, `ti.iloc[0]`, `ti.loc['a']`

# Series Operations

---

- Like numpy: elementwise / broadcasting
  - `Series([1,2,3]) + Series([1,2,3]) # Series([2,4,6])`
  - `Series([1,2,3]) + 4 # Series([5,6,7])`
- ...but for pandas, with custom indexes, the operations **align** on the index:
  - `pd.Series([1,2,3], index=list('abc')) +  
pd.Series([1,2,3], index=list('cba'))  
# pd.Series([4,4,4], index=['a','b','c'])`
  - also have `.add`, `.subtract`, ... with `fill_value` argument

# DataFrame

---

- A collection of Series (uniquely named)
  - Similar to a table in a database
  - Similar to a sheet in a spreadsheet
- ```
df = DataFrame({'state': ['Ohio', 'Ohio', 'Ohio', 'Nevada'],  
               'year': [2000, 2001, 2002, 2001],  
               'pop': [1.5, 1.7, 3.6, 2.4]})
```
- In pandas:
  - Has an index shared with each series
  - Index is automatically assigned just as with a series but can be passed in as well via `index` kwarg

# Assignment 7

---

- Concurrency, System Integration, and Structural Pattern Matching
- Download System Logs
- Locate Logs of Interest
- Read JSON & Binary Data
- Filter suspicious events using structural pattern matching
- Process all files using threading

# Schedule

---

- Next week (April 20 and 22):
  - I am at a conference
  - Recorded lectures for **all** sections, no in-person lectures
  - No In-person Office hours: Email questions, virtual meetings by request
  - Assignment 8 released, due May 1
- April 27 and April 29: Normal lectures (in-person and online)
- May 6: Final Exam (starts at **8:00am**)

# DataFrame Indexing and Slicing

---

- polars:
  - `df[0]`, `df[0:1]` # equivalent, data frame with single row
- pandas:
  - `dfa[0]` # error
  - `dfa.loc[0]` # a Series!
  - `dfa[0:2]` # a data frame with two rows
- pandas with an index (`dfi = dfa.set_index('state')`)
  - `dfi['Texas']`, `dfi['Ohio']` # a Series, a DataFrame!
  - `dfi.loc['Ohio':'Texas']` # inclusive slice!
  - `dfi.iloc[0:2]` # not inclusive!

# pandas DataFrame Indexing and Slicing

---

- Same as with NumPy arrays but can use index labels
- Slicing with labels: NumPy is **exclusive**, Pandas is **inclusive!**
  - `s = Series(np.arange(4))`  
`s[0:2]` # gives two values like numpy
  - `s = Series(np.arange(4), index=['a', 'b', 'c', 'd'])`  
`s['a':'c']` # gives three values, not two!
- Obtaining data subsets
  - `loc`: get rows/cols by label
  - `iloc`: get rows/cols by position (integer index)

# DataFrame Filtering

---

- polars:

- `df['pop'] > 2` # boolean Series
- `df.filter(pl.col('pop') > 2)` # subset of dataframe

- pandas:

- `dfa['pop'] > 2` # boolean Series
- `dfa[dfa['pop'] > 2]` # subset of dataframe
- `dfa.query('pop > 2')` # subset of dataframe

- Multiple criteria, use `&`, `|`, and `~`; remember parentheses!

- `df.filter((pl.col('year') < 2002) & (pl.col('pop') > 2))`
- `dfa[(dfa['year'] < 2002) & (dfa['pop'] > 2)]`

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

344 rows x 17 columns

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

Column Names

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

344 rows x 17 columns

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

Column Names

| studyName | Sample Number | Species | Region | Island | Stage | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----------|---------------|---------|--------|--------|-------|---------------|-------------------|----------|--------------------|
|-----------|---------------|---------|--------|--------|-------|---------------|-------------------|----------|--------------------|

Index

|     |         |     |                                     |        |           |                    |       |     |          |      |
|-----|---------|-----|-------------------------------------|--------|-----------|--------------------|-------|-----|----------|------|
| 0   | PAL0708 | 1   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1  | Yes | 11/11/07 | 39.1 |
| 1   | PAL0708 | 2   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2  | Yes | 11/11/07 | 39.5 |
| 2   | PAL0708 | 3   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1  | Yes | 11/16/07 | 40.3 |
| 3   | PAL0708 | 4   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2  | Yes | 11/16/07 | NaN  |
| 4   | PAL0708 | 5   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1  | Yes | 11/16/07 | 36.7 |
| ... | ...     | ... | ...                                 | ...    | ...       | ...                | ...   | ... | ...      | ...  |
| 339 | PAL0910 | 120 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2 | No  | 12/1/09  | NaN  |
| 340 | PAL0910 | 121 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1 | Yes | 11/22/09 | 46.8 |
| 341 | PAL0910 | 122 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2 | Yes | 11/22/09 | 50.4 |
| 342 | PAL0910 | 123 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1 | Yes | 11/22/09 | 45.2 |
| 343 | PAL0910 | 124 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2 | Yes | 11/22/09 | 49.9 |

344 rows x 17 columns

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

Column Names

| studyName | Sample Number | Species | Region | Island | Stage | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----------|---------------|---------|--------|--------|-------|---------------|-------------------|----------|--------------------|
|-----------|---------------|---------|--------|--------|-------|---------------|-------------------|----------|--------------------|

Index

|     |         |     |                                     |        |           |                    |       |     |          |      |
|-----|---------|-----|-------------------------------------|--------|-----------|--------------------|-------|-----|----------|------|
| 0   | PAL0708 | 1   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1  | Yes | 11/11/07 | 39.1 |
| 1   | PAL0708 | 2   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2  | Yes | 11/11/07 | 39.5 |
| 2   | PAL0708 | 3   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1  | Yes | 11/16/07 | 40.3 |
| 3   | PAL0708 | 4   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2  | Yes | 11/16/07 | NaN  |
| 4   | PAL0708 | 5   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1  | Yes | 11/16/07 | 36.7 |
| ... | ...     | ... | ...                                 | ...    | ...       | ...                | ...   | ... | ...      | ...  |
| 339 | PAL0910 | 120 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2 | No  | 12/1/09  | NaN  |
| 340 | PAL0910 | 121 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1 | Yes | 11/22/09 | 46.8 |
| 341 | PAL0910 | 122 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2 | Yes | 11/22/09 | 50.4 |
| 342 | PAL0910 | 123 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1 | Yes | 11/22/09 | 45.2 |
| 343 | PAL0910 | 124 | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2 | Yes | 11/22/09 | 49.9 |

344 rows x 17 columns

Column: df['Island']

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

Column Names

|            | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|------------|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| <b>0</b>   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| <b>1</b>   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| <b>2</b>   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| <b>3</b>   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| <b>4</b>   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ...        | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| <b>339</b> | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| <b>340</b> | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| <b>341</b> | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| <b>342</b> | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| <b>343</b> | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

Row: `df.loc[2]`

Index

344 rows x 17 columns

Column: `df['Island']`

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

Column Names

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

Row: `df.loc[2]`

Index

Cell: `df.loc[341, 'Species']`

Column: `df['Island']`

344 rows x 17 columns

# pandas DataFrame

```
df = pd.read_csv('penguins_lter.csv')
```

Column Names

| studyName | Sample Number | Species | Region | Island | Stage | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----------|---------------|---------|--------|--------|-------|---------------|-------------------|----------|--------------------|
|-----------|---------------|---------|--------|--------|-------|---------------|-------------------|----------|--------------------|

Row: df.loc[2]

|   |         |   |                                     |        |           |                    |      |     |          |      |
|---|---------|---|-------------------------------------|--------|-----------|--------------------|------|-----|----------|------|
| 0 | PAL0708 | 1 | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1 | Yes | 11/11/07 | 39.1 |
| 1 | PAL0708 | 2 | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2 | Yes | 11/11/07 | 39.5 |
| 2 | PAL0708 | 3 | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1 | Yes | 11/16/07 | 40.3 |

Index

|     |         |     |                                     |        |           |                    |      |     |          |     |
|-----|---------|-----|-------------------------------------|--------|-----------|--------------------|------|-----|----------|-----|
| 3   | PAL0708 | 4   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2 | Yes | 11/16/07 | NaN |
| 4   | PAL0708 | 5   | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1 | Yes | 11/16/07 |     |
| ... | ...     | ... | ...                                 | ...    | ...       | ...                | ...  | ... | ...      | ... |

Missing Data

Cell: df.loc[341, 'Species']

|     |         |     |                                   |        |        |                    |       |     |          |      |
|-----|---------|-----|-----------------------------------|--------|--------|--------------------|-------|-----|----------|------|
| 339 | PAL0910 | 120 | Gentoo penguin (Pygoscelis papua) | Anvers | Biscoe | Adult, 1 Egg Stage | N38A2 | No  | 12/1/09  | NaN  |
| 340 | PAL0910 | 121 | Gentoo penguin (Pygoscelis papua) | Anvers | Biscoe | Adult, 1 Egg Stage | N39A1 | Yes | 11/22/09 | 46.8 |
| 341 | PAL0910 | 122 | Gentoo penguin (Pygoscelis papua) | Anvers | Biscoe | Adult, 1 Egg Stage | N39A2 | Yes | 11/22/09 | 50.4 |
| 342 | PAL0910 | 123 | Gentoo penguin (Pygoscelis papua) | Anvers | Biscoe | Adult, 1 Egg Stage | N43A1 | Yes | 11/22/09 | 45.2 |
| 343 | PAL0910 | 124 | Gentoo penguin (Pygoscelis papua) | Anvers | Biscoe | Adult, 1 Egg Stage | N43A2 | Yes | 11/22/09 | 49.9 |

344 rows x 17 columns

Column: df['Island']



# polars DataFrame

shape: (344, 10)

| studyName | Sample Number | Species                             | Region   | Island      | Stage                | Individual ID | Clutch Completion | Date Egg   | Culmen Length (mm) |
|-----------|---------------|-------------------------------------|----------|-------------|----------------------|---------------|-------------------|------------|--------------------|
| str       | i64           | str                                 | str      | str         | str                  | str           | str               | str        | f64                |
| "PAL0708" | 1             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A1"        | "Yes"             | "11/11/07" | 39.1               |
| "PAL0708" | 2             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A2"        | "Yes"             | "11/11/07" | 39.5               |
| "PAL0708" | 3             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A1"        | "Yes"             | "11/16/07" | 40.3               |
| "PAL0708" | 4             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A2"        | "Yes"             | "11/16/07" | null               |
| "PAL0708" | 5             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N3A1"        | "Yes"             | "11/16/07" | 36.7               |
| ...       | ...           | ...                                 | ...      | ...         | ...                  | ...           | ...               | ...        | ...                |
| "PAL0910" | 120           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N38A2"       | "No"              | "12/1/09"  | null               |
| "PAL0910" | 121           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A1"       | "Yes"             | "11/22/09" | 46.8               |
| "PAL0910" | 122           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A2"       | "Yes"             | "11/22/09" | 50.4               |
| "PAL0910" | 123           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A1"       | "Yes"             | "11/22/09" | 45.2               |
| "PAL0910" | 124           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A2"       | "Yes"             | "11/22/09" | 49.9               |

# polars DataFrame

## Column Names & Types

shape: (344, 10)

| studyName | Sample Number | Species                             | Region   | Island      | Stage                | Individual ID | Clutch Completion | Date Egg   | Culmen Length (mm) |
|-----------|---------------|-------------------------------------|----------|-------------|----------------------|---------------|-------------------|------------|--------------------|
| str       | i64           | str                                 | str      | str         | str                  | str           | str               | str        | f64                |
| "PAL0708" | 1             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A1"        | "Yes"             | "11/11/07" | 39.1               |
| "PAL0708" | 2             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A2"        | "Yes"             | "11/11/07" | 39.5               |
| "PAL0708" | 3             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A1"        | "Yes"             | "11/16/07" | 40.3               |
| "PAL0708" | 4             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A2"        | "Yes"             | "11/16/07" | null               |
| "PAL0708" | 5             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N3A1"        | "Yes"             | "11/16/07" | 36.7               |
| ...       | ...           | ...                                 | ...      | ...         | ...                  | ...           | ...               | ...        | ...                |
| "PAL0910" | 120           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N38A2"       | "No"              | "12/1/09"  | null               |
| "PAL0910" | 121           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A1"       | "Yes"             | "11/22/09" | 46.8               |
| "PAL0910" | 122           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A2"       | "Yes"             | "11/22/09" | 50.4               |
| "PAL0910" | 123           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A1"       | "Yes"             | "11/22/09" | 45.2               |
| "PAL0910" | 124           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A2"       | "Yes"             | "11/22/09" | 49.9               |

# polars DataFrame

## Column Names & Types

shape: (344, 10)

| studyName | Sample Number | Species                             | Region   | Island      | Stage                | Individual ID | Clutch Completion | Date Egg   | Culmen Length (mm) |
|-----------|---------------|-------------------------------------|----------|-------------|----------------------|---------------|-------------------|------------|--------------------|
| str       | i64           | str                                 | str      | str         | str                  | str           | str               | str        | f64                |
| "PAL0708" | 1             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A1"        | "Yes"             | "11/11/07" | 39.1               |
| "PAL0708" | 2             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A2"        | "Yes"             | "11/11/07" | 39.5               |
| "PAL0708" | 3             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A1"        | "Yes"             | "11/16/07" | 40.3               |
| "PAL0708" | 4             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A2"        | "Yes"             | "11/16/07" | null               |
| "PAL0708" | 5             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N3A1"        | "Yes"             | "11/16/07" | 36.7               |
| ...       | ...           | ...                                 | ...      | ...         | ...                  | ...           | ...               | ...        | ...                |
| "PAL0910" | 120           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N38A2"       | "No"              | "12/1/09"  | null               |
| "PAL0910" | 121           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A1"       | "Yes"             | "11/22/09" | 46.8               |
| "PAL0910" | 122           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A2"       | "Yes"             | "11/22/09" | 50.4               |
| "PAL0910" | 123           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A1"       | "Yes"             | "11/22/09" | 45.2               |
| "PAL0910" | 124           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    |                      |               |                   |            | 49.9               |

Column: df['Island']

# polars DataFrame

Column Names  
& Types

shape: (344, 10)

| studyName | Sample Number | Species                             | Region   | Island      | Stage                | Individual ID | Clutch Completion | Date Egg   | Culmen Length (mm) |
|-----------|---------------|-------------------------------------|----------|-------------|----------------------|---------------|-------------------|------------|--------------------|
| str       | i64           | str                                 | str      | str         | str                  | str           | str               | str        | f64                |
| "PAL0708" | 1             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A1"        | "Yes"             | "11/11/07" | 39.1               |
| "PAL0708" | 2             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A2"        | "Yes"             | "11/11/07" | 39.5               |
| "PAL0708" | 3             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A1"        | "Yes"             | "11/16/07" | 40.3               |
| "PAL0708" | 4             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A2"        | "Yes"             | "11/16/07" | null               |
| "PAL0708" | 5             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N3A1"        | "Yes"             | "11/16/07" | 36.7               |
| ...       | ...           | ...                                 | ...      | ...         | ...                  | ...           | ...               | ...        | ...                |
| "PAL0910" | 120           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N38A2"       | "No"              | "12/1/09"  | null               |
| "PAL0910" | 121           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A1"       | "Yes"             | "11/22/09" | 46.8               |
| "PAL0910" | 122           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A2"       | "Yes"             | "11/22/09" | 50.4               |
| "PAL0910" | 123           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A1"       | "Yes"             | "11/22/09" | 45.2               |
| "PAL0910" | 124           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    |                      |               |                   |            | 49.9               |

Row: df [ 2 ]

Column: df [ 'Island' ]

# polars DataFrame

Column Names  
& Types

shape: (344, 10)

| studyName | Sample Number | Species                             | Region   | Island      | Stage                | Individual ID | Clutch Completion | Date Egg   | Culmen Length (mm) |
|-----------|---------------|-------------------------------------|----------|-------------|----------------------|---------------|-------------------|------------|--------------------|
| str       | i64           | str                                 | str      | str         | str                  | str           | str               | str        | f64                |
| "PAL0708" | 1             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A1"        | "Yes"             | "11/11/07" | 39.1               |
| "PAL0708" | 2             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A2"        | "Yes"             | "11/11/07" | 39.5               |
| "PAL0708" | 3             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A1"        | "Yes"             | "11/16/07" | 40.3               |
| "PAL0708" | 4             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A2"        | "Yes"             | "11/16/07" | null               |
| "PAL0708" | 5             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N3A1"        | "Yes"             | "11/16/07" | 36.7               |
| ...       | ...           | ...                                 | ...      | ...         | ...                  | ...           | ...               | ...        | ...                |
| "PAL0910" | 120           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N38A2"       | "No"              | "12/1/09"  | null               |
| "PAL0910" | 121           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A1"       | "Yes"             | "11/22/09" | 46.8               |
| "PAL0910" | 122           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A2"       | "Yes"             | "11/22/09" | 50.4               |
| "PAL0910" | 123           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A1"       | "Yes"             | "11/22/09" | 45.2               |
| "PAL0910" | 124           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A2"       | "Yes"             | "11/22/09" | 49.9               |

Row: df[2]

Cell: df['Species'][341]

Column: df['Island']

# polars DataFrame

Column Names  
& Types

shape: (344, 10)

| studyName | Sample Number | Species                             | Region   | Island      | Stage                | Individual ID | Clutch Completion | Date Egg   | Culmen Length (mm) |
|-----------|---------------|-------------------------------------|----------|-------------|----------------------|---------------|-------------------|------------|--------------------|
| str       | i64           | str                                 | str      | str         | str                  | str           | str               | str        | f64                |
| "PAL0708" | 1             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A1"        | "Yes"             | "11/11/07" | 39.1               |
| "PAL0708" | 2             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N1A2"        | "Yes"             | "11/11/07" | 39.5               |
| "PAL0708" | 3             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A1"        | "Yes"             | "11/16/07" | 40.3               |
| "PAL0708" | 4             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N2A2"        | "Yes"             | "11/16/07" | null               |
| "PAL0708" | 5             | "Adelie Penguin (Pygoscelis ade..." | "Anvers" | "Torgersen" | "Adult, 1 Egg Stage" | "N3A1"        | "Yes"             | "11/16/07" | ...                |
| ...       | ...           | ...                                 | ...      | ...         | ...                  | ...           | ...               | ...        | ...                |
| "PAL0910" | 120           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N38A2"       | "No"              | "12/1/09"  | null               |
| "PAL0910" | 121           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A1"       | "Yes"             | "11/22/09" | 46.8               |
| "PAL0910" | 122           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N39A2"       | "Yes"             | "11/22/09" | 50.4               |
| "PAL0910" | 123           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A1"       | "Yes"             | "11/22/09" | 45.2               |
| "PAL0910" | 124           | "Gentoo penguin (Pygoscelis pap..." | "Anvers" | "Biscoe"    | "Adult, 1 Egg Stage" | "N43A2"       | "Yes"             | "11/22/09" | 49.9               |

Row: df[2]

Missing Data

Cell: df['Species'][341]

Column: df['Island']

# pandas Filtering

```
df[df['Culmen Length (mm)'] > 40]
```

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

344 rows x 17 columns

# pandas Filtering

```
df[df['Culmen Length (mm)'] > 40]
```

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

344 rows x 17 columns

# polars Filtering

```
df.filter(pl.col('Culmen Length (mm)') > 40)
```

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

344 rows x 17 columns

# polars Filtering

```
df.filter(pl.col('Culmen Length (mm)') > 40)
```

|     | studyName | Sample Number | Species                             | Region | Island    | Stage              | Individual ID | Clutch Completion | Date Egg | Culmen Length (mm) |
|-----|-----------|---------------|-------------------------------------|--------|-----------|--------------------|---------------|-------------------|----------|--------------------|
| 0   | PAL0708   | 1             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A1          | Yes               | 11/11/07 | 39.1               |
| 1   | PAL0708   | 2             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N1A2          | Yes               | 11/11/07 | 39.5               |
| 2   | PAL0708   | 3             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A1          | Yes               | 11/16/07 | 40.3               |
| 3   | PAL0708   | 4             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N2A2          | Yes               | 11/16/07 | NaN                |
| 4   | PAL0708   | 5             | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | N3A1          | Yes               | 11/16/07 | 36.7               |
| ... | ...       | ...           | ...                                 | ...    | ...       | ...                | ...           | ...               | ...      | ...                |
| 339 | PAL0910   | 120           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N38A2         | No                | 12/1/09  | NaN                |
| 340 | PAL0910   | 121           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A1         | Yes               | 11/22/09 | 46.8               |
| 341 | PAL0910   | 122           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N39A2         | Yes               | 11/22/09 | 50.4               |
| 342 | PAL0910   | 123           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A1         | Yes               | 11/22/09 | 45.2               |
| 343 | PAL0910   | 124           | Gentoo penguin (Pygoscelis papua)   | Anvers | Biscoe    | Adult, 1 Egg Stage | N43A2         | Yes               | 11/22/09 | 49.9               |

344 rows x 17 columns

# Sorting

---

- polars: `df.sort('pop')`
- pandas: `dfa.sort_values('pop')`
- Can sort by multiple columns, too
- pandas also has a `sort_index` method to sort by the index
  - `dfa.sort_index()`

# Statistics

---

- Many common statistical methods can be used (min, max, median, etc.)
- `describe`: shortcut for easy stats!

```
In [204]: df.describe()
```

```
Out[204]:
```

|       | one      | two       |
|-------|----------|-----------|
| count | 3.000000 | 2.000000  |
| mean  | 3.083333 | -2.900000 |
| std   | 3.493685 | 2.262742  |
| min   | 0.750000 | -4.500000 |
| 25%   | 1.075000 | -3.700000 |
| 50%   | 1.400000 | -2.900000 |
| 75%   | 4.250000 | -2.100000 |
| max   | 7.100000 | -1.300000 |

```
In [205]: obj = Series(['a', 'a', 'b', 'c'] * 4)
```

```
In [206]: obj.describe()
```

```
Out[206]:
```

|        |        |
|--------|--------|
| count  | 16     |
| unique | 3      |
| top    | a      |
| freq   | 8      |
| dtype: | object |

# Unique Values and Value Counts

---

- polars: `unique()` returns a Series/DataFrame with duplicates dropped
- pandas is more complicated
  - Series `unique()` returns an array with only the unique values (no index)
    - `s = Series(['c', 'a', 'd', 'a', 'a', 'b', 'b', 'c', 'c'])`  
`s.unique()` # `array(['c', 'a', 'd', 'b'])`
  - Data Frame `drop_duplicates` returns a DataFrame with duplicates dropped
- Also `nunique()` / `n_unique()` to count number of unique entries
- `value_counts` returns a Series/DataFrame with index frequencies:
  - `s.value_counts()` # `Series({'c': 3, 'a': 3, 'b': 2, 'd': 1})`