

Section 9 Introduction to Data Science

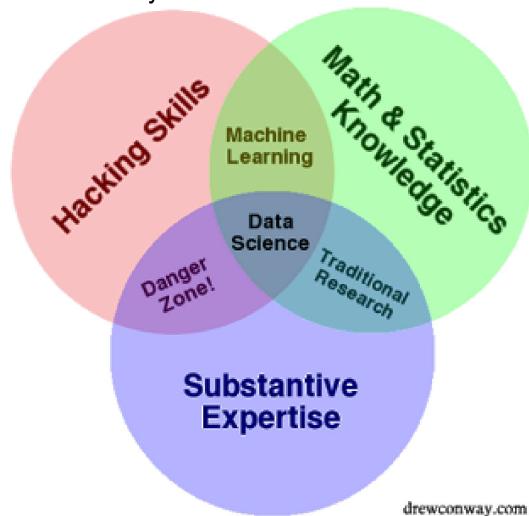
What is Data Science?

Three correlated concepts:

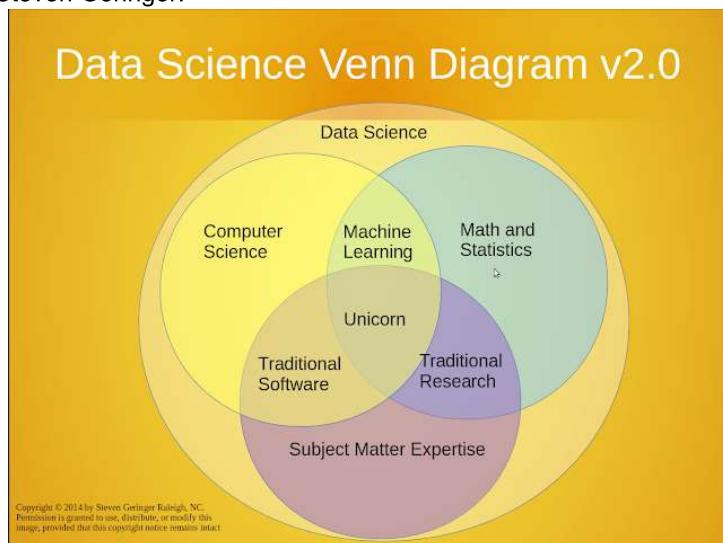
- Data Science
- Artificial Intelligence
- Machine Learning

[Battle of the Data Science Venn Diagrams \(<https://www.kdnuggets.com/2016/10/battle-data-science-venn-diagrams.html>\)](https://www.kdnuggets.com/2016/10/battle-data-science-venn-diagrams.html)

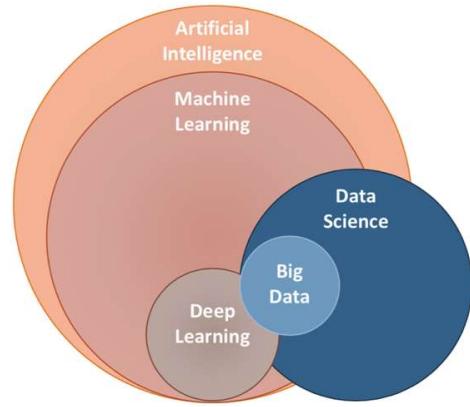
The original Venn diagram from Drew Conway:



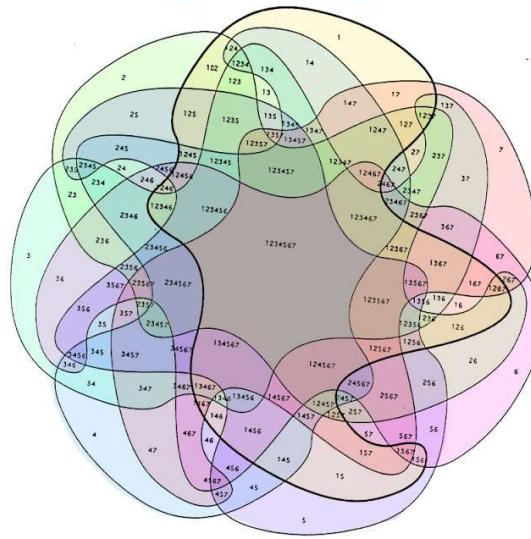
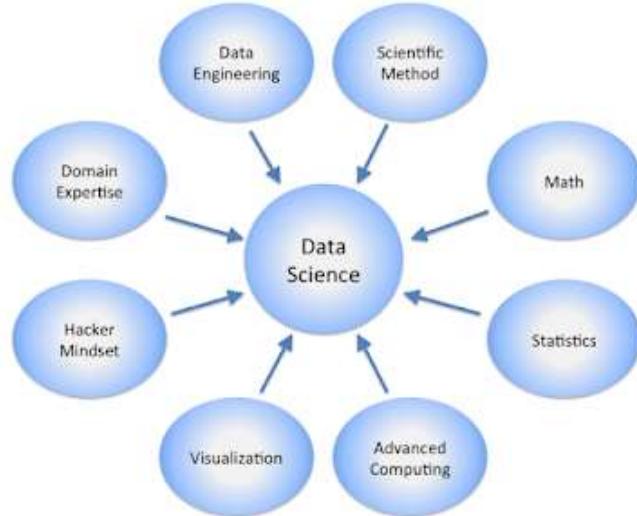
Another diagram from Steven Geringer:



Another version:



Perhaps the reality should be:



[David Robinson's Auto-pilot example \(<http://varianceexplained.org/r/ds-ml-ai/>\):](http://varianceexplained.org/r/ds-ml-ai/)

- machine learning: **predict** whether there is a stop sign in the camera
- artificial intelligence: design the **action** of applying brakes (either by rules or from data)
- data science: provide the **insights** why the system does not work well after sunrise

Peijie's Definition: Data Science is the science

- of the data -- what

- *by* the data -- how
- *for* the data -- why

Mathematics of Data

Representation of Data

What data do we have, and how to relate it with math objects?

Tabular Data

```
In [6]: ┆ import pandas as pd
      import numpy as np
      df_house = pd.read_csv('./data/kc_house_data.csv')
      print(df_house.shape)
      df_house.head()
```

(21613, 21)

Out[6]:

| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront |
|---|------------|-----------------|----------|----------|-----------|-------------|----------|--------|------------|
| 0 | 7129300520 | 20141013T000000 | 221900.0 | 3 | 1.00 | 1180 | 5650 | 1.0 | 0 |
| 1 | 6414100192 | 20141209T000000 | 538000.0 | 3 | 2.25 | 2570 | 7242 | 2.0 | 0 |
| 2 | 5631500400 | 20150225T000000 | 180000.0 | 2 | 1.00 | 770 | 10000 | 1.0 | 0 |
| 3 | 2487200875 | 20141209T000000 | 604000.0 | 4 | 3.00 | 1960 | 5000 | 1.0 | 0 |
| 4 | 1954400510 | 20150218T000000 | 510000.0 | | 3 | 2.00 | 1680 | 8080 | 1.0 |

5 rows × 21 columns

- A structured data table, with n observations and p variables.
- **Mathematical representation:** The data matrix $X \in \mathbb{R}^{n \times p}$. For notations we write

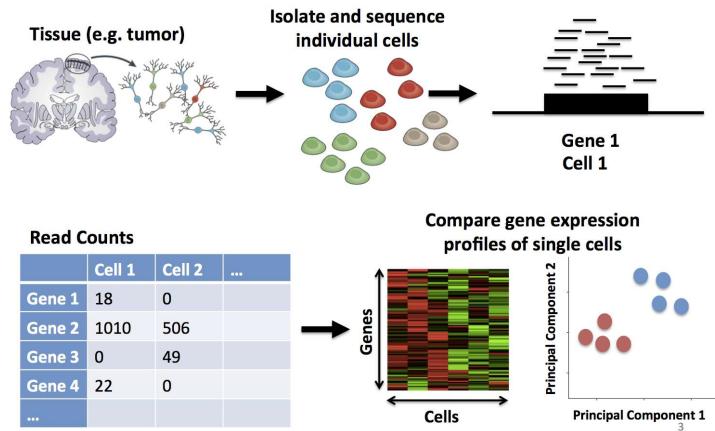
$$X = \begin{pmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(n)} \end{pmatrix}, \text{ where the } i\text{-th row vector represents } i\text{-th observation, } \mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)}) \in \mathbb{R}^p.$$

To really emphasize that each element is a row, we can also write X as:

$$X = \begin{pmatrix} \mathbf{x}^{(1)} \longrightarrow \\ \mathbf{x}^{(2)} \longrightarrow \\ \vdots \\ \mathbf{x}^{(n)} \longrightarrow \end{pmatrix}$$

- Example: Precision Medicine and Single-cell Sequencing. (<https://learn.gencore.bio.nyu.edu/single-cell-rnaseq/>).

Single-cell RNA-Seq (scRNA-Seq)



- Roughly speaking, big data -- large n , high-dimensional data -- large p .

Time-series Data

```
In [7]: # import matplotlib.pyplot as plt
ts_tesla = pd.read_csv('./data/Tesla.csv')
print(ts_tesla.head())

ts_tesla['Date'] = pd.to_datetime(ts_tesla['Date'])
ts_tesla.set_index('Date', inplace=True)

# Suppose we only focus on the time-series of close price
plt.figure(dpi=80)
plt.title('Closing Price History')
plt.plot(ts_tesla['Close'], color='red')
plt.xlabel('Date', fontsize=18)
plt.ylabel('Closing Price USD', fontsize = 18)
plt.show()
# this is only about tesla -- we can also have the time-series of apple,amazon,facebook.
```

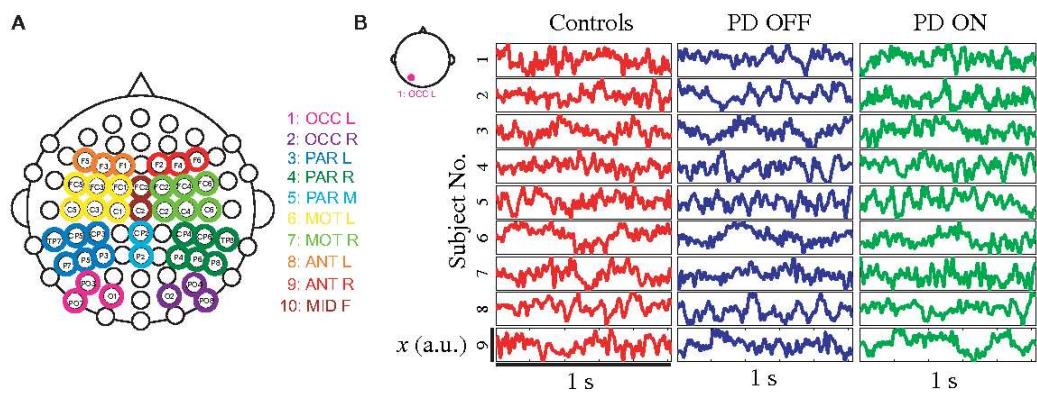
| | Date | Open | High | Low | Close | Volume | Adj Close |
|---|-----------|-----------|-------|-----------|-----------|----------|-----------|
| 0 | 6/29/2010 | 19.000000 | 25.00 | 17.540001 | 23.889999 | 18766300 | 23.889999 |
| 1 | 6/30/2010 | 25.790001 | 30.42 | 23.299999 | 23.830000 | 17187100 | 23.830000 |
| 2 | 7/1/2010 | 25.000000 | 25.92 | 20.270000 | 21.959999 | 8218800 | 21.959999 |
| 3 | 7/2/2010 | 23.000000 | 23.10 | 18.709999 | 19.200001 | 5139800 | 19.200001 |
| 4 | 7/6/2010 | 20.000000 | 20.00 | 15.830000 | 16.110001 | 6866900 | 16.110001 |



- Simple case: N one-dimensional trajectories with each sampled at T time points.
- **Mathematical representation I:** Still use the data matrix $X \in \mathbb{R}^{N \times T}$. For notations we write

$$X = \begin{pmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(N)} \end{pmatrix}, \text{ where the } i\text{-th row vector represents } i\text{-th trajectory, } \mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_T^{(i)}) \in \mathbb{R}^T.$$

- Question: The difference with tabular data?
- **Mathematical representation II:** Each trajectory is a *function* of time t . The whole dataset can be represented as $z = f(\omega, t)$ where ω represents the sample and t represents the time. In probability theory, this is called *stochastic process*.
 - For fixed ω , we have a trajectory, which is the function of time.
 - For fixed t , we obtain an ensemble drawn from random distribution.
- Question: How about N d -dimensional trajectories with each sampled at T time points?
- Example: Electroencephalography (EEG) data and Parkinson's disease (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3858815/>).



Images

Example: [MNIST handwritten digits data \(<http://yann.lecun.com/exdb/mnist/>\)](http://yann.lecun.com/exdb/mnist/): Each image is 28x28 matrix

```
In [16]: ➤ import pandas as pd
mnist = pd.read_csv('./data/train.csv') # stored as data table
mnist.sample(5)
```

Out[16]:

| | label | pixel0 | pixel1 | pixel2 | pixel3 | pixel4 | pixel5 | pixel6 | pixel7 | pixel8 | ... | pixel774 | pixel775 | pix |
|-------|-------|--------|--------|--------|--------|--------|--------|--------|--------|--------|-----|----------|----------|-----|
| 4673 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | |
| 7588 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | |
| 15050 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | |
| 41202 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | |
| 5520 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | |

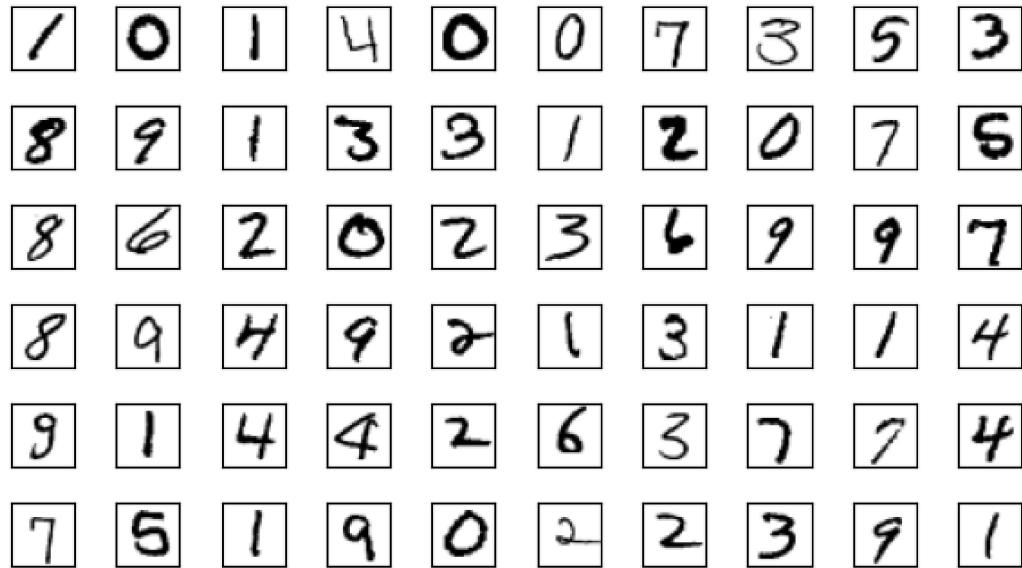
5 rows × 785 columns

```
In [17]: ➤ mnist.shape
```

Out[17]: (42000, 785)

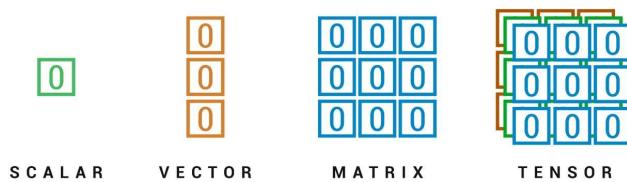
```
In [18]: # target = mnist['label']
# mnist = mnist.drop("label",axis=1)

import matplotlib.pyplot as plt
plt.figure(dpi=100)
for i in range(0,70): #plot the first 70 images
    plt.subplot(7,10,i+1)
    grid_data = mnist.iloc[i,:].to_numpy().reshape(28,28) # reshape from 1d to 2d pixel
    plt.imshow(grid_data,cmap='gray_r', vmin=0, vmax=255)
    plt.xticks([])
    plt.yticks([])
plt.tight_layout()
```



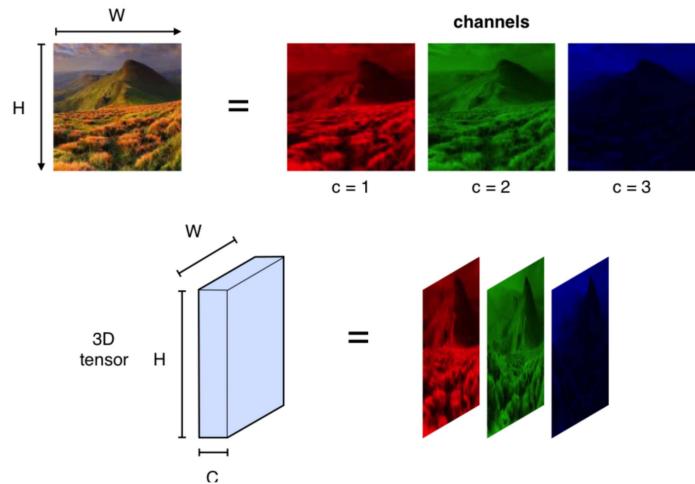
- Simple case: N grayscale images with $m \times n$ pixels each.
- **Mathematical Representation I:** Each image can be represented by a matrix $I \in \mathbb{R}^{m \times n}$, whose elements denotes the intensities of pixels. The whole datasets have N matrices of m by n , or represented by a $N \times m \times n$ tensor.

[Illustrated Introduction to Linear Algebra using NumPy](https://medium.com/@kaaanishk/illustrated-introduction-to-linear-algebra-using-NumPy_(https://medium.com/@kaaanishk/illustrated-introduction-to-linear-algebra-using-numpy-11d503d244a1)) (<https://medium.com/@kaaanishk/illustrated-introduction-to-linear-algebra-using-numpy-11d503d244a1>).



- **Mathematical representation II:** Random field model $z = \mathbf{f}(\omega, x, y)$.
- **Color images:** Decompose into RGB (red, green and blue) channels and
 - use three matrices (or three-dimensional tensor) to represent one image, or
 - build the random field model with vector-valued functions $z = \mathbf{f}(\omega, x, y) \in \mathbb{R}^3$

[convolutional neural networks](https://www.esantus.com/blog/2019/1/31/convolutional-neural-networks-a-quick-guide-for-newbies) (<https://www.esantus.com/blog/2019/1/31/convolutional-neural-networks-a-quick-guide-for-newbies>).



- Question: Can image datasets also be transformed into tabular data? What are the pros/cons?

In [19]: `mnist.head()`

Out[19]:

| | pixel0 | pixel1 | pixel2 | pixel3 | pixel4 | pixel5 | pixel6 | pixel7 | pixel8 | pixel9 | ... | pixel774 | pixel775 | pixel7 |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|-----|----------|----------|--------|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 |

5 rows \times 784 columns

Videos

- Time-series of images, or random field model $z = \mathbf{f}(\omega, x, y, t)$

Texts

```
In [3]: └─▶ from sklearn.feature_extraction.text import CountVectorizer

corpus = ['He is a good person',
          'He is bad student',
          'He is hardworking']
df = pd.DataFrame(data=corpus, columns=['sentences'])
print(df)
vectorizer = CountVectorizer(vocabulary=['he', 'is', 'a', 'good', 'person', 'bad', 'student', 'hardworking'],
                             stop_words=frozenset(), token_pattern=r"(?u)\b\w+\b")
X = vectorizer.fit_transform(df['sentences'].values)
result = pd.DataFrame(data=X.toarray(), columns=vectorizer.get_feature_names())
result.head()
```

sentences

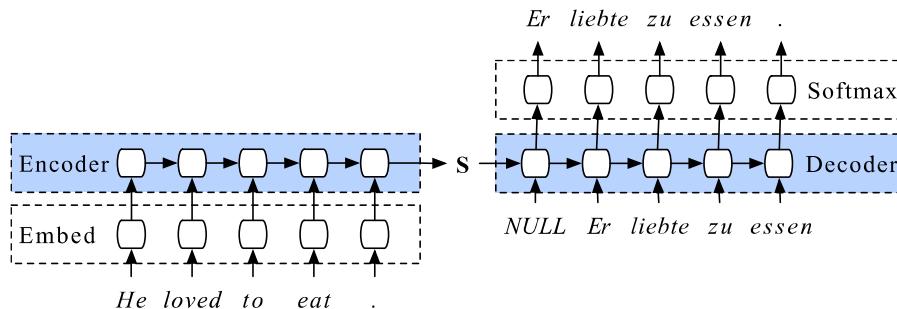
| | |
|---|---------------------|
| 0 | He is a good person |
| 1 | He is bad student |
| 2 | He is hardworking |

Out[3]:

| | he | is | a | good | person | bad | student | hardworking |
|---|----|----|---|------|--------|-----|---------|-------------|
| 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| 2 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

- **Proposal I:** Tabular data by extracting key words. "Document-Term Matrix"
 - useful in sentiment analysis, document clustering, topic modelling
 - popular algorithms include tf-idf, Word2Vec, bag of words, etc.
- **Proposal II:** Time-series of individual words.
 - useful in machine translation

[Recurrent neural network model for machine translations
\(\[https://smerity.com/articles/2016/google_nmt_arch.html\]\(https://smerity.com/articles/2016/google_nmt_arch.html\)\)](https://smerity.com/articles/2016/google_nmt_arch.html)



Networks

- Concepts: node/edge/weight, directed/undirected
- **Mathematical Representation:** adjacency matrix
- Question: what about the whole datasets of networks, and time-evolving networks?

Tasks with Data: Machine Learning

The tasks with data can often be transformed into *machine learning* problems, which can be generally classified as:

- Supervised Learning -- "learning with training";
- Unsupervised Learning -- "learning without training";
- Reinforcement Learning -- "learning by doing".

Our course will focus on the first two categories.

Supervised Learning

- Given the *training dataset* $(x^{(i)}, y^{(i)})$ with $y^{(i)} \in \mathbb{R}^q$ denotes the *labels*, the supervised learning aims to find a mapping $\mathbf{f} : \mathbb{R}^p \rightarrow \mathbb{R}^q$ such that $y^{(i)} \approx \mathbf{f}(x^{(i)})$. Then with a new observation $x^{(new)}$, we can predict that $y^{(new)} = \mathbf{f}(x^{(new)})$.
 - when $y \in \mathbb{R}$ is continuous, the problem is also called as *regression*. **Example:** Housing price prediction
 - when $y \in \mathbb{R}$ is discrete, the problem is also called as *classification*. **Example:** Handwritten digit recognition
- **Practical Strategy:** Limit the mapping \mathbf{f} to certain space by parametrization $\mathbf{f}(\mathbf{x}; \theta)$. Then define the loss function of θ

$$L(\theta) = \sum_{i=1}^n \ell(y^{(i)}, \mathbf{f}(x^{(i)})),$$

where ℓ quantifies the "distance" between $y^{(i)}$ and $\mathbf{f}(x^{(i)})$, and a common choice is mean square error (MSE) for continuous data $\ell(y^{(i)}, \mathbf{f}(x^{(i)})) = \|y^{(i)} - \mathbf{f}(x^{(i)})\|^2$. We then seek to choose the optimal θ that minimizes the loss function

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta),$$

which can be tackled numerically by optimization methods (including the popular stochastic gradient descent).

- Difference choice of $\mathbf{f}(\mathbf{x}; \theta)$ leads to various supervised learning models:
 - Linear function : Linear Regression (for regression)/Logistic Regression (for classification)
 - Composition of linear + nonlinear functions: Neural Network
- **Important Terms:**
 - **Training Data:** Both X and y are provided. The dataset which we use to fit the function.
 - **Test Data:** In principle, only X is provided (some times y^{test} is also provided as the ground-truth to verify). The dataset which we generate new predictions y^{pred} . -- This is the final judgement of your unsupervised ML model!
 - **Validation Data:** A good-fit model on training data does not guarantee the good performance on test data. To gain more confidence before really applying to test data, we "fake" some test data as the "sample exam". To do this, we further split the original training data into new training data and validation data, and then learn the mapping on new training data, and judge on the validation data. We may make some adjustment if the model does not perform well in the "sample exam".
 - Intuitive Understanding: Training data is like quizzes -- you want to learn the "mapping" between the question and correct answer. Test data is like your exam. Validation is like you take a sample exam before the real exam and make some "clinics" about your weakpoints.
 - See the illustration [here](https://towardsdatascience.com/train-validation-and-test-sets-72cb40cba9e7) (<https://towardsdatascience.com/train-validation-and-test-sets-72cb40cba9e7>).

Example: The [Wisconsin breast cancer dataset](#)

([https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))) and low-code ML package [pycaret](https://pycaret.org/) (<https://pycaret.org/>).

Install pycaret -- it's a new package, not included with Anaconda

In [33]: ➜ pip install --upgrade pycaret

```
Requirement already satisfied: querystring-parser in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (1.2.4)
Requirement already satisfied: packaging in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (20.9)
Requirement already satisfied: protobuf>=3.7.0 in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (3.17.3)
Requirement already satisfied: click>=7.0 in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (7.1.2)
Requirement already satisfied: databricks-cli>=0.8.7 in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (0.14.3)
Requirement already satisfied: sqlalchemy in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (1.4.7)
Requirement already satisfied: Flask in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (1.1.2)
Requirement already satisfied: entrypoints in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (0.3)
Requirement already satisfied: sqlparse>=0.3.1 in c:\users\lukea\anaconda3\lib\site-packages (from mlflow->pycaret) (0.4.1)
Requirement already satisfied: Mako in c:\users\lukea\anaconda3\lib\site-packages (from alembic<=1.4.1->mlflow->pycaret) (1.1.4)
```

In [9]: ➜ `from sklearn.datasets import load_breast_cancer # Load the dataset
X,y = load_breast_cancer(as_frame = True,return_X_y = True)`

In [10]: ➜ X

Out[10]:

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | me frac dimensi |
|-----|-------------|--------------|----------------|-----------|-----------------|------------------|----------------|---------------------|---------------|-----------------|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.30010 | 0.14710 | 0.2419 | 0.078 |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.08690 | 0.07017 | 0.1812 | 0.056 |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.19740 | 0.12790 | 0.2069 | 0.059 |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.24140 | 0.10520 | 0.2597 | 0.097 |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.19800 | 0.10430 | 0.1809 | 0.058 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 564 | 21.56 | 22.39 | 142.00 | 1479.0 | 0.11100 | 0.11590 | 0.24390 | 0.13890 | 0.1726 | 0.056 |
| 565 | 20.13 | 28.25 | 131.20 | 1261.0 | 0.09780 | 0.10340 | 0.14400 | 0.09791 | 0.1752 | 0.055 |
| 566 | 16.60 | 28.08 | 108.30 | 858.1 | 0.08455 | 0.10230 | 0.09251 | 0.05302 | 0.1590 | 0.056 |
| 567 | 20.60 | 29.33 | 140.10 | 1265.0 | 0.11780 | 0.27700 | 0.35140 | 0.15200 | 0.2397 | 0.070 |
| 568 | 7.76 | 24.54 | 47.92 | 181.0 | 0.05263 | 0.04362 | 0.00000 | 0.00000 | 0.1587 | 0.058 |

569 rows × 30 columns

```
In [11]: █ y
```

```
Out[11]: 0      0
         1      0
         2      0
         3      0
         4      0
         ..
        564     0
        565     0
        566     0
        567     0
        568     1
Name: target, Length: 569, dtype: int32
```

In this dataset, all labels are known. To mimic a real situation, we manually create train and test datasets.

```
In [12]: █ from sklearn.model_selection import train_test_split # manually split into train and test
          X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)
```

```
In [13]: █ X_train.shape
```

```
Out[13]: (381, 30)
```

```
In [14]: █ y_test.shape
```

```
Out[14]: (188,)
```

```
In [15]: ➜ import pandas as pd  
data_train = pd.concat([X_train,y_train],axis=1) # the whole data table of training  
data_train
```

Out[15]:

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimens |
|-----|-------------|--------------|----------------|-----------|-----------------|------------------|----------------|---------------------|---------------|---------------------|
| 56 | 19.210 | 18.57 | 125.50 | 1152.0 | 0.10530 | 0.12670 | 0.13230 | 0.089940 | 0.1917 | 0.058 |
| 144 | 10.750 | 14.97 | 68.26 | 355.3 | 0.07793 | 0.05139 | 0.02251 | 0.007875 | 0.1399 | 0.056 |
| 60 | 10.170 | 14.88 | 64.55 | 311.9 | 0.11340 | 0.08061 | 0.01084 | 0.012900 | 0.2743 | 0.068 |
| 6 | 18.250 | 19.98 | 119.60 | 1040.0 | 0.09463 | 0.10900 | 0.11270 | 0.074000 | 0.1794 | 0.057 |
| 8 | 13.000 | 21.82 | 87.50 | 519.8 | 0.12730 | 0.19320 | 0.18590 | 0.093530 | 0.2350 | 0.073 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 277 | 18.810 | 19.98 | 120.90 | 1102.0 | 0.08923 | 0.05884 | 0.08020 | 0.058430 | 0.1550 | 0.048 |
| 9 | 12.460 | 24.04 | 83.97 | 475.9 | 0.11860 | 0.23960 | 0.22730 | 0.085430 | 0.2030 | 0.082 |
| 359 | 9.436 | 18.32 | 59.82 | 278.6 | 0.10090 | 0.05956 | 0.02710 | 0.014060 | 0.1506 | 0.068 |
| 192 | 9.720 | 18.22 | 60.73 | 288.1 | 0.06950 | 0.02344 | 0.00000 | 0.000000 | 0.1653 | 0.064 |
| 559 | 11.510 | 23.93 | 74.52 | 403.5 | 0.09261 | 0.10210 | 0.11120 | 0.041050 | 0.1388 | 0.068 |

381 rows × 31 columns

In [16]:

```
from pycaret.classification import setup
from pycaret.classification import compare_models

bc = setup(data=data_train, target='target') # target is the y column name we want to predict
```

| | Description | Value |
|----|--|------------------|
| 0 | session_id | 1198 |
| 1 | Target | target |
| 2 | Target Type | Binary |
| 3 | Label Encoded | None |
| 4 | Original Data | (381, 31) |
| 5 | Missing Values | False |
| 6 | Numeric Features | 30 |
| 7 | Categorical Features | 0 |
| 8 | Ordinal Features | False |
| 9 | High Cardinality Features | False |
| 10 | High Cardinality Method | None |
| 11 | Transformed Train Set | (266, 29) |
| 12 | Transformed Test Set | (115, 29) |
| 13 | Shuffle Train-Test | True |
| 14 | Stratify Train-Test | False |
| 15 | Fold Generator | StratifiedKFold |
| 16 | Fold Number | 10 |
| 17 | CPU Jobs | -1 |
| 18 | Use GPU | False |
| 19 | Log Experiment | False |
| 20 | Experiment Name | clf-default-name |
| 21 | USI | 8ad9 |
| 22 | Imputation Type | simple |
| 23 | Iterative Imputation Iteration | None |
| 24 | Numeric Imputer | mean |
| 25 | Iterative Imputation Numeric Model | None |
| 26 | Categorical Imputer | constant |
| 27 | Iterative Imputation Categorical Model | None |
| 28 | Unknown Categoricals Handling | least_frequent |
| 29 | Normalize | False |
| 30 | Normalize Method | None |
| 31 | Transformation | False |
| 32 | Transformation Method | None |
| 33 | PCA | False |
| 34 | PCA Method | None |
| 35 | PCA Components | None |

| | Description | Value |
|----|------------------------------|---------|
| 36 | Ignore Low Variance | False |
| 37 | Combine Rare Levels | False |
| 38 | Rare Level Threshold | None |
| 39 | Numeric Binning | False |
| 40 | Remove Outliers | False |
| 41 | Outliers Threshold | None |
| 42 | Remove Multicollinearity | False |
| 43 | Multicollinearity Threshold | None |
| 44 | Remove Perfect Collinearity | True |
| 45 | Clustering | False |
| 46 | Clustering Iteration | None |
| 47 | Polynomial Features | False |
| 48 | Polynomial Degree | None |
| 49 | Trigonometry Features | False |
| 50 | Polynomial Threshold | None |
| 51 | Group Features | False |
| 52 | Feature Selection | False |
| 53 | Feature Selection Method | classic |
| 54 | Features Selection Threshold | None |
| 55 | Feature Interaction | False |
| 56 | Feature Ratio | False |
| 57 | Interaction Threshold | None |
| 58 | Fix Imbalance | False |
| 59 | Fix Imbalance Method | SMOTE |

```
In [17]: ┏ best = compare_models() # pycaret automatically fits different ML models for you, and co
```

| | Model | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC | TT (Sec) |
|-----------------|---------------------------------|----------|--------|--------|--------|--------|--------|--------|-------------|
| et | Extra Trees Classifier | 0.9590 | 0.9868 | 0.9750 | 0.9583 | 0.9659 | 0.9144 | 0.9168 | 0.0500 |
| ada | Ada Boost Classifier | 0.9588 | 0.9756 | 0.9812 | 0.9539 | 0.9664 | 0.9132 | 0.9168 | 0.0300 |
| qda | Quadratic Discriminant Analysis | 0.9548 | 0.9888 | 0.9688 | 0.9598 | 0.9632 | 0.9048 | 0.9078 | 0.0060 |
| rf | Random Forest Classifier | 0.9473 | 0.9819 | 0.9562 | 0.9587 | 0.9550 | 0.8915 | 0.8971 | 0.0620 |
| nb | Naive Bayes | 0.9440 | 0.9864 | 0.9688 | 0.9412 | 0.9537 | 0.8826 | 0.8867 | 0.0050 |
| ridge | Ridge Classifier | 0.9439 | 0.0000 | 0.9938 | 0.9221 | 0.9558 | 0.8795 | 0.8869 | 0.0050 |
| lightgbm | Light Gradient Boosting Machine | 0.9436 | 0.9835 | 0.9562 | 0.9525 | 0.9527 | 0.8828 | 0.8871 | 0.1170 |
| lda | Linear Discriminant Analysis | 0.9363 | 0.9815 | 0.9812 | 0.9216 | 0.9495 | 0.8637 | 0.8696 | 0.0050 |
| gbc | Gradient Boosting Classifier | 0.9360 | 0.9855 | 0.9375 | 0.9560 | 0.9459 | 0.8677 | 0.8699 | 0.0420 |
| lr | Logistic Regression | 0.9359 | 0.9845 | 0.9625 | 0.9355 | 0.9479 | 0.8648 | 0.8683 | 0.4830 |
| knn | K Neighbors Classifier | 0.9172 | 0.9499 | 0.9562 | 0.9148 | 0.9331 | 0.8249 | 0.8321 | 0.0090 |
| dt | Decision Tree Classifier | 0.9138 | 0.9135 | 0.9188 | 0.9399 | 0.9279 | 0.8208 | 0.8246 | 0.0050 |
| svm | SVM - Linear Kernel | 0.8533 | 0.0000 | 0.8625 | 0.9100 | 0.8731 | 0.6974 | 0.7264 | 0.0050 |

```
In [18]: ┏ best # the best model selected by pycaret
```

```
Out[18]: ExtraTreesClassifier(bootstrap=False, ccp_alpha=0.0, class_weight=None,
                           criterion='gini', max_depth=None, max_features='auto',
                           max_leaf_nodes=None, max_samples=None,
                           min_impurity_decrease=0.0, min_impurity_split=None,
                           min_samples_leaf=1, min_samples_split=2,
                           min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=-1,
                           oob_score=False, random_state=1198, verbose=0,
                           warm_start=False)
```

```
In [19]: ┏ from pycaret.classification import predict_model
predict_model(best); # predict on the validation data that pycaret have selected -- samp
```

| | Model | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC |
|----------|------------------------|----------|--------|--------|--------|--------|--------|--------|
| 0 | Extra Trees Classifier | 0.9565 | 0.9963 | 0.9474 | 0.9863 | 0.9664 | 0.9048 | 0.9063 |

```
In [20]: ┏ from pycaret.classification import finalize_model
best_final = finalize_model(best) # re-train the dataset with whole input training data
```

In [21]:

```
from pycaret.classification import predict_model
predictions = predict_model(best_final, data = X_test) # make new predictions on new-cor
predictions
```

Out[21]:

| mean ctness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ... | worst perimeter | worst area | worst smoothness | worst compactness | worst texture | wors concavity |
|----------------|-------------------|---------------------------|------------------|------------------------------|-----|--------------------|---------------|---------------------|----------------------|------------------|-------------------|
| .14690 | 0.14450 | 0.08172 | 0.2116 | 0.07325 | ... | 113.30 | 844.4 | 0.15740 | 0.38560 | 0.51060 | ... |
| .05205 | 0.02772 | 0.02068 | 0.1619 | 0.05584 | ... | 91.29 | 632.9 | 0.12890 | 0.10630 | 0.13900 | ... |
| .05581 | 0.02087 | 0.02652 | 0.1589 | 0.05586 | ... | 96.53 | 688.9 | 0.10340 | 0.10170 | 0.06260 | ... |
| .05220 | 0.02475 | 0.01374 | 0.1635 | 0.05586 | ... | 105.80 | 819.7 | 0.09445 | 0.21670 | 0.15650 | ... |
| .03766 | 0.02562 | 0.02923 | 0.1467 | 0.05863 | ... | 84.46 | 545.9 | 0.09701 | 0.04619 | 0.04833 | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| .08511 | 0.08625 | 0.04489 | 0.1609 | 0.05871 | ... | 108.60 | 906.5 | 0.12650 | 0.19430 | 0.31690 | ... |
| .09965 | 0.03738 | 0.02098 | 0.1652 | 0.07238 | ... | 87.38 | 576.0 | 0.11420 | 0.19750 | 0.14500 | ... |
| .06797 | 0.02495 | 0.01875 | 0.1695 | 0.06556 | ... | 70.10 | 362.7 | 0.11430 | 0.08614 | 0.04158 | ... |
| .07780 | 0.04608 | 0.03528 | 0.1521 | 0.05912 | ... | 114.20 | 880.8 | 0.12200 | 0.20090 | 0.21510 | ... |
| .10850 | 0.05928 | 0.03279 | 0.1943 | 0.06612 | ... | 86.67 | 552.0 | 0.15800 | 0.17510 | 0.18890 | ... |

In [22]:

```
df_compare = pd.concat([predictions['Label'],y_test],axis = 1) # compare with the ground truth
df_compare
```

Out[22]:

| | Label | target |
|------------|-------|--------|
| 512 | 0 | 0 |
| 457 | 1 | 1 |
| 439 | 1 | 1 |
| 298 | 1 | 1 |
| 37 | 1 | 1 |
| ... | ... | ... |
| 100 | 0 | 0 |
| 336 | 1 | 1 |
| 299 | 1 | 1 |
| 347 | 1 | 1 |
| 502 | 1 | 1 |

188 rows × 2 columns

```
In [23]: └─▶ import numpy as np  
      np.mean(predictions['Label'].to_numpy() == y_test.to_numpy()) # calculate the percentage  
      #mean of the number of matches, using a boolean test on the array.
```

Out[23]: 0.9574468085106383

```
In [34]: └─▶ from pycaret.classification import create_model  
      lr = create_model('lr') # what if we only want the Logistic regression model?
```

| | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC |
|------|----------|--------|--------|--------|--------|--------|--------|
| 0 | 0.9630 | 0.9886 | 1.0000 | 0.9412 | 0.9697 | 0.9222 | 0.9250 |
| 1 | 0.9259 | 1.0000 | 1.0000 | 0.8889 | 0.9412 | 0.8421 | 0.8528 |
| 2 | 0.9630 | 0.9943 | 1.0000 | 0.9412 | 0.9697 | 0.9222 | 0.9250 |
| 3 | 0.9259 | 1.0000 | 1.0000 | 0.8889 | 0.9412 | 0.8421 | 0.8528 |
| 4 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 5 | 0.8889 | 0.9375 | 0.9375 | 0.8824 | 0.9091 | 0.7666 | 0.7689 |
| 6 | 0.9231 | 0.9750 | 0.9375 | 0.9375 | 0.9375 | 0.8375 | 0.8375 |
| 7 | 0.9615 | 1.0000 | 0.9375 | 1.0000 | 0.9677 | 0.9202 | 0.9232 |
| 8 | 0.9615 | 0.9875 | 0.9375 | 1.0000 | 0.9677 | 0.9202 | 0.9232 |
| 9 | 0.8462 | 0.9625 | 0.8750 | 0.8750 | 0.8750 | 0.6750 | 0.6750 |
| Mean | 0.9359 | 0.9845 | 0.9625 | 0.9355 | 0.9479 | 0.8648 | 0.8683 |
| SD | 0.0419 | 0.0197 | 0.0415 | 0.0484 | 0.0339 | 0.0886 | 0.0886 |

```
In [25]: ► predict_model(lr) # validation dataset -- sample exam!
```

| Model | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC |
|-----------------------|----------|--------|--------|--------|--------|--------|--------|
| 0 Logistic Regression | 0.9652 | 0.9970 | 0.9605 | 0.9865 | 0.9733 | 0.9234 | 0.9240 |

Out[25]:

| mean radius | mean texture | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ra |
|-------------|--------------|-------------|-----------------|------------------|----------------|---------------------|---------------|------------------------|-----|
| 9.173000 | 13.860000 | 260.899994 | 0.07721 | 0.08751 | 0.05988 | 0.021800 | 0.2341 | 0.06963 | 0.4 |
| 11.680000 | 16.170000 | 420.500000 | 0.11280 | 0.09263 | 0.04279 | 0.031320 | 0.1853 | 0.06401 | 0.4 |
| 16.129999 | 17.879999 | 807.200012 | 0.10400 | 0.15590 | 0.13540 | 0.077520 | 0.1998 | 0.06515 | 0.4 |
| 12.180000 | 14.080000 | 461.399994 | 0.07734 | 0.03212 | 0.01123 | 0.005051 | 0.1673 | 0.05649 | 0.4 |
| 9.667000 | 18.490000 | 289.100006 | 0.08946 | 0.06258 | 0.02948 | 0.015140 | 0.2238 | 0.06413 | 0.4 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 14.710000 | 21.590000 | 656.900024 | 0.11370 | 0.13650 | 0.12930 | 0.081230 | 0.2027 | 0.06758 | 0.4 |
| 14.400000 | 26.990000 | 646.099976 | 0.06995 | 0.05223 | 0.03476 | 0.017370 | 0.1707 | 0.05433 | 0.4 |
| 20.290001 | 14.340000 | 1297.000000 | 0.10030 | 0.13280 | 0.19800 | 0.104300 | 0.1809 | 0.05883 | 0.4 |
| 10.290000 | 27.610001 | 321.399994 | 0.09030 | 0.07658 | 0.05999 | 0.027380 | 0.1593 | 0.06127 | 0.4 |
| 12.580000 | 18.400000 | 489.000000 | 0.08393 | 0.04216 | 0.00186 | 0.002924 | 0.1697 | 0.05855 | 0.4 |

ws × 32 columns

```
In [35]: ► final_lr = finalize_model(lr)
```

```
In [36]: ► predictions_lr = predict_model(final_lr, data = X_test)
np.mean(predictions_lr['Label'].to_numpy() == y_test.to_numpy())
```

Out[36]: 0.9627659574468085

In [37]: ┶

```
from pycaret.classification import tune_model
tuned_lr = tune_model(lr) # fine-tuning the parameters in Logistic regression
```

| | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC |
|------|----------|--------|--------|--------|--------|--------|--------|
| 0 | 0.9630 | 0.9830 | 1.0000 | 0.9412 | 0.9697 | 0.9222 | 0.9250 |
| 1 | 0.9630 | 1.0000 | 1.0000 | 0.9412 | 0.9697 | 0.9222 | 0.9250 |
| 2 | 0.9630 | 0.9886 | 1.0000 | 0.9412 | 0.9697 | 0.9222 | 0.9250 |
| 3 | 0.9259 | 1.0000 | 1.0000 | 0.8889 | 0.9412 | 0.8421 | 0.8528 |
| 4 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 5 | 0.8889 | 0.9318 | 0.9375 | 0.8824 | 0.9091 | 0.7666 | 0.7689 |
| 6 | 0.9231 | 0.9750 | 0.9375 | 0.9375 | 0.9375 | 0.8375 | 0.8375 |
| 7 | 0.9615 | 1.0000 | 0.9375 | 1.0000 | 0.9677 | 0.9202 | 0.9232 |
| 8 | 0.9231 | 0.9875 | 0.9375 | 0.9375 | 0.9375 | 0.8375 | 0.8375 |
| 9 | 0.8846 | 0.9625 | 0.8750 | 0.9333 | 0.9032 | 0.7607 | 0.7632 |
| Mean | 0.9396 | 0.9828 | 0.9625 | 0.9403 | 0.9505 | 0.8731 | 0.8758 |
| SD | 0.0348 | 0.0208 | 0.0415 | 0.0363 | 0.0287 | 0.0728 | 0.0725 |

In [38]: ┶

```
predict_model(tuned_lr) # still doing the sample exam -- validation dataset
```

| | Model | Accuracy | AUC | Recall | Prec. | F1 | Kappa | MCC |
|---|---------------------|----------|--------|--------|--------|--------|--------|--------|
| 0 | Logistic Regression | 0.9478 | 0.9956 | 0.9342 | 0.9861 | 0.9595 | 0.8864 | 0.8890 |

Out[38]:

| | mean radius | mean texture | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimens |
|-----|-------------|--------------|-------------|-----------------|------------------|----------------|---------------------|---------------|---------------------|
| 0 | 9.173000 | 13.860000 | 260.899994 | 0.07721 | 0.08751 | 0.05988 | 0.021800 | 0.2341 | 0.06 |
| 1 | 11.680000 | 16.170000 | 420.500000 | 0.11280 | 0.09263 | 0.04279 | 0.031320 | 0.1853 | 0.06 |
| 2 | 16.129999 | 17.879999 | 807.200012 | 0.10400 | 0.15590 | 0.13540 | 0.077520 | 0.1998 | 0.06 |
| 3 | 12.180000 | 14.080000 | 461.399994 | 0.07734 | 0.03212 | 0.01123 | 0.005051 | 0.1673 | 0.05 |
| 4 | 9.667000 | 18.490000 | 289.100006 | 0.08946 | 0.06258 | 0.02948 | 0.015140 | 0.2238 | 0.06 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 110 | 14.710000 | 21.590000 | 656.900024 | 0.11370 | 0.13650 | 0.12930 | 0.081230 | 0.2027 | 0.06 |
| 111 | 14.400000 | 26.990000 | 646.099976 | 0.06995 | 0.05223 | 0.03476 | 0.017370 | 0.1707 | 0.05 |
| 112 | 20.290001 | 14.340000 | 1297.000000 | 0.10030 | 0.13280 | 0.19800 | 0.104300 | 0.1809 | 0.05 |
| 113 | 10.290000 | 27.610001 | 321.399994 | 0.09030 | 0.07658 | 0.05999 | 0.027380 | 0.1593 | 0.06 |
| 114 | 12.580000 | 18.400000 | 489.000000 | 0.08393 | 0.04216 | 0.00186 | 0.002924 | 0.1697 | 0.05 |

115 rows × 32 columns

```
In [39]: final_tuned_lr = finalize_model(tuned_lr) #retrain with the whole dataset
```

```
In [40]: predictions_tuned_lr = predict_model(final_tuned_lr, data = X_test)  
np.mean(predictions_tuned_lr['Label'].to_numpy() == y_test.to_numpy())
```

```
Out[40]: 0.9574468085106383
```

Let's recap the workflow above (or about general supervised learning)

- The **minimum requirement** is that we have a training dataset with both X and y (also called labels, targets...). We want to **fit the mapping** between x and y with **training dataset** (the process is indeed called training), and making predictions about the new y given new X in the test dataset.
 - *Remark 1:* The true y in test dataset sometimes can also be known, so that we can know the performance the model immediately. But in general, we won't expect this.
 - *Remark 2:* In our course, just to mimic a real-world situation, sometimes we manually create (split) the train or test data.
- (Optional) We may train multiple models or one model with multiple parameters. How can we compare them and gain more confidence about the final test? Sometimes we further split the training dataset into (real) training dataset and **validation dataset** (imagine it as the sample exam), so that we can get instant feedback because we know the true label in validation dataset.
- (Optional) During training, to be more cautious, sometimes we even make more "quizzes" -- that is called **cross-validation** (will talk about the details in the next lecture)
- (Optional) With 10 "quizzes" (10-fold cross-validation) and "one sample exam" (validation data), for instance, we finally pick up the best candidate model. Before applying to the real test dataset, we don't want to waste any sample. Therefore we **finalize** training by picking up the winner model, while updating it with all the samples (including the validation data) in the training dataset.
- Finally, applying the model to test data -- wait and see!

Of course, as a math course, we are not satisfied with merely calling functions in pycaret. In the rest of lectures this quarter, we are going to dig into details of some algorithms and learn more underlying math -- turn the black box of ML into white (at least gray) one!

Unsupervised Learning

It is still challenging to give a general and rigorous definition for unsupervised learning mathematically. Let's focus on more specific tasks.

- Dimension Reduction

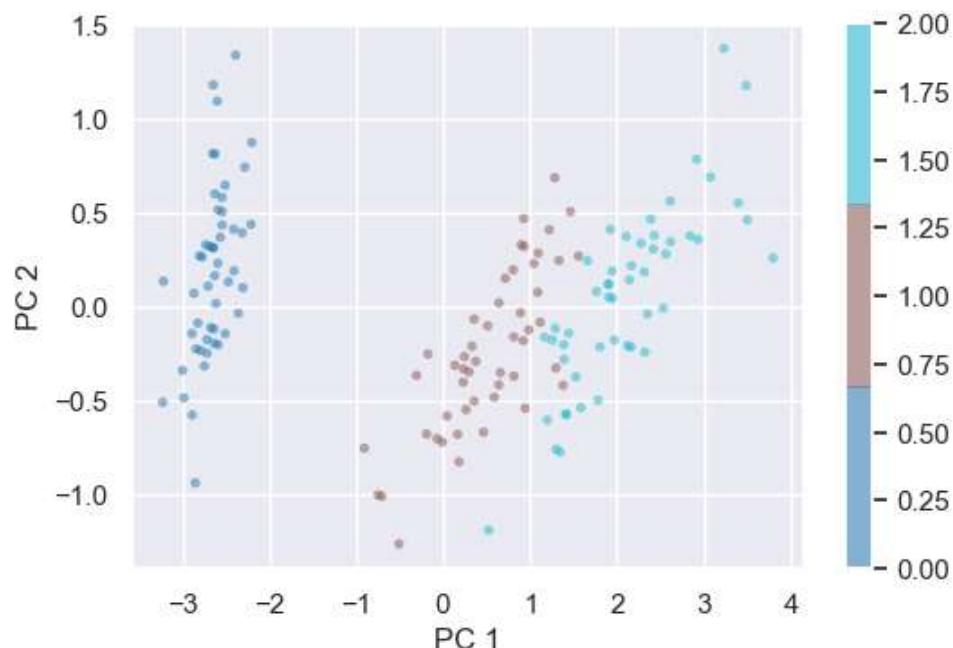
Given $X \in \mathbb{R}^{n \times p}$, finding a mapping function $\mathbf{f} : \mathbb{R}^p \rightarrow \mathbb{R}^q (q \ll p)$ such that the low-dimensional coordinates $z^{(i)} = \mathbf{f}(x^{(i)})$ "preserve the information" about $x^{(i)}$.

- Question: Difference with supervised learning?
- Linear mapping: Principle Component Analysis (PCA)
- Nonlinear mapping: Manifold Learning, Autoencoder

```
In [41]: └─▶ from sklearn.datasets import load_iris  
      X,y = load_iris(return_X_y = True) # Note that in the hw this week, it's not allowed to  
      X
```

```
In [42]: ┆ from sklearn.decomposition import PCA  
pca = PCA(n_components=2) # principle component analysis, reduce 4-dimenional data to 2-  
X_pca = pca.fit_transform(X)  
X_pca
```

```
In [43]: # import matplotlib.pyplot as plt
# import seaborn as sns
sns.set() # set the seaborn theme style
figure = plt.figure(dpi=100)
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, s=15, edgecolor='none', alpha=0.5,cmap=plt.cm
plt.xlabel('PC 1')
plt.ylabel('PC 2')
plt.colorbar();
```



- Clustering

Given $X \in \mathbb{R}^{n \times p}$, finding a partition of the dataset into K groups such that

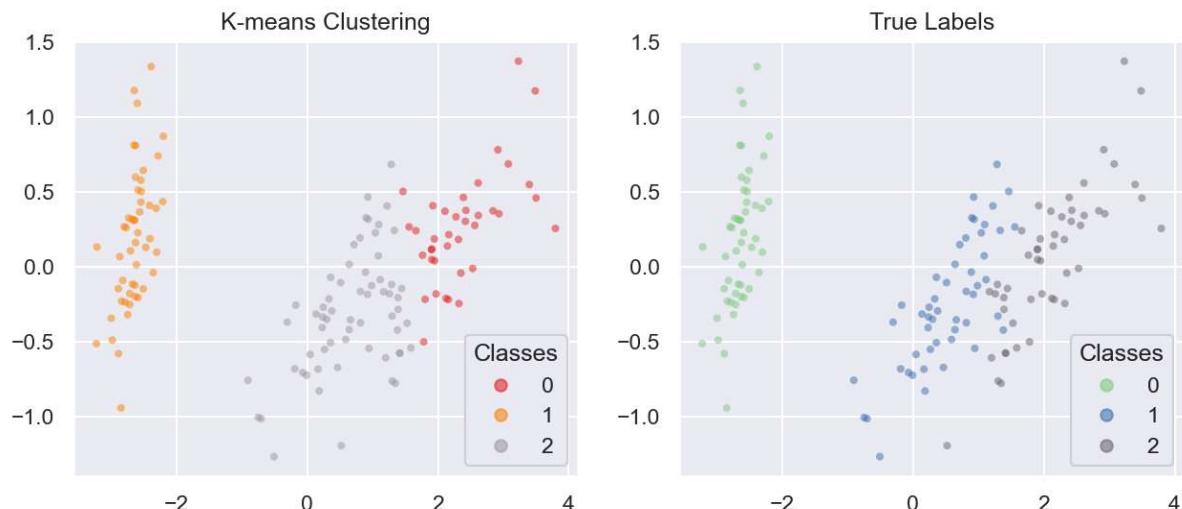
- data within the same group are similar;
- data from different groups are dissimilar.

```
In [44]: ┆ from sklearn.cluster import KMeans  
kmeans = KMeans(n_clusters=3, random_state=0) #call k-means clustering algorithm  
y_km = kmeans.fit_predict(X)  
y_km # the groups assigned by algorithm
```

```
In [45]: import matplotlib.pyplot as plt
import seaborn as sns; sns.set()
fig, (ax1, ax2) = plt.subplots(1, 2,dpi=150, figsize=(10,4))

fig1 = ax1.scatter(X_pca[:, 0], X_pca[:, 1],c=y_km, s=15, edgecolor='none', alpha=0.5,cr
fig2 = ax2.scatter(X_pca[:, 0], X_pca[:, 1],c=y, s=15, edgecolor='none', alpha=0.5,cmap=
ax1.set_title('K-means Clustering')
legend1 = ax1.legend(*fig1.legend_elements(), loc="best", title="Classes")
ax1.add_artist(legend1)
ax2.set_title('True Labels')
legend2 = ax2.legend(*fig2.legend_elements(), loc="best", title="Classes")
ax2.add_artist(legend2)
```

Out[45]: <matplotlib.legend.Legend at 0x16bfed19d90>



Question: What is the difference between clustering and classification? Can you try classification on Iris data with pycaret right now?

```
In [ ]: # try classification with pycaret for Iris data by yourself!
```