

ICDS Spring 2025

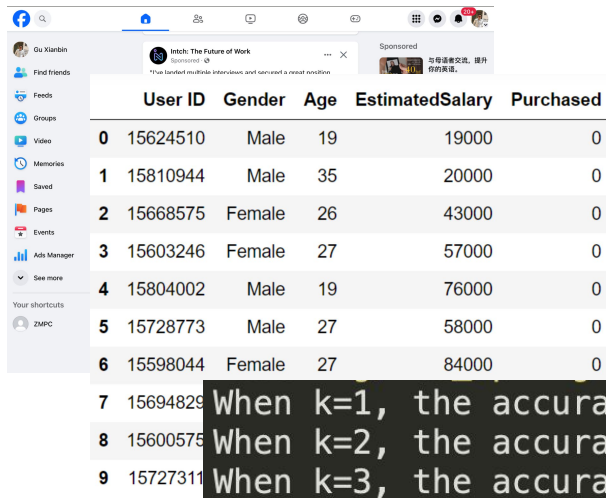
Data Science: Part 2

Supervised and Unsupervised Learning

Recap: Machine Learning

- Definition: **PET!**

- A computer program that learns from **experience E** with respect to some class of **task T** and performance **measure P**



The screenshot shows a Facebook interface with a table of user data. The table has columns: User ID, Gender, Age, EstimatedSalary, and Purchased. The data is as follows:

	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	19	19000	0
1	15810944	Male	35	20000	0
2	15668575	Female	26	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	19	76000	0
5	15728773	Male	27	58000	0
6	15598044	Female	27	84000	0
7	15694829				
8	15600575				
9	15727311				

Below the table, there is a black box with white text containing the following information:

- When $k=1$, the accuracy is 0.68.
- When $k=2$, the accuracy is 0.68.
- When $k=3$, the accuracy is 0.62.

Task: predicting whether a user will make purchase or not

Experience: the observations/dataset on how existing users responding to the ads

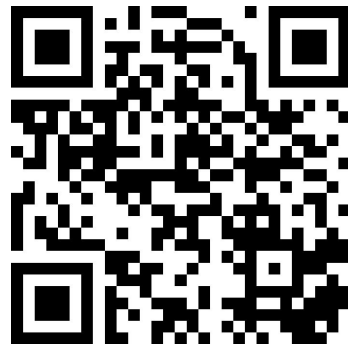
Performance: prediction accuracy

Exercise 1: Types of ML Tasks

We used KNN algorithm in the previous example. What type of machine learning task did it perform?

- A. Clustering
- B. Classification
- C. Regression
- D. Dimension reduction

Scan to answer!



Exercise 2: KNN Algorithm

Scan to answer!



Which of the problem below can be solved by KNN?

- A. Partitioning Genomic Data: Given a dataset containing gene expression profiles, partition genes into clusters based on their expression patterns to uncover underlying biological processes.
- B. Movie Recommendation: Given a dataset of user ratings for different movies, recommend movies to a new user based on the ratings of similar users.
- C. Segmenting Satellite Images: Given a collection of satellite images, group similar regions together based on pixel intensities to identify different land cover types.
- D. Travel Time Estimation: Given a set of records of traffics, weather conditions, and the time spent traveling from your apartment to the campus, predict tomorrow's commuting time.

Recap: Machine Learning Models

- 3 fundamental components of a ML model:
 - Model design (learning/inductive bias)
 - Assumptions to generalize to unseen data samples
 - Principle for choosing the model to use
 - Data representation (features)
 - Feature engineering: deriving task-appropriate features
 - Measurement (distance metrics)
 - Definition for similarity: Euclidean, non-euclidean, etc.

Exercise 3: Feature Engineering

To build a predictor for travel time estimation, what feature can also be helpful besides historical traffics and weathers?

- A. Time: time-of-day, day-of-week, whether-holiday
- B. Road attributes (e.g., speed limit, #lanes) in the route
- C. Historical records of traffic accident occurrences
- D. All of above

Scan to answer!



Agenda

- Regression and gradient descent
 - Regression analysis
 - Gradient descent
- Supervised learning
- Clustering and K-Means
- Unsupervised learning
- Appendix: using the sklearn's k-means

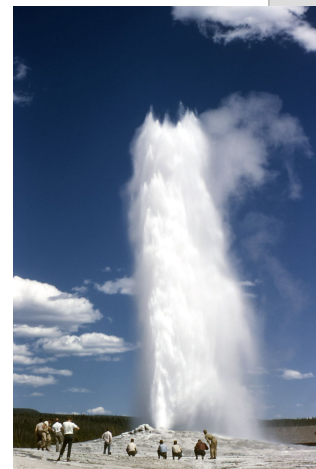
Linear Regression and Gradient Descent

Predicting the waiting time

The old faithful geyser dataset

- Waiting time between eruptions and the duration of the eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA.
- 272 observations on 2 variables

Task: given an eruption duration, predict the waiting time.

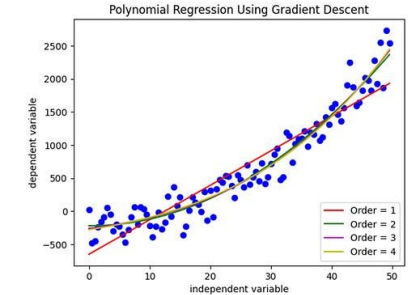
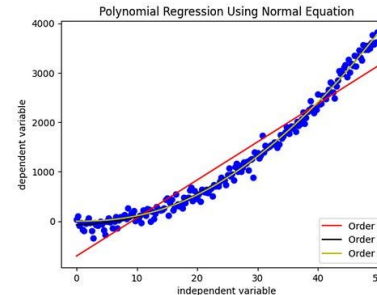
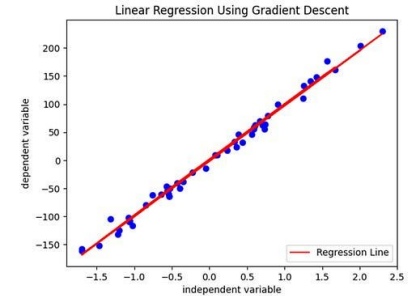
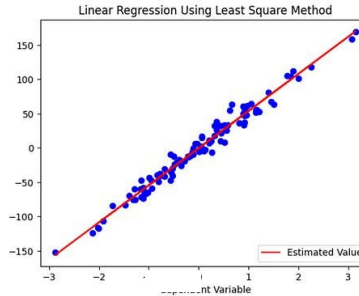


faithful		
index	eruptions	waiting
1	3.6	79
2	1.8	54
3	3.333	74
4	2.283	62
5	4.533	85
6	2.883	55
7	4.7	88
8	3.6	85
9	1.95	51
10	4.35	85

Regression Analysis

- It is a statistical method for estimating the relationship between a dependent variable and one or more independent variables.
- A tool for making predictions, understanding relationships between variables, and making decisions
- Widely used in economics, finance, engineering, and social science

Regression Analysis



Regression model

A general form of a regression model,

$$Y_i = f(X_i, \beta) + e_i \quad (1)$$

where,

- Y_i is the dependent variable;
 - X_i is the independent variable;
 - β is the model parameter(s);
 - e_i is the error terms, representing the observation errors.
- It is a mathematical equation that represents the relationship between Y and X.
 - β represents the change of Y for a one-unit change of X.
 - e represents the error term (also called, residuals), representing the variability in Y that cannot explained by X.

Linear regression

If we assume the form of f is linear combination of the parameters β , we have the linear regression. For example, let $f(X_i, \beta) = \beta_0 + \beta_1 X_i$, we have

$$Y_i = \beta_0 + \beta_1 X_i + e_i \quad (1)$$

where

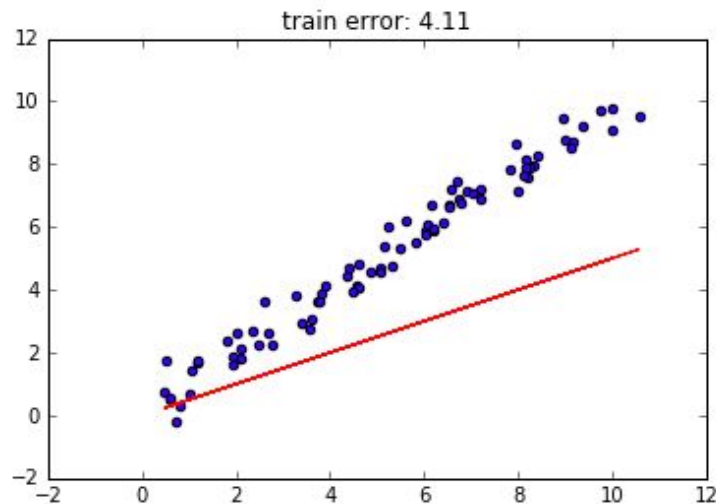
- β_0 is the intercept, the value of Y when $X = 0$;
 - β_1 is the slope;
 - e_i is the error item; we assume $e_i \sim N(0, 1)$.
- Linear regression is the most common model used in regression analysis; it assumes the relationship between Y and X can be represented by a straight line.
 - In practice, when we have no knowledge about what the relationship between Y and X , we can firstly try the linear regression.

Find the parameters of a regression model


We want the best fit line (making as less error as possible), i.e., the line covers as many points as possible.

In the old faithful geyser example:

- Our model: $\hat{y} = \mathbf{w}^T \mathbf{x}$
- Goal: find a \mathbf{w} that minimizes the prediction error, $(\hat{y} - y)$.

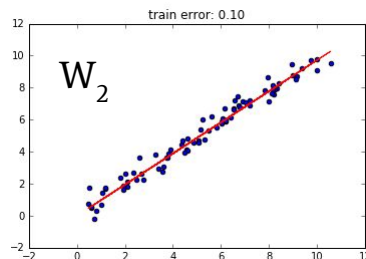
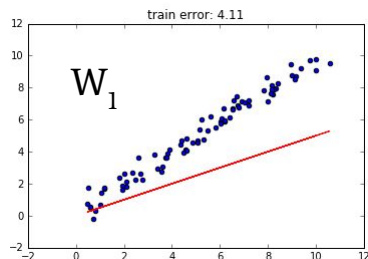


Loss function for the task

- Model: $y' = wx$
- Dependent variable: y (i.e., the observed true value)
- Loss function: $E(w) = \frac{1}{2n} \sum_{i=1}^n (\boxed{wx_i} - y_i)^2$
 i-th prediction, y_i'
- Task: find a proper w such that E is minimized

Solving the minimization task

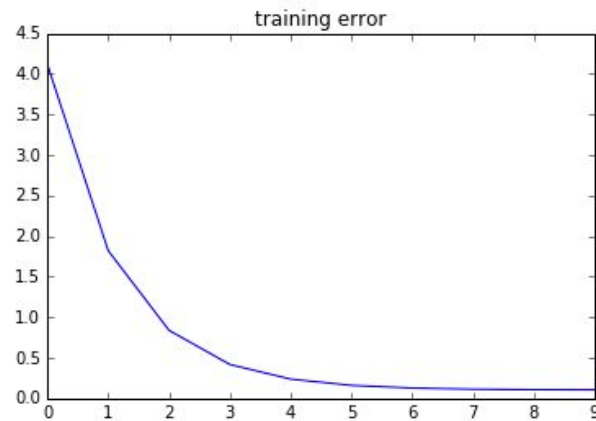
- Adjust w such that E decreases monotonically



What method can achieve this process?




- Gradient descent

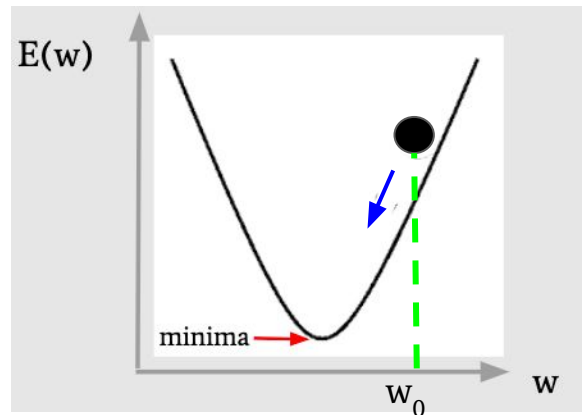
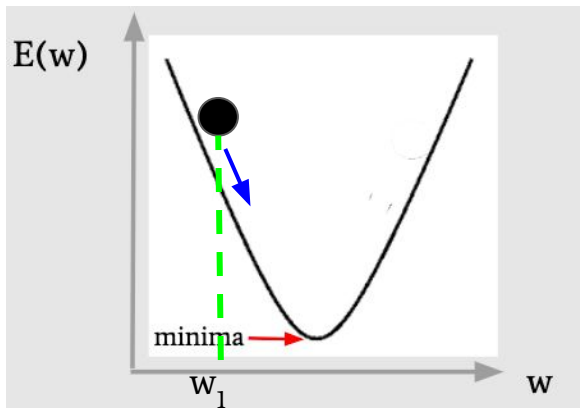


Gradient descent

Gradient: $\frac{dE}{dw}$, is the ratio of changes in E in response to the changes in w .

We have $\Delta E = \frac{dE}{dw} \Delta w$  It represents how much E changes when w changes

Can we find a Δw that always makes $\Delta E < 0$?



Gradient descent

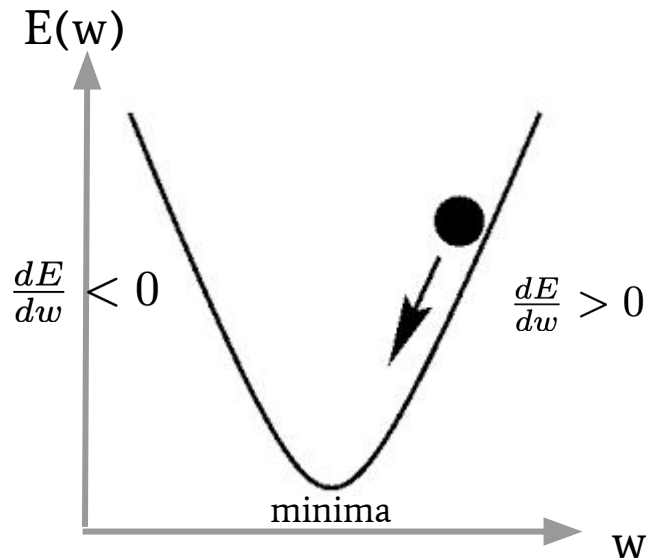
It is always opposite to the gradient so guarantees $\Delta E \leq 0$

$$\Delta w = -\lambda \frac{dE}{dw}$$

Learning rate, a hyperparameter

$$\Delta E = \frac{dE}{dw} \Delta w$$

ΔE is always non-positive.



A little bit calculus

Q: What is the gradient of our loss function?

A: The first-Order derivative.

$$E(w) = \frac{1}{2n} \sum_{i=1}^n (wx_i - y_i)^2$$

$$\frac{dE}{dw} = ?$$

The gradient of the loss function

$$E(w) = \frac{1}{2n} \sum_{i=1}^n (wx_i - y_i)^2$$

$$\frac{dE}{dw} = \frac{1}{n} \sum_{i=1}^n (wx_i - y_i)x_i$$

$$\frac{dE}{dw} = \frac{1}{n} \sum_{i=1}^n \underline{(y'_i - y_i)x_i}$$

- Averaging the prediction errors that are weighted by input contribution (i.e., x_i)

Solving the minimization task using gradient descent:

1. Initial w
2. Calculate y' by $y'_i = wx_i$ for $i = 1, 2, 3, \dots, n$
3. Calculate gradient $\frac{dE}{dw}$ by $\frac{dE}{dw} = \frac{1}{n} \sum_{i=1}^n (y'_i - y_i)x_i$;
4. Increase w by $\Delta w = -\lambda \frac{dE}{dw}$, then, goto step 2.
 - Pick a random **w** to start, update **w** until it converges.
 - We would like to record the E in each round.
 - The algorithm stops when E decreases very little (i.e., $|\Delta E| < \text{threshold}$; converge).

Regression with gradient descent

```
80 # training loop
81 while True:
82
83     # predict and compute error
84     pred_y = [w*x for x in d_x]
85     # print(pred_y)
86     error = compute_error(pred_y, d_y)
87     error_history.append(error)
88     #compute the change of error
89     delta_error = error - error_last_iter
90     error_last_iter = error
91
92     plt.cla()
93     plt.scatter(d_x, d_y)
94     plt.plot(d_x, pred_y, 'r')
95     title_str = "train error: %0.2f" % error
96     plt.title(title_str)
97     plt.show()
98
99     steps += 1
100     if abs(delta_error) <= margin or steps >= num_iter:
101         break
102
103     # gradient descent
104     grad = compute_grad(pred_y, d_y, d_x)
105     w -= learning_rate * grad
```

Computing predictions

Recording errors

Checking the convergence

Computing gradients and
updating w

Regression with gradient descent

Complete `compute_error` and `compute_grad`:

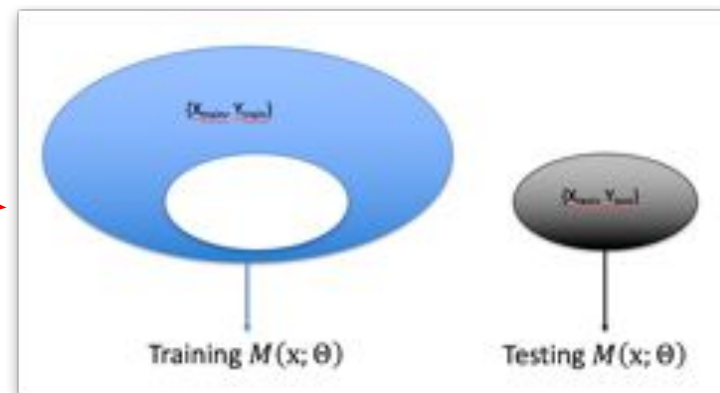
$$E(w) = \frac{1}{2n} \sum_{i=1}^n (wx_i - y_i)^2$$

$$\frac{dE}{dw} = \frac{1}{n} \sum_{i=1}^n (pred_i - y_i)x_i$$

```
24 def compute_error(pred, truth):
25     """ mean square error:
26         sum of (pred[i] - truth[i])^2, divided by 2*length
27     """
28     return regression_helper.compute_error(pred, truth)
29
30 def compute_grad(pred, truth, x):
31     """ gradient is mean of (pred[i] - truth[i]) * x[i]
32     """
33     return regression_helper.compute_grad(pred, truth, x)
```

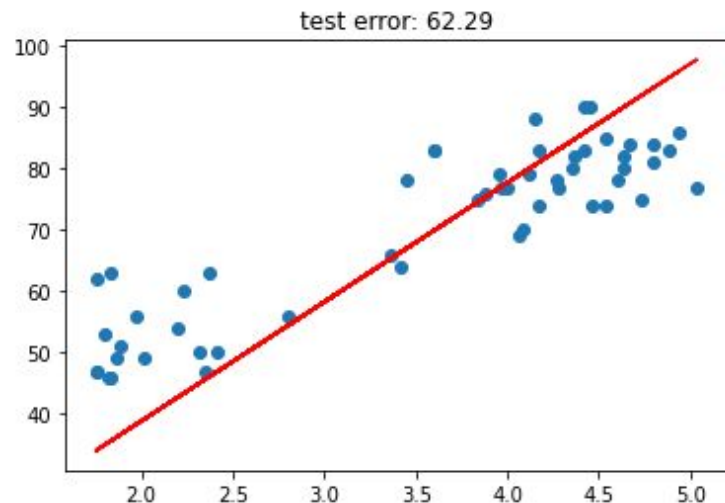
Data processing

```
31 if __name__ == "__main__":
32
33     path = './faithful.csv'
34     # loading data
35     f = open(path, 'r')
36     raw_data = f.readlines()
37     data = []
38     for item in raw_data[1:]:
39         item = item.strip().split(',')
40         data.append([float(item[1]), float(item[2])])
41     print(f"There are {len(data)} data items.")
42     # Make a random split of the data
43     random.shuffle(data)
44     train_data = data[:round(0.8*len(data))]
45     test_data = data[round(0.8*len(data)):]
46
47     # Plotting the randomly generated data
48     d_x = [d[0] for d in train_data]
49     d_y = [d[1] for d in train_data]
50     # plt.ion()
51     plt.figure(1)
52     plt.scatter(d_x, d_y)
53     plt.show()
```




The model never looks at test data

```
87
88     # testing on test data
89     d_x = [d.getFeatures()[0] for d in test_data]
90     d_y = [d.getFeatures()[1] for d in test_data]
91     pred_y = [w*x for x in d_x]
92     error = compute_error(pred_y, d_y)
93
94     plt.scatter(d_x, d_y)
95     plt.plot(d_x, pred_y, 'r')
96     title_str = "test error: %0.2f" % error
97     plt.title(title_str)
98     plt.show()
99
100    plt.plot(error_history)
101    plt.title('training error')
102    plt.show()
```



Settings of hyperparameters

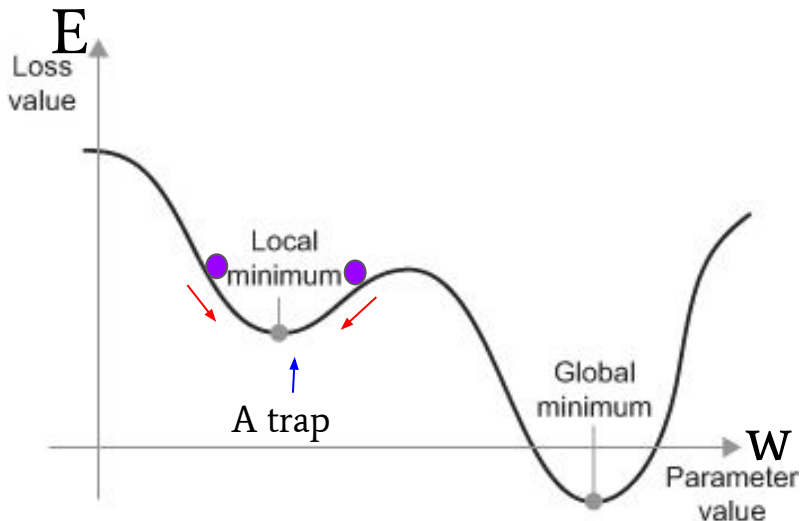
```
69     num_iter = 100
70     steps = 0
71     learning_rate = 0.01
72     n = len(train_data)
73     error_history = []
74     error_last_iter = 0
75     margin = 0.01
76
77     # w is the parameter to be learned
78     w = 0.5
```



The learning rate is a hyperparameter in gradient descent. It has large influence on the learning process.

- An improper learning rate may ruin the learning process. (You may change it to 1, and run the program to see its impact.)
- In practice, we often normalize the raw data before using them to train models which can help us find a suitable learning rate.

Traps of local minimum



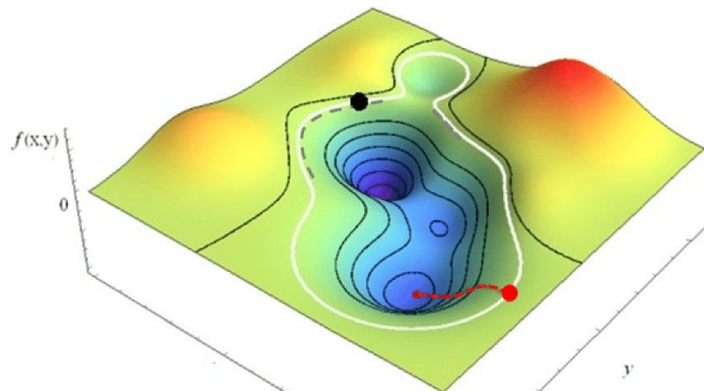
$$\Delta W = -\lambda \frac{dE}{dw}$$

- The art of setting λ
 - increasing λ , roll across the local minima.
 - An oversize λ results in non convergence (i.e., jump over the global minima)

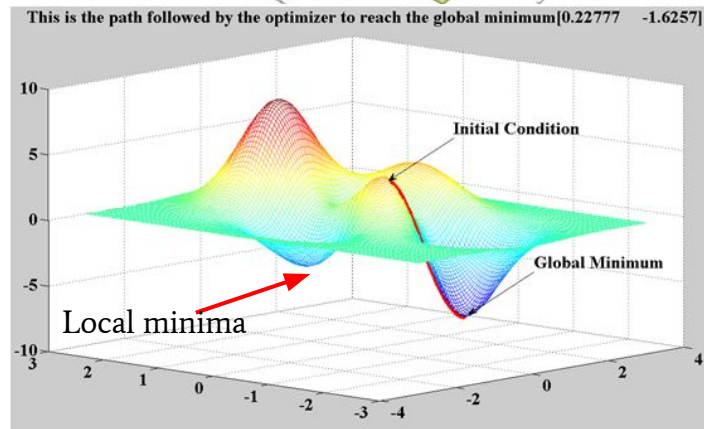
Model's performance is influenced by the hyperparameters. For **each** model, we need to tune them for fitting the dataset.

- How can we find the best hyperparameters? We may use grid search, but no perfect solution yet.
- This is why some researchers alleged machine learning is alchemy.

Solving optimization problems



- Gradient descent is widely used in solving optimization problems.
- When there are multiple parameters in the loss function, the surfaces will be intricate.
- Gradient descent cannot promise to find the global optima.



Agenda

- Regression and gradient descent
 - Regression analysis
 - Gradient descent
- **Supervised learning**
- Clustering and K-Means
- Unsupervised learning
- Appendix: using the sklearn's k-means

Supervised learning

- the most common approach in ML
 - “right answers” are given
 - models learn relations between inputs and answers
 - task: make predictions (label/value)

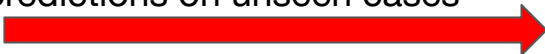


We always need to test the model (validation)

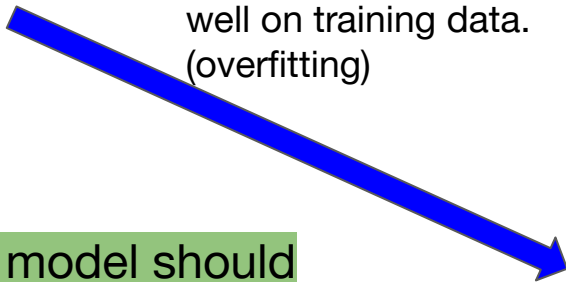
It is important to know how many correct/incorrect predictions from the model.



Good learner makes correct predictions on unseen cases



Bad learner only works well on training data.
(overfitting)

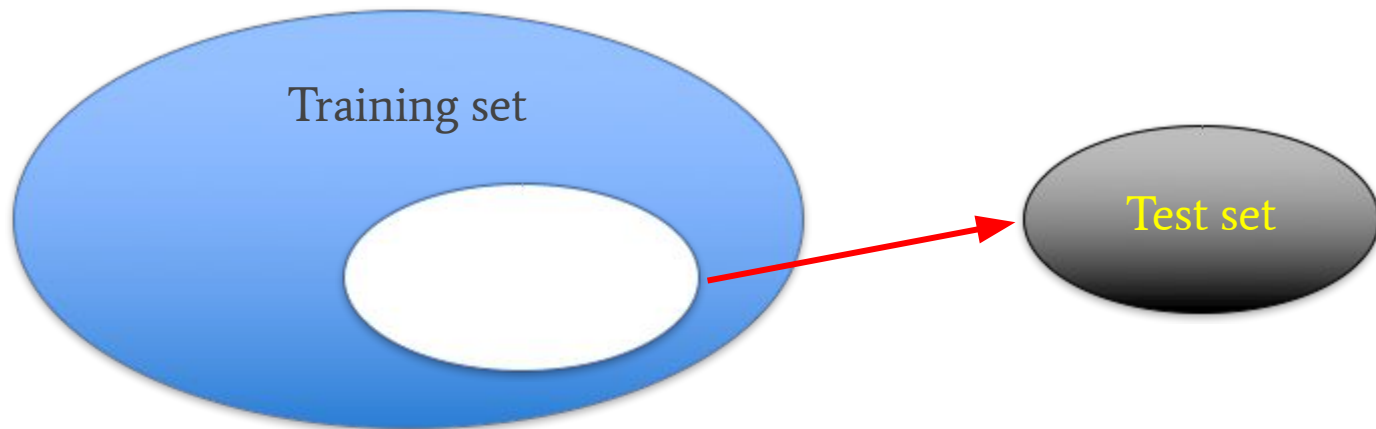


- Generalizability: a good model should work well not only on the train data but also the unseen data (i.e., testing data)



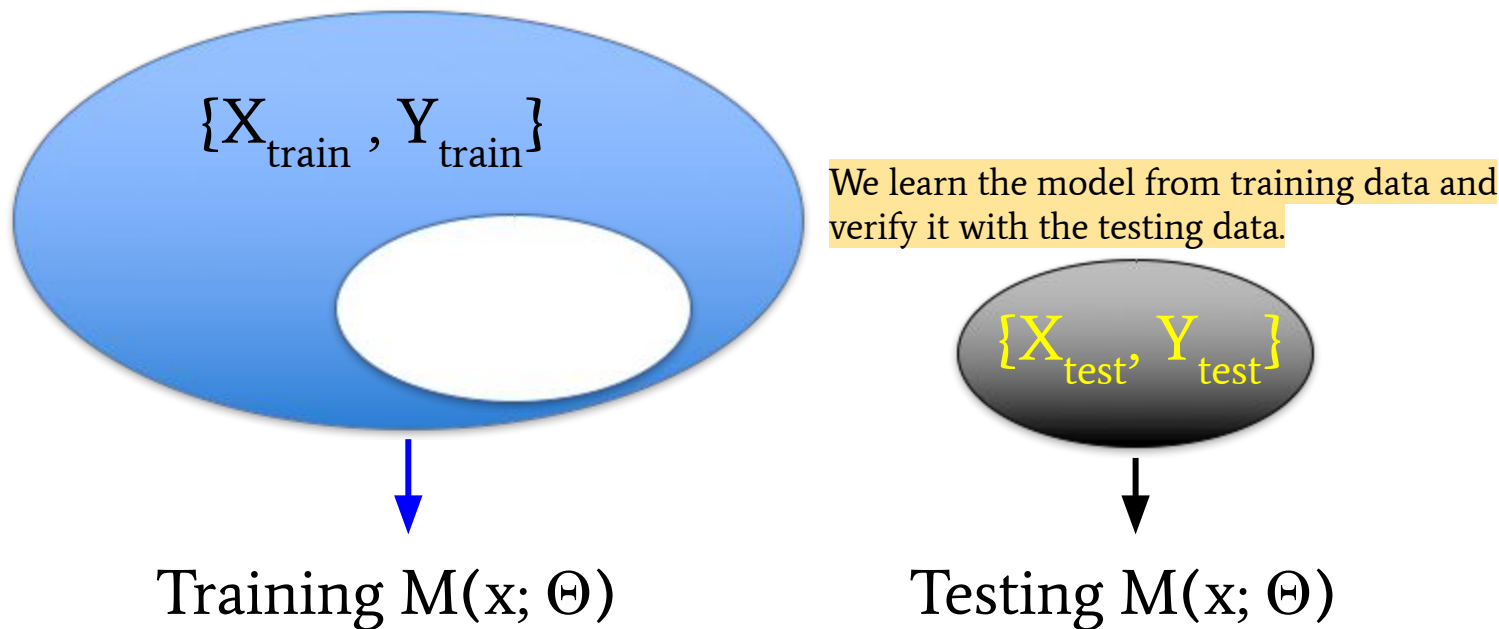
How to test?

- Splitting the samples into two parts:
 - One for training (inferring model parameters)
 - One for testing (verifying)



Why this makes sense?

Inductive bias: “the unseen samples are similar to the given samples”.



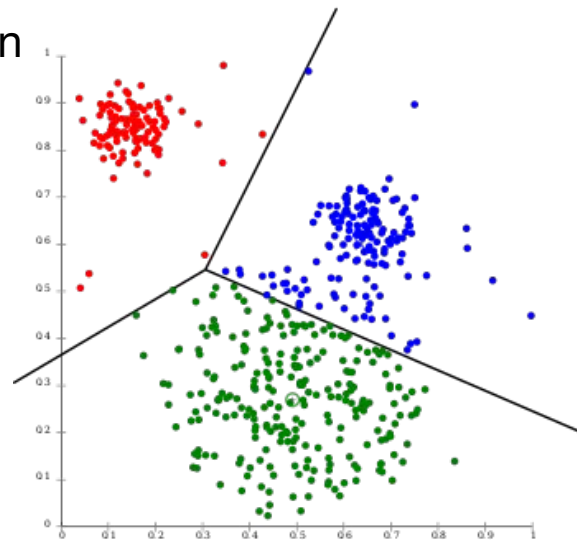
Quantify the Performance

- Main idea: counting the correct predictions on the test dataset
- Statistics for performance measurement:
 - Confusion matrix/precision-recall
 - [Confusion matrix - Wikipedia](#)
 - [Precision and recall](#)
 - [Accuracy and precision](#)

Clustering analysis

Clustering analysis

- It is a type of unsupervised learning technique for grouping similar data points together based on certain features.
- **Task:** to divide a dataset into clusters, where data points within the same cluster are more similar to each other than those in other clusters.
 - It helps in discovering hidden structures within datasets.
 - Applications: data mining, pattern recognition, image analysis, market segmentation, and etc.



Goals of Clustering analysis

The approach:

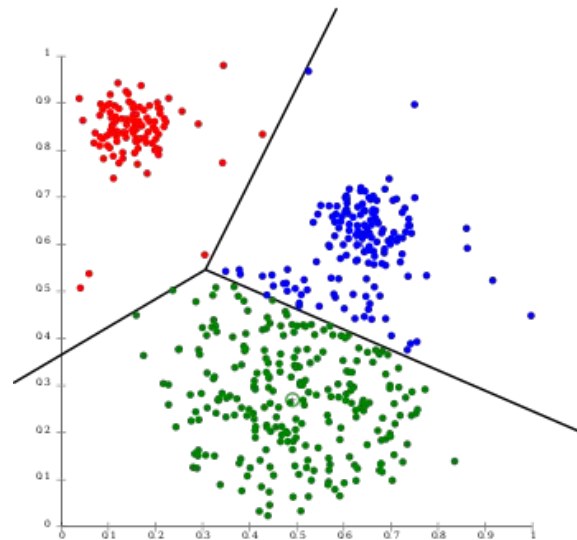
- “Similar” members go into the same group

Key metrics:

- Intra-cluster variance (difference within a group)
- Inter-cluster variance (difference across groups)

Common objectives:

- Minimize intra-cluster variance
- Maximize inter-cluster variance



The k-means algorithm

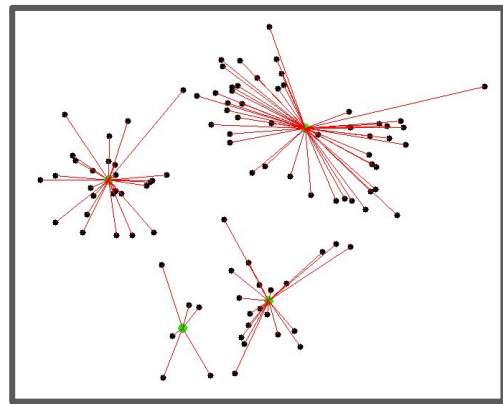
How to achieve the following objectives?

- Minimize intra-cluster variance
- Maximize inter-cluster variance

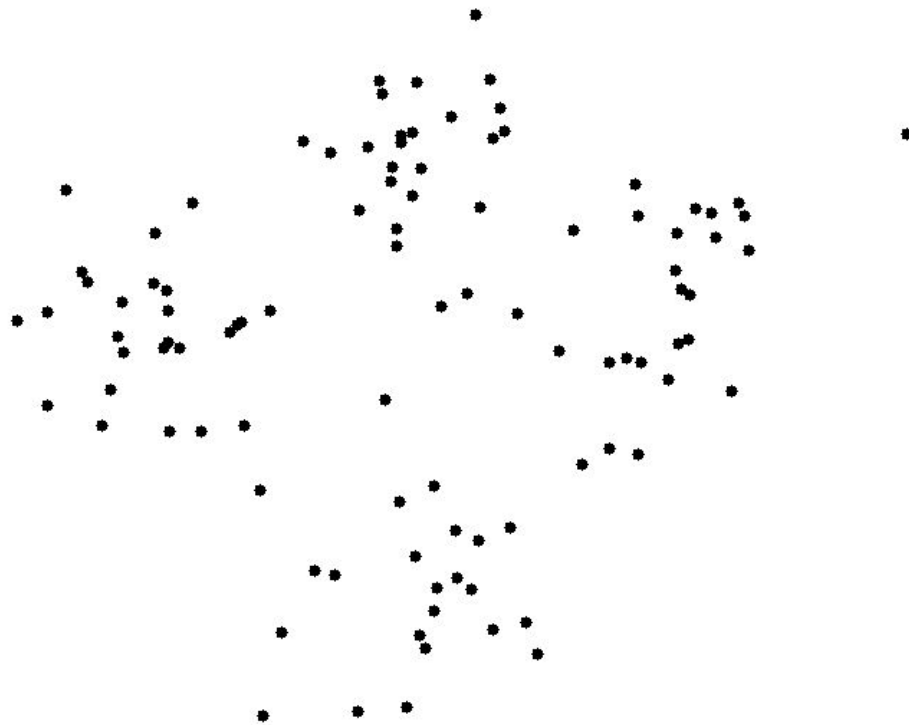
Assume that a dataset has k clusters; each cluster has a cluster center.

- If we put every data item to a cluster whose centroid is the closest one to it, we will obtain a clustering that meets the objectives.
- k-means is the algorithm to realize the above operations.

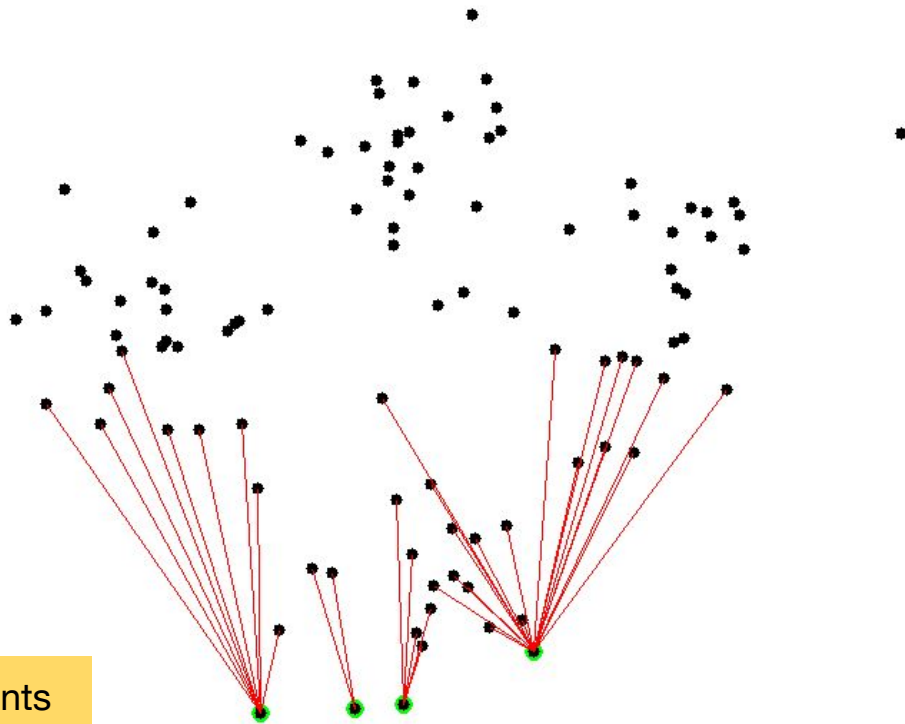
Demo: [K-means clustering, starting with 4 points](#)



k-means: given a dataset

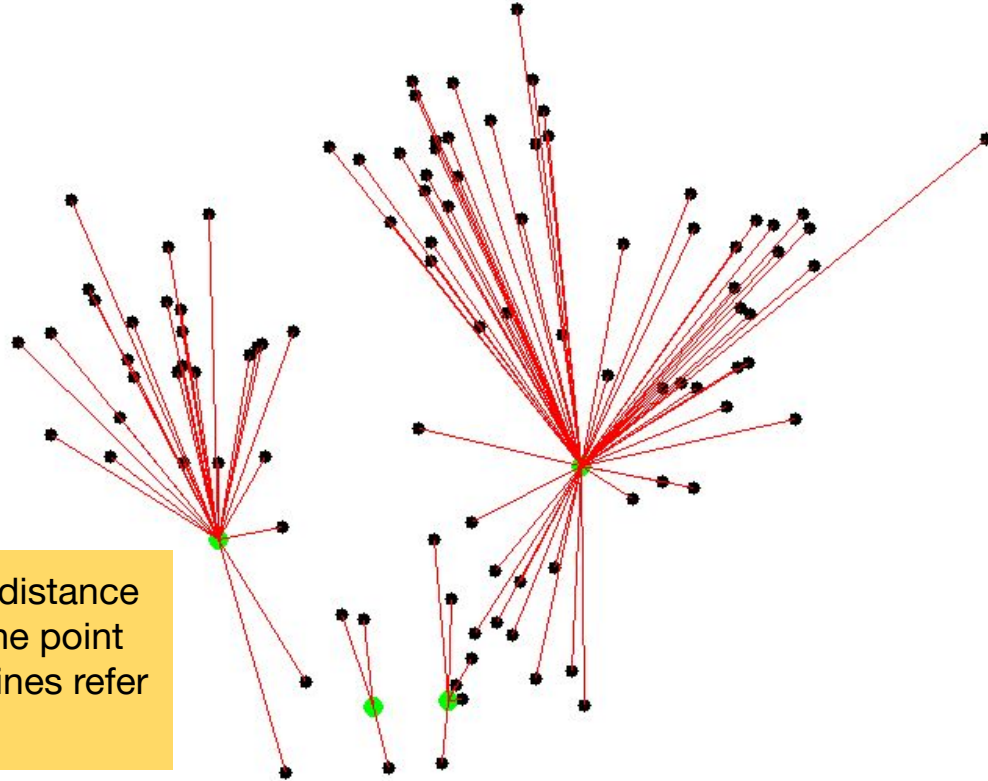


k-means: generating k centroids



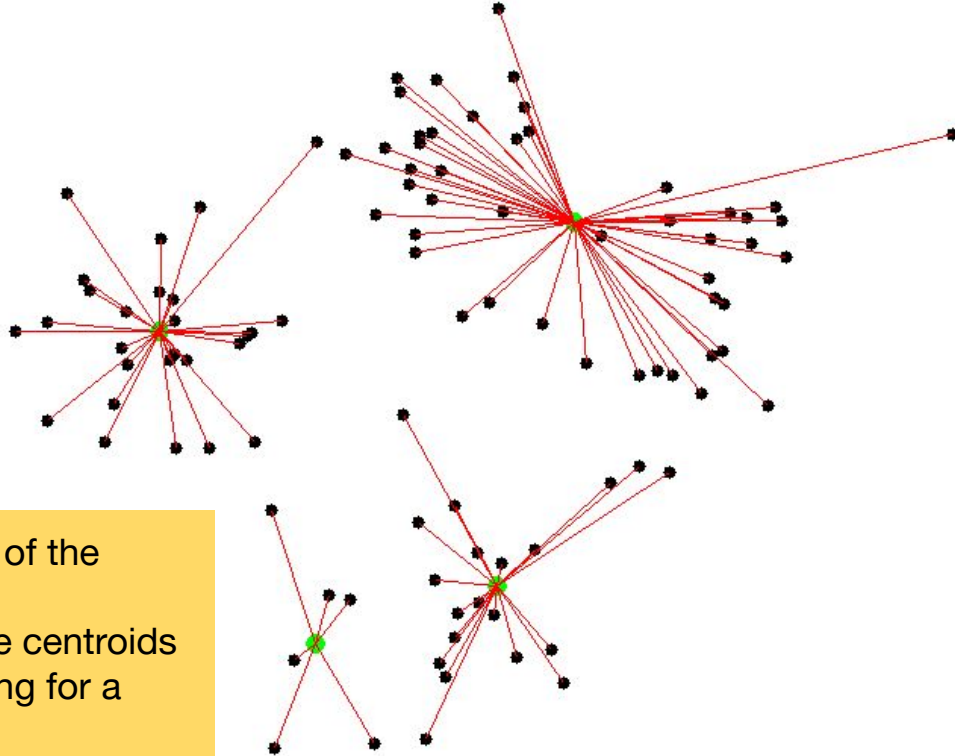
We can randomly pick k points as the cluster centroids.

k-means: assigning points to a centroid



For each point, compute the distance to each centroid; associate the point to the nearest centroid. (red lines refer to the distances)

k-means: updating centroids



Update the centers to the average of the points belonging it.

- We stop updating until **all** the centroids don't move (or, doing updating for a number of times)

k-means: the algorithm (Lloyd's algorithm)

Initialization:

- **k** empty clusters, each centroid is at one of **k** randomly picked members

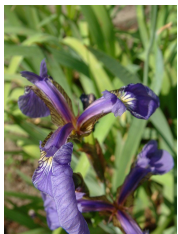
Iterate until converge (i.e. no cluster moves its centroid).

1. Each sample finds its nearest cluster, add itself to the cluster
2. Each cluster computes its new centroid
3. **If** no cluster changes its centroid – done! (condition)
4. Otherwise, all clusters clear their samples, **go to** 1. (repetition)

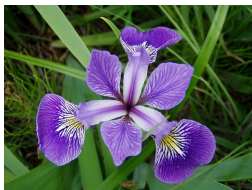
Let's play!

The Iris Dataset: [Iris flower data set - Wikipedia](#)

- 150 samples of three species of Iris (setosa, virginica, and versicolor)
- Each sample has four features: length and width of the sepals and petals respectively.



setosa



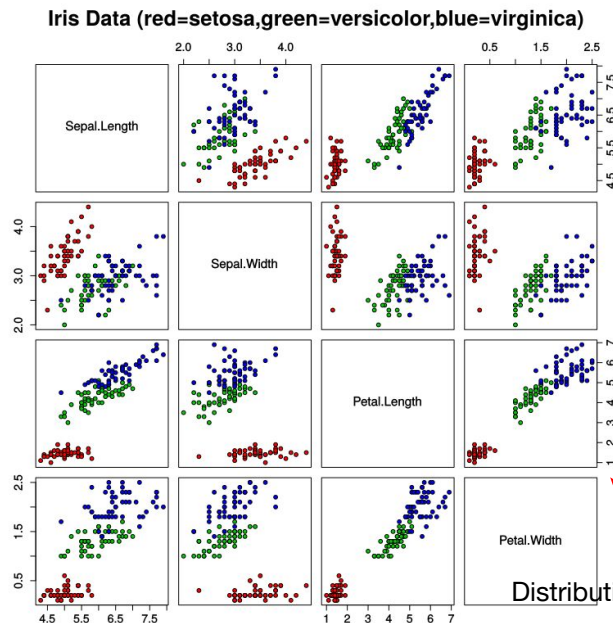
versicolor



virginica

Dataset:

https://docs.google.com/spreadsheets/d/1jpUmKBL_eW_U76f_Efozp9YItER5MkxCRLnUsVOI8xg/edit?usp=sharing



iris.csv: the raw dataset

We can load it by `open('iris.csv', 'r')`

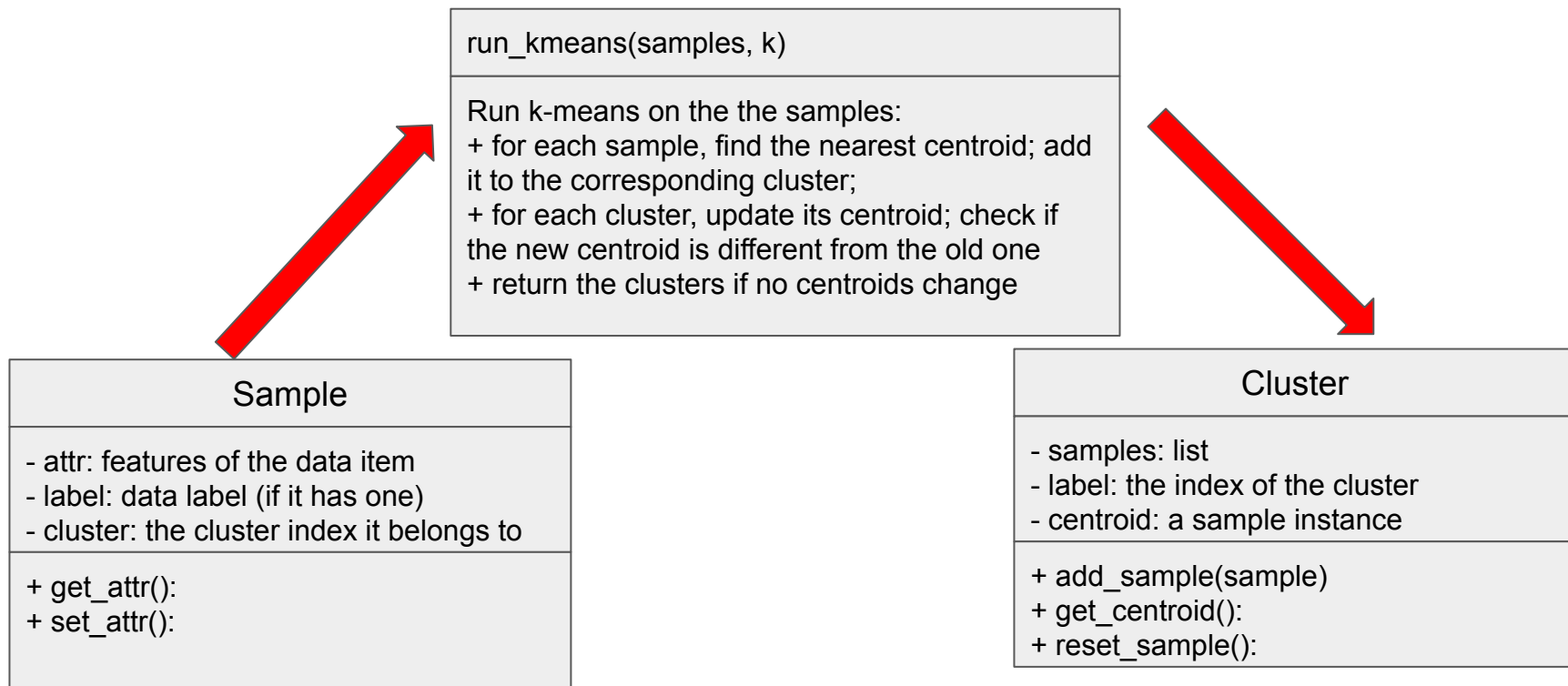
sepal.length	sepal.width	petal.length	petal.width	variety
5.1	3.5	1.4	0.2	Setosa
4.9	3	1.4	0.2	Setosa
4.7	3.2	1.3	0.2	Setosa
4.6	3.1	1.5	0.2	Setosa
5	3.6	1.4	0.2	Setosa
5.4	3.9	1.7	0.4	Setosa
4.6	3.4	1.4	0.3	Setosa
5	3.4	1.5	0.2	Setosa
4.4	2.9	1.4	0.2	Setosa

```
88     f = open('iris.csv', 'r')
89     raw_data = f.readlines()
```

```
"sepal.length","sepal.width","petal.length","petal.width","variety"\n
'5.1,3.5,1.4,.2,"Setosa"\n', '4.9,3,1.4,.2,"Setosa"\n',
'.7,3.2,1.3,.2,"Setosa"\n', '4.6,3.1,1.5,.2,"Setosa"\n',
```

Note: we can build a class to represent the data.

Abstract Data Types (ADT) for the iris dataset



K-means code

Initialization: **k** empty clusters, each centroid is at one of **k** randomly picked members

Iterate until converge:

1. Remove the samples of each cluster to be an empty list (reset);
2. Each sample is assigned to its nearest cluster;
3. Each cluster computes its new centroid;
4. **If** all centroids don't change, done! (converged)
5. Otherwise, **go to** 1. (repetition)

```
103 def kmeans(samples:list, k=2):
104
105     clusters = []
106     for i in range(k):
107         cluster = Cluster([], label=i)
108         clusters.append(cluster)
109
110     centroids = []
111     for c in clusters:
112         init_centroid = random.choice(samples)
113         centroids.append(init_centroid)
114
115     centroid_changed = True
116     max_iteration = 100
117     iterations = 0
118
119     while centroid_changed or iterations < max_iteration:
120         for c in clusters:
121             c.reset_samples()
122
123         for s in samples:
124             cluster_idx = find_nearest_centroid(s, centroids)
125             clusters[cluster_idx].add_sample(s)
126
127         new_centroids = []
128         for c in clusters:
129             new_centroid = centroid_update(c)
130             new_centroids.append(new_centroid)
131
132         for i in range(k):
133             dist = new_centroids[i].distance(centroids[i])
134             if dist > 0.001:
135                 centroid_changed = True
136                 centroids = new_centroids
137                 break
138             else:
139                 centroid_changed = False
140         iterations += 1
141     return clusters
```

Assigning samples to the nearest cluster

Updating centroids

Checking convergence

Sample class

```
12 class Sample:
13
14     def __init__(self):
15
16         self.attr = []
17         self.label = None
18         self.cluster = None
19
20     def get_attr(self):
21         return self.attr
22
23     def set_attr(self, attr):
24         self.attr = attr
25
26     def set_label(self, newLabel:str):
27         self.label = newLabel
28
29     def get_label(self):
30         return self.label
31
32     def set_cluster(self, newCluster):
33         self.cluster = newCluster
34
35     def get_cluster(self):
36         return self.cluster
37
```

```
38     def distance(self, otherSample):
39         '''compute the Eculidian distatnce between
40            this sample to the other.
41            ...
42            d = 0
43
44            for i in range(len(self.attr)):
45                d += (self.attr[i] - otherSample.get_attr()[i])**2
46
47            return d**0.5
48
49
50     def __add__(self, sample):
51         new_s = Sample()
52         # print(self.data)
53         attr = []
54         for i in range(len(self.attr)):
55             attr.append(self.attr[i]+sample.get_attr()[i])
56         new_s.set_attr(attr)
57         return new_s
58
59
60     def __truediv__(self, n):
61
62         new_s = Sample()
63         attr = []
64         for i in range(len(self.attr)):
65             attr.append(self.attr[i]/n)
66
67         new_s.set_attr(attr)
68         return new_s
69
```

Operator overloading of
“+”

Operator overloading of
“/”

Cluster class

```
9
10 class Cluster:
11
12     def __init__(self, index:int):
13
14         self.index = index
15         self.samples = []
16         self.centroid = None
17
18
19     def add_sample(self, sample):
20         self.samples.append(sample)
21
22
23     def get_centroid(self):
24
25         '''compute the centroid of the cluster'''
26
27         return None
28
29     def reset_samples(self):
30         self.sample = []
```

To get the centroid, we need to sum all samples in the cluster and divide by the number of samples. [Hint: you should use the + and / of the samples.]

Complexity of k-Means

n: number of points

k: number of clusters

d: feature dimension of each point

i: number of total iterations

The time complexity of k-mean is $O(nkdi)$; can you explain why?

Unsupervised learning

- It is a type of machine learning where the model learns patterns or structures from unlabeled data without explicit supervision or guidance.
- When the task lacks of labeled data, one may use unsupervised learning methods to discover the intrinsic structures of the dataset
- Applications:
 - clustering similar documents
 - segmenting customers based on their purchasing behavior
 - detecting anomalies in data
 - reducing the dimensionality of high-dimensional data
 -

Homeworks

[ML] Hands-on - Finish in-class exercises:

- Regression with Gradient Descent: page 22
- Clustering with K-means: page 48

[ML] Reading - Chapter 22-24: ML, Clustering, Classification

- Brightspace/Reference Books/*Introduction to Computation and Programming Using Python*, by John V. Guttag

Appendix: k-means using sklearn

sklearn: Scikit-learn

- Sklearn is a free machine learning library, which provide you with a quick way to implement data science tools in your practices.
- It has various built-in functions for classification, regression, and clustering
 - knn, regression, and k-mean are all available
 - official documents are under good maintenance, with a variety of examples showing how to use them.

We will use the K-Means algorithm in sklearn.cluster: [sklearn.cluster.KMeans](#) — [scikit-learn 1.1.1 documentation](#) . Without too many modifications, we can replace our logic of k-means with the sklearn k-means.

Using the K-Means in Sklearn

Examples

```
>>> from sklearn.cluster import KMeans
>>> import numpy as np
>>> X = np.array([[1, 2], [1, 4], [1, 0],
...               [10, 2], [10, 4], [10, 0]])
>>> kmeans = KMeans(n_clusters=2, random_state=0, n_init="auto").fit(X)
>>> kmeans.labels_
array([1, 1, 1, 0, 0, 0], dtype=int32)
>>> kmeans.predict([[0, 0], [12, 3]])
array([1, 0], dtype=int32)
>>> kmeans.cluster_centers_
array([[10.,  2.],
       [ 1.,  2.]])
```

- Need to specify `n_clusters` (i.e., `k`);
- Call `.fit(X)` to run k-means on `X`;
- `X` is a list (or numpy array); each element is a list;
- Call `.labels_` to access the labels of elements in `X`;

run_kmeans

```
18 def run_kmeans(samples, attributes, k=2):
19     data = []
20     for s in samples:
21         attrs = []
22         for attr in attributes:
23             attrs.append(s.get(attr))
24         data.append(attrs)
25
26     kmeans = KMeans(k).fit(data)
27     labels = kmeans.labels_
28     clusters = []
29     '''
30     define a Cluster class; it is a group of samples
31     create k clusters; each cluster should have a label
32     assign each sample to the cluster according to its label
33     '''
34     return clusters
```

K-Means: a class; please use help to view the manual

- .fit(x) method will run k-means on X
- X is an array like dataset; each row is a list of features of a sample
- .labels_: the labels of each row of X

We need a Cluster class to represent the group of samples!