# CSCI 285
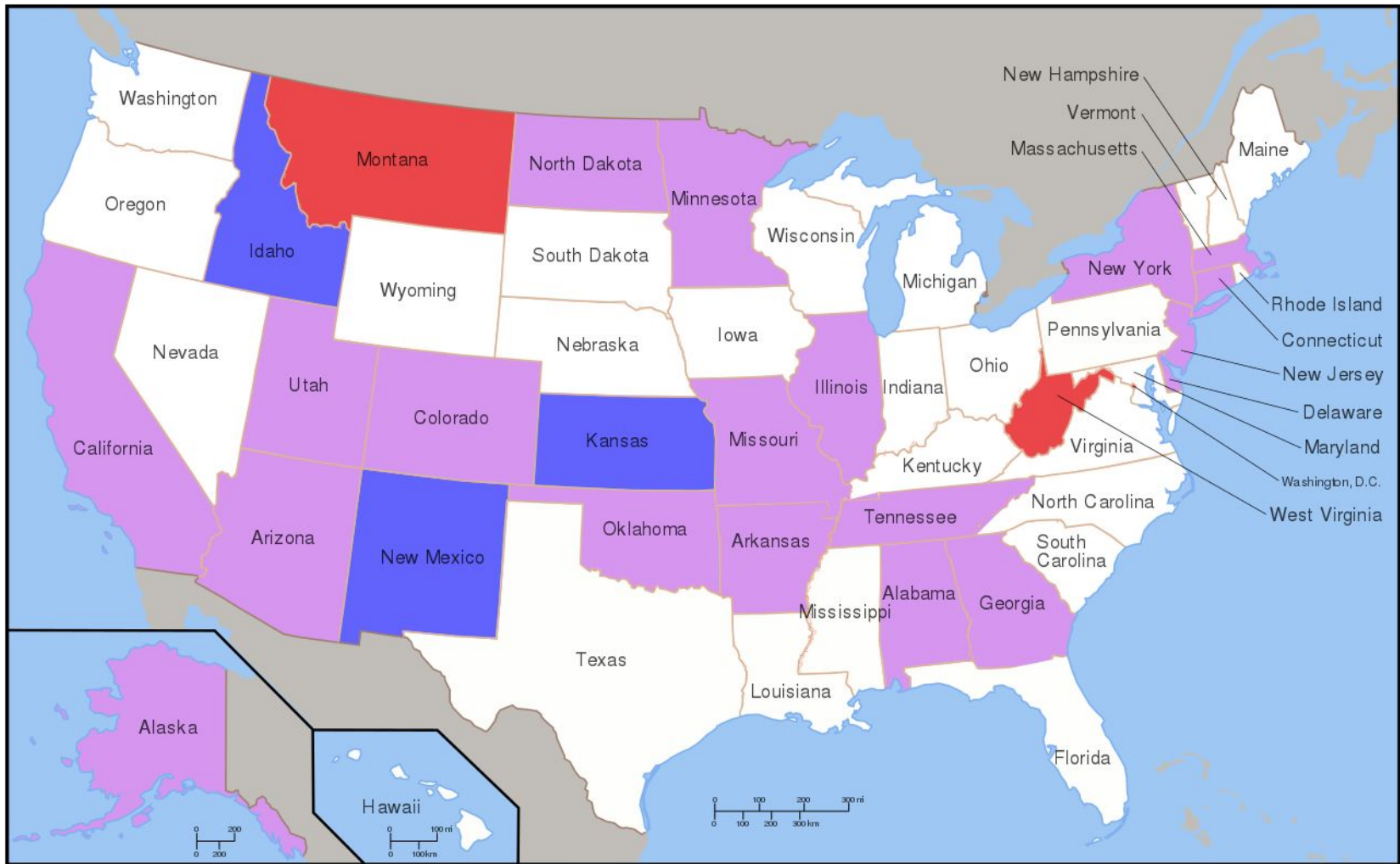# Scientific Computing
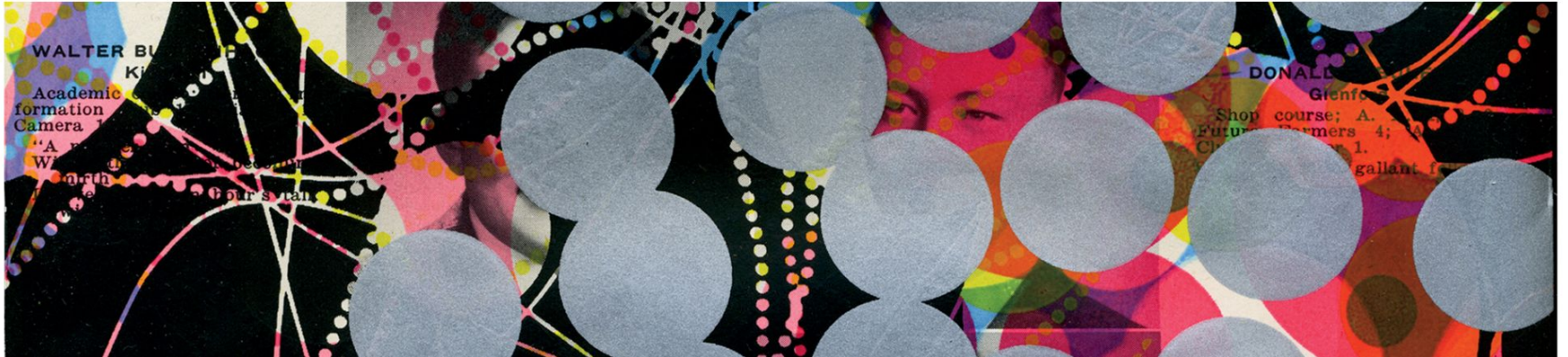
# Data Scientist: The Sexiest Job of the 21st Century

Meet the people who can coax treasure out of messy, unstructured data. by Thomas H. Davenport and DJ Patil
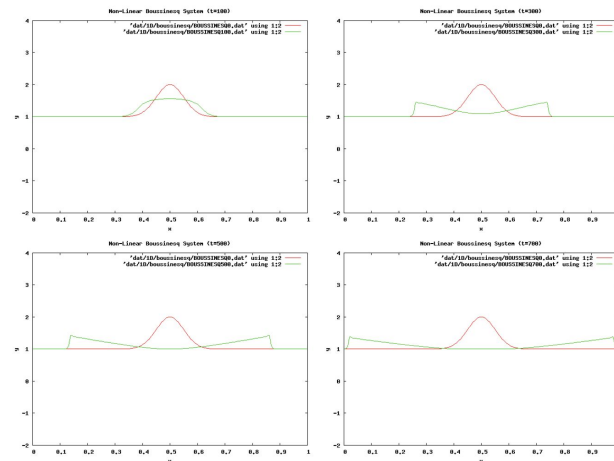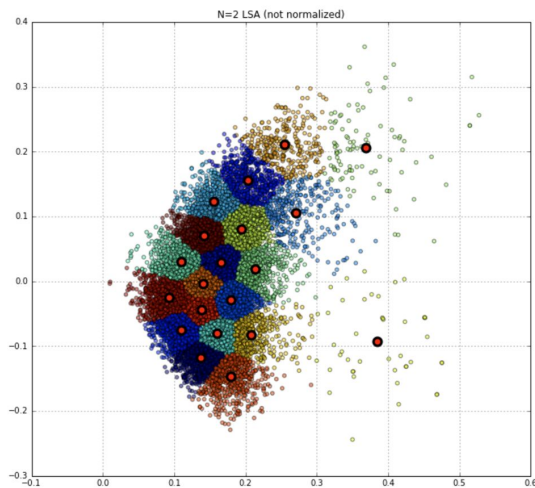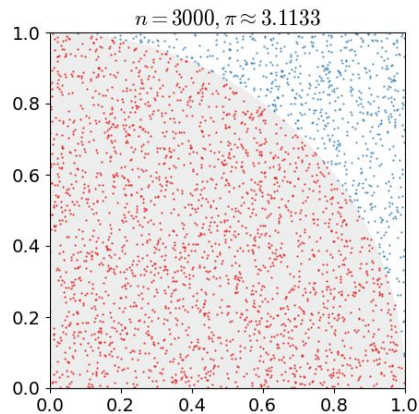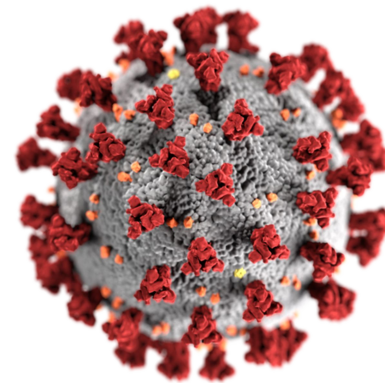
# Why did you sign up?



Figure 4: Initial gaussian wave with zero velocity. $\mu = 0.001$, $\varepsilon = 1$, $L = 1$, $N = 256$. cfl# = 0.1, $\sigma = 0.05$. Nonlinear shallow water.



$n = 3000$, $\pi \approx 3.1133$



N=2 LSA (not normalized)

# CSCI 285 Learning Goals

**Module #1: Data Analysis**

- Analyze & visualize data sets from a variety of sources.
- Learn several analysis techniques including clustering and regression.

**Module #2: Modeling**

- Model and solve system dynamics problems.
- Construct a Monte-Carlo simulation model.
- Develop agent-based models for complex simulations.

**Module #3: Numerical Techniques**

- Approximate the roots of continuous functions.
- Understand the strengths and limitations of numerical techniques.

Write idiomatic python and use scientific python libraries.

# CSCI 285 Course Overview

https://hendrix-cs.github.io/csci285/index.html

**Policies**
- Attendance
- Check ins / Office Hours
- Late Work

**Coursework / In-class**
- Lecture (36%)
- Labs (27%)
- Exams / Module Review (20%)
- Final Project (17%)

**More Info**
- Course Calendar / Class Notes / Project Timeline
- W2 Requirement
- Grading scale
- Prerequisites: MATH 130 & CSCI 150
- Teams - comms / submitting assignments

**Commitments**
- Active Participation
- Constructive Feedback
- Academic Integrity
- Learning Accommodation
- Physical & Mental Health

# CSCI 285 Grading Scale

# CSCI 285 Development Environment

1. Visit https://www.anaconda.com/

2. Download the open source distribution.

3. Follow the Anaconda3 installer instructions.

4. Launch Anaconda-Navigator (Mac, Windows, Linux)

5. Create new environments, launch processes, surf learning resources, etc.

(Alternatively, check out miniconda if you prefer a more lightweight approach)

# Module #1
# Data Analysis

# import pandas as pd

pandas is an open-source python library built for data manipulation and analysis. It is part of the standard library for many teams of data scientists and engineers. pandas introduces new types that have special syntax for data manipulation that are not shared with python's builtin types (e.g. list, dict). Some of the new syntax can look jarring at first, but is *lingua franca* for many data researchers.

**Getting Started with pandas**
- https://pandas.pydata.org/pandas-docs/stable/user_guide/10min.html
- https://chrisalbon.com/
- https://www.datacamp.com/courses/data-manipulation-with-pandas

# pd.Series

pandas

| | Series.1 |
|---|---|
| 0 | Item 1.1 |
| 1 | Item 1.2 |
| 2 | Item 1.3 |
| 3 | Item 1.4 |

| | Series.2 |
|---|---|
| 0 | Item 2.1 |
| 1 | Item 2.2 |
| 2 | Item 2.3 |
| 3 | Item 2.4 |

# pd.DataFrame

# Components of a DataFrame

## The Columns, Index, and Data

Columns



Index

Data

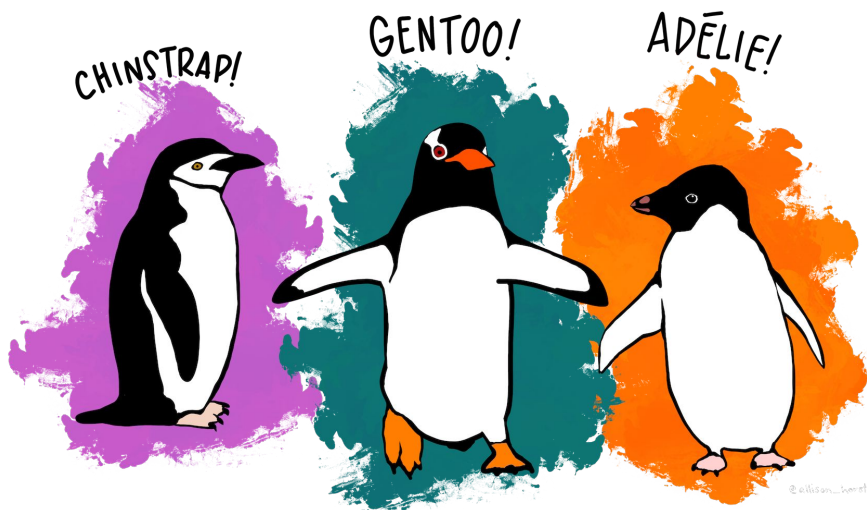| Description | Alternative Names | Axis Number |
|---|---|---|
| • Columns - label each column<br>• Index - label each row<br>• Data - actual values in DataFrame | • Columns - column names/labels, column index<br>• Index - index names/labels, row names/labels<br>• Data - values | • Columns: 1<br>• Index: 0 |

# Application: Pandas Intro

# Application: Palmer Penguins



Artwork by @allison_horst".

# Lab #1: Lake Trout

# Module #1
# Data Visualization

# Lab Report Format

- Palmer Penguins notebook (reference)
- Professor Wilson's trout lab (reference, next week)
- Mixture of Markdown, Code, and Figures
- Submitted via Teams (zip file)
- Must Include
    a. Any input data (data used to produce the report)
    b. Any output data (e.g. CSV files)
    c. Notebook with relative paths to load the data
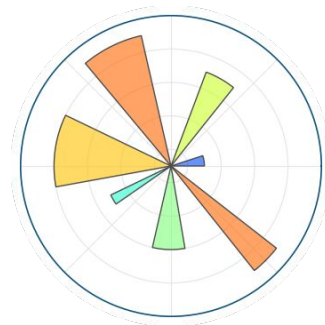
Any Questions about Lab #1?

# import seaborn as sns

Seaborn is a Python data visualization library based on matplotlib. It provides a

high-level interface for drawing attractive and informative statistical graphics

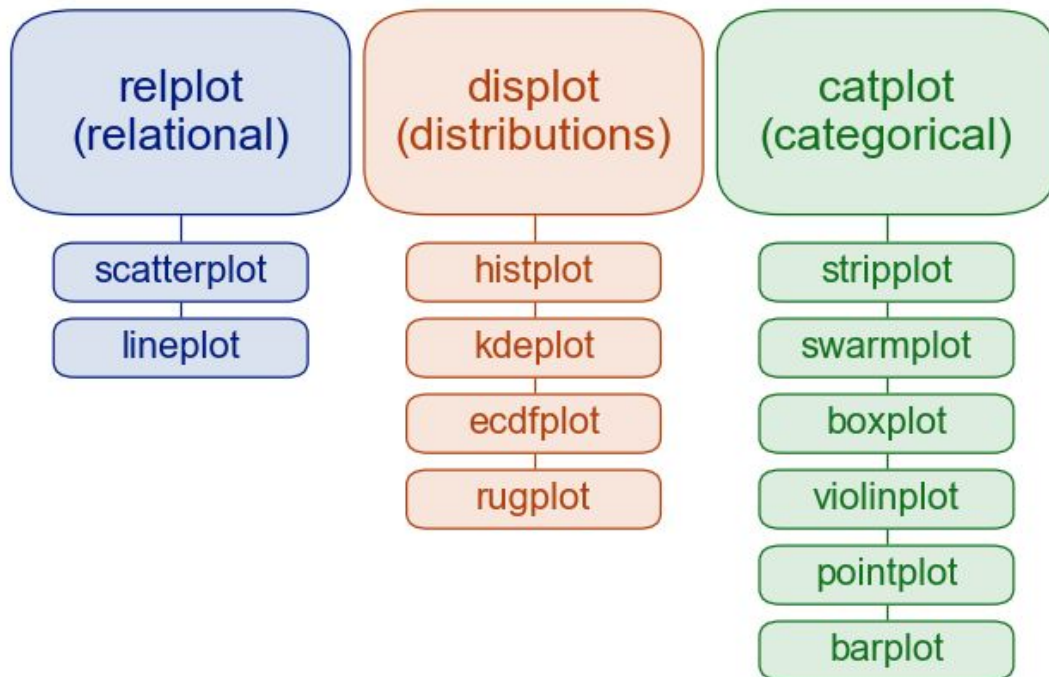**Getting Started with seaborn**

- https://seaborn.pydata.org/introduction.html
- https://seaborn.pydata.org/tutorial/function_overview.html
- https://chrisalbon.com/
- https://www.datacamp.com/courses/intermediate-data-visualization-with-seaborn

# import seaborn as sns

seaborn

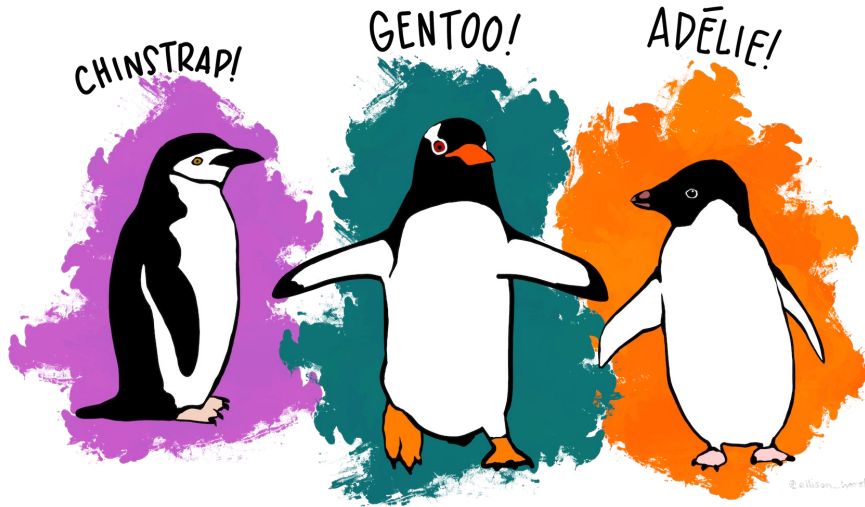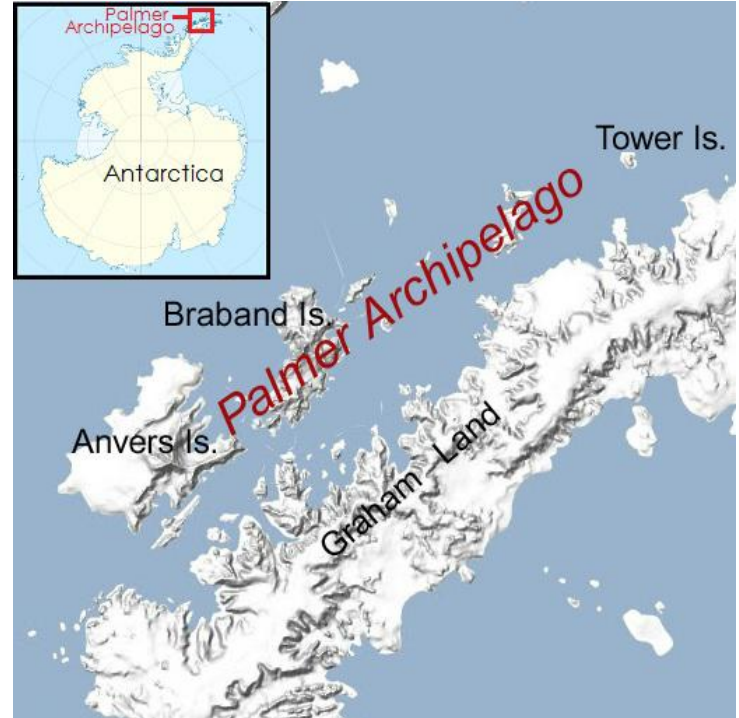| relplot (relational) | displot (distributions) | catplot (categorical) |
|---|---|---|
| scatterplot | histplot | stripplot |
| lineplot | kdeplot | swarmplot |
| | ecdfplot | boxplot |
| | rugplot | violinplot |
| | | pointplot |
| | | barplot |

# Application: Palmer Penguins



Artwork by @allison_horst".

# Final Project

- [Project Description](#)
- [COMAP](#)
- Markdown ([https://dillinger.io/](https://dillinger.io/))
- LaTeX / [Overleaf](#)

# Lab #2: Data Visualization

- Due Date: 9/12 (midnight).
- FEV notebook
- Linear Regression

# Module #1
# Machine Learning

# What is Machine Learning?

From Wikipedia:

*Machine learning, a branch of artificial intelligence, is about the construction and study of systems that can learn from data…*

*The core of machine learning deals with **representation** and **generalization***

- Representation == extracting structure from data
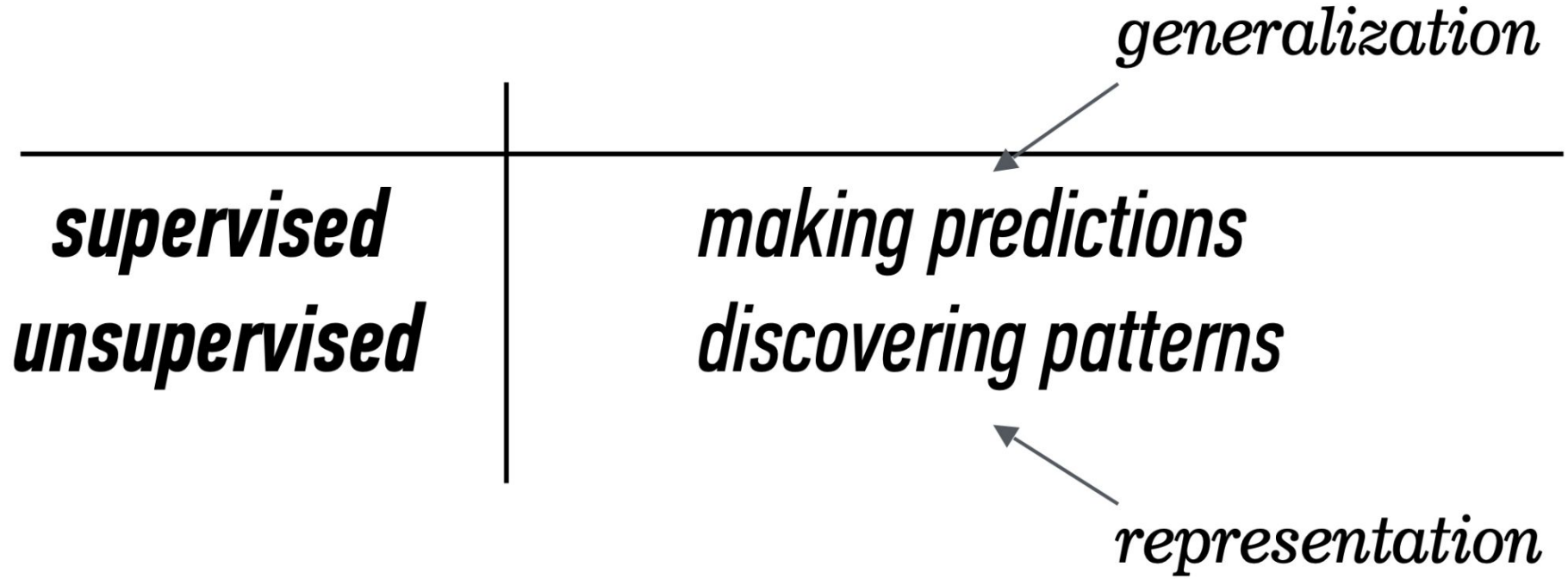
- Generalization == making predictions from data

# Types of Learning Problems

| | |
|---|---|
| *supervised* | *labeled examples* |
| *unsupervised* | *no labeled examples* |

# Types of Learning Problems

| | |
|---|---|
| *supervised* | *making predictions* |
| *unsupervised* | *discovering patterns* |

# Types of Learning Problems

*generalization*

supervised
unsupervised

*making predictions*
*discovering patterns*

*representation*

# Types of Data

|  | **continuous** | **categorical** |
|---|---|---|
|  | *quantitative* | *qualitative* |
|  | age, salary, height, etc. | city, yes/no, vote, etc. |

# Types of ML

|  | *continuous* | *categorical* |
|---|---|---|
| *supervised* | regression | classification |
| *unsupervised* | dimension reduction | clustering |

# from sklearn import

- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
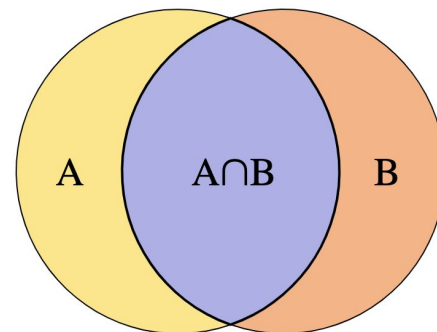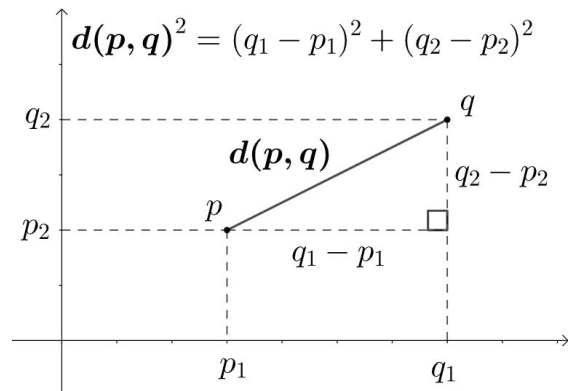- Open source, commercially usable - BSD license

**Getting Started with scikit-learn**
- https://scikit-learn.org/stable/index.html
- https://www.datacamp.com/courses/machine-learning-with-scikit-learn
- https://chrisalbon.com/

# Distance Measures (sklearn.metrics)

- Euclidean Space
  a. L1 Dist (manhattan)
  b. L2 Dist (pythagorean)
  c. LR Dist (general formulation)
- Non-Euclidean Space
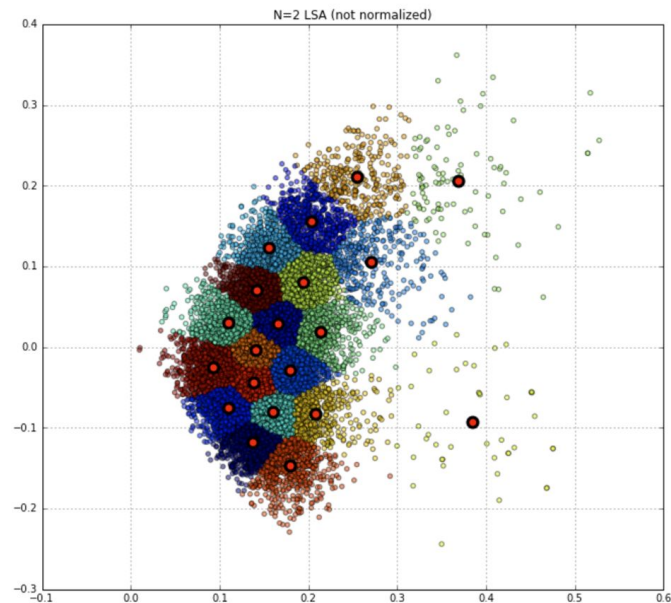  a. Jaccard Dist (sets)
  b. Edit Dist (strings)

$$d(p, q)^2 = (q_1 - p_1)^2 + (q_2 - p_2)^2$$
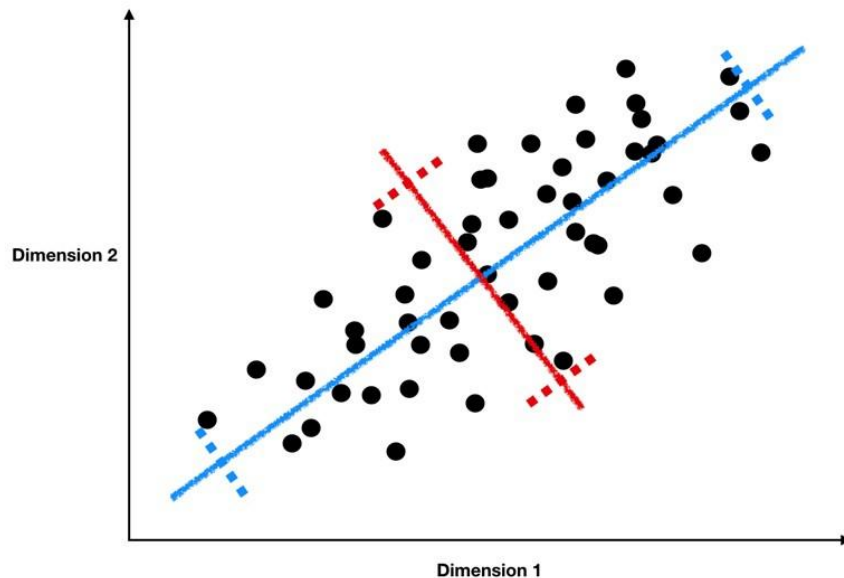
# Clustering (sklearn.cluster)

- K-Means Algorithm (step-by-step)
- scikit-learn K-Means module
- `from` `sklearn.datasets` `import` make_blobs
- Picking the right value for *k*
- Clustering Palmer Penguins
- Comparing predictions to ground truth
  - Confusion matrix
  - Seaborn heatmaps
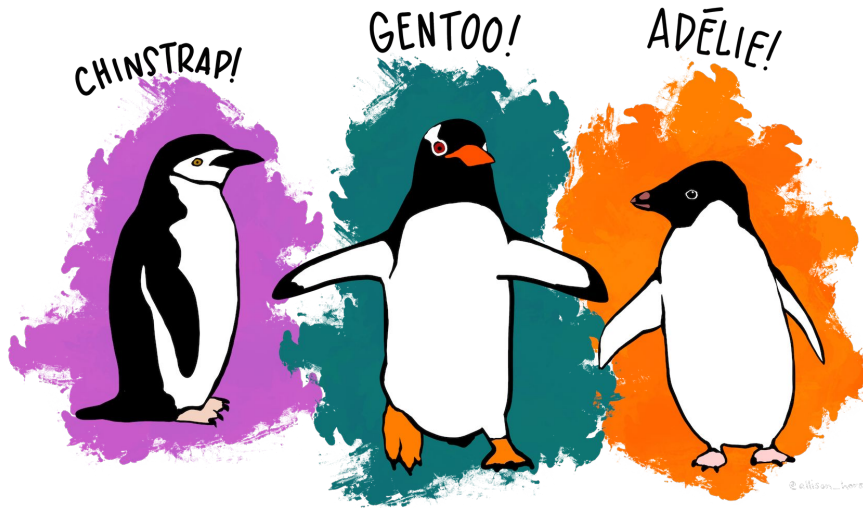


N=2 LSA (not normalized)

# Decomposition & Preprocessing
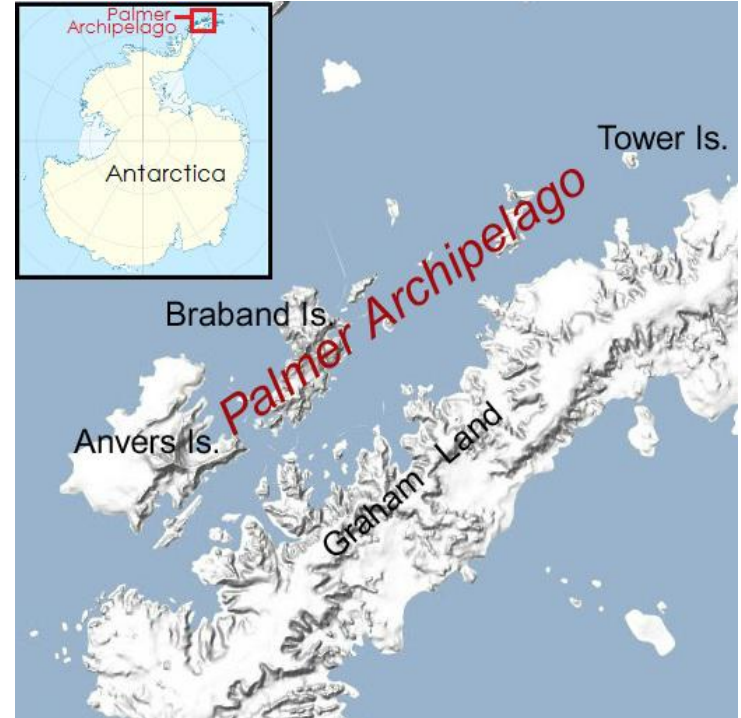
- Feature Scaling
  - StandardScaler
- Curse of dimensionality
- PCA (principal component analysis)

# Application: Palmer Penguins



CHINSTRAP! GENTOO! ADÉLIE!

Artwork by @allison_horst".

# from sklearn.datasets

- [sklearn.datasets](sklearn.datasets)

# from sklearn.linear_model

- [sklearn.linear_model](#)
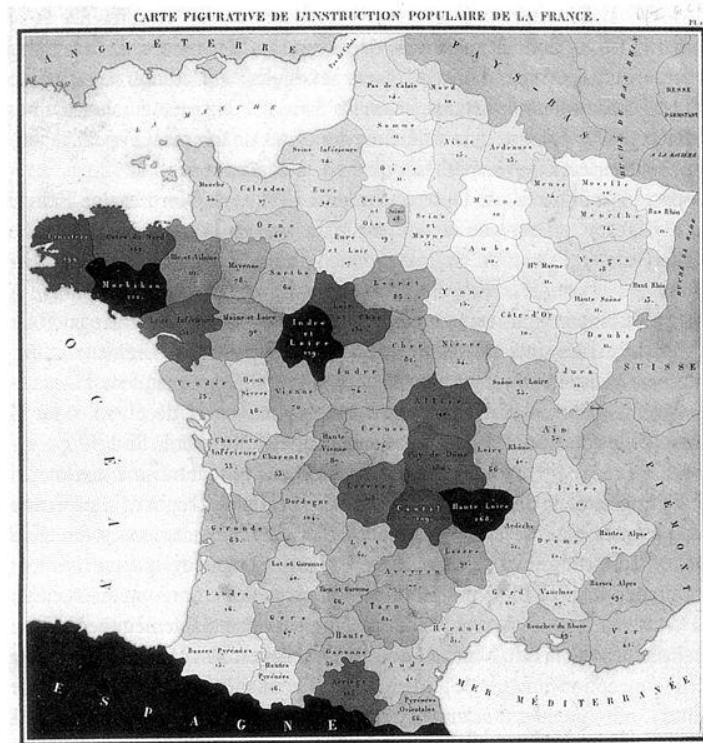


Flipper Length (mm) vs Bill Length (mm)

# Module #1
## Choropleths & More Data Viz

# Module #1 Timeline

- More Data Visualizations (today!)

- Module #1 Review / Exam Prep (9/20)
- Class Cancelled (9/22)

- Exam #1 (9/27)
- Lab #4 (9/29)

- System Dynamics Modeling (10/4)

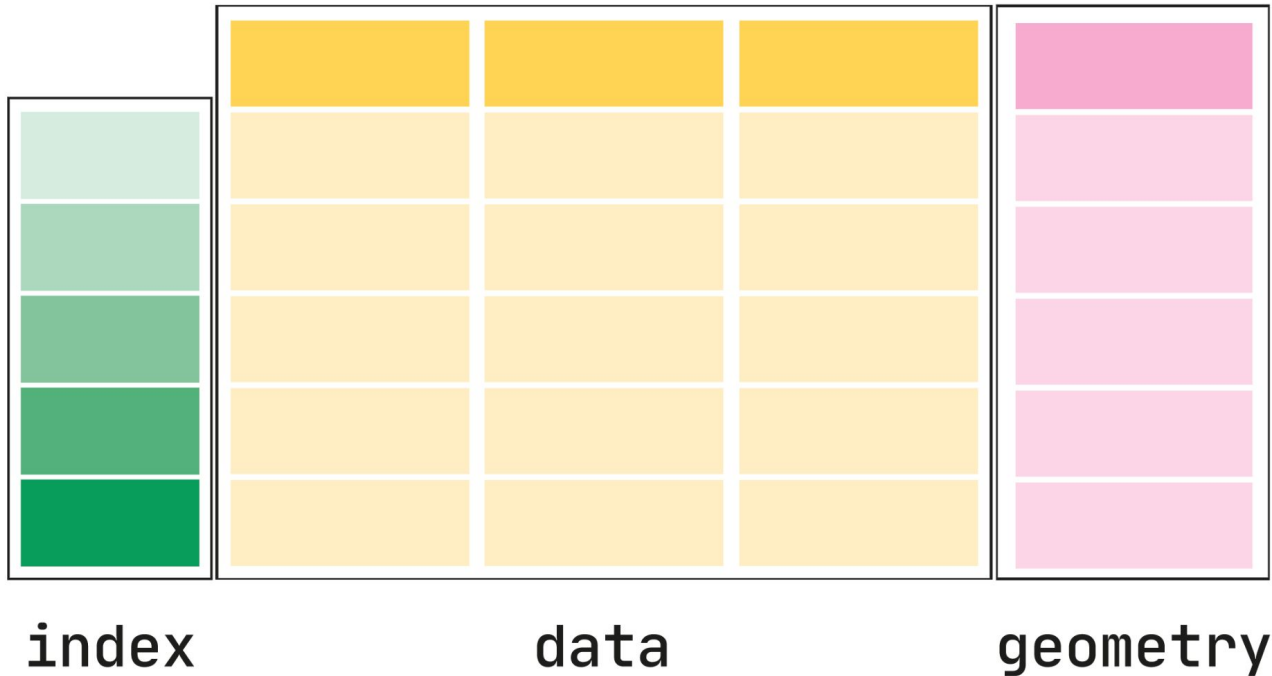# Choropleths

# import geopandas as gpd

GeoPandas, as the name suggests, extends the popular data science library pandas by adding support for geospatial data.

- pandas.DataFrame → geopandas.GeoDataFrame
- pandas.Series → geopandas.GeoSeries

**Getting Started with geopandas**
- https://geopandas.org/en/stable/getting_started/introduction.html
- https://www.kaggle.com/

# import geopandas as gpd

GeoPandas

index      data      geometry

# Barcharts & Heatmaps

seaborn

# Exam #1: 9/27 (Next Tuesday)

- A: **Complete** two **Exams** and at least **Partially Complete** the remaining **Exam**.
- B: **Complete** one **Exam** and **Partially Complete** the remaining **Exams**.
- C: **Partially Complete** all three **Exams**.
- D: **Partially Complete** two **Exams**.

1. Tuesday, 9/27 Exam #1 will be "opened"
2. Due by beginning of class on Thursday (9/29)

# CSCI 285 Learning Goals

**Module #1: Data Analysis**

- Analyze & visualize data sets from a variety of sources.
- Learn several analysis techniques including clustering and regression.

**Module #2: Modeling**

- Model and solve system dynamics problems.
- Construct a Monte-Carlo simulation model.
- Develop agent-based models for complex simulations.

**Module #3: Numerical Techniques**

- Approximate the roots of continuous functions.
- Understand the strengths and limitations of numerical techniques.

Write idiomatic python and use scientific python libraries.

# Module #2
# System Dynamics