

20210929

Wednesday, September 29, 2021 3:56 PM

- All of us should have received emails if you are doing a Spark project
- Team projects are difficult
 - o Personalities
 - o Delegate and work together
- Meet with your team regularly:
 - o Even if for 5 minutes so that we can be aware of issues as fast as possible
 - HW, midterm, illnesses, etc
 - So you can budget for it
 - Be agile with what people can handle
 - If someone can't do something, then you need to adjust the scope of the project rather than just redistributing work
 - Talk with project manager (not the client) for these kinds of things
 - ◆ Ask for advice on regular basis
- Can escalate to the Prof if necessary
- If you have questions for general purpose of project
 - o Talk to project manager
 - o Minimize the impact on the client
 - o Project manager already interfacing with client on regular basis
- HW1 is due soon
 - o You should have all started
 - o It is challenging

=====

DB-SCAN: note I have a bug in the following,

```
C:\Users\Wolfs\anaconda3\envs\cs506_20210927lec\python.exe C:/Users/Wolfs/PycharmProjects/cs506_20210927lec/05-dbscan/main.py
C:\Users\Wolfs\PycharmProjects\cs506_20210927lec\05-dbscan\main.py:17: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths) from a 1-d array will result in an array of objects with length n. To avoid this, set np.RaggedTensor.CreateInstanceFromNestedSequencesEnabled=False.
  plt.scatter(X[:, 0], X[:, 1], color=colors[clustering].tolist(), s=10, alpha=0.8)
Traceback (most recent call last):
  File "C:\Users\Wolfs\PycharmProjects\cs506_20210927lec\05-dbscan\main.py", line 17, in <module>
    plt.scatter(X[:, 0], X[:, 1], color=colors[clustering].tolist(), s=10, alpha=0.8)
IndexError: only integers, slices ('[:]'), ellipsis ('...'), numpy.newaxis ('None') and integer or boolean arrays are valid indices
Process finished with exit code 1
```

Looks like you're using NumPy
Would you like to turn scientific mode on?
Use scientific mode Keep current layout...

```

1 # Always want to minimize debugging space
2 # Try to get ahead of what he's coding and find flaws
# before he does
3 # Bredford???? search
4 # Top down approach
5
6
7 # Get distance:
8 from .sim import euclidean_dist
9
10 class DBC():
11
12     def __init__(self, dataset, min_pts, epsilon):
13         self.dataset = dataset
14         self.min_pts = min_pts
15         self.epsilon = epsilon
16
17     def _get_epsilon_neighborhood(self, idx):
18         neighborhood = []
19         for i, y in enumerate(self.dataset):
20             # There could be duplicate points
21             # We may not want to compare a point with
# literally itself
22             if i == idx:
23                 continue
24             if euclidean_dist(self.dataset[idx], y
) <= self.epsilon: # Then the point is in the epsilon
# neighborhood of x
25                 neighborhood.append(i) # Just
# tracking the index of the point
26         return neighborhood
27
28     def _make_bfs_assignments(self, queue, assignment
, assignments):
29         # We need to explore to and know if something
# is a core point and therefore it's neighborhood needs
# to be
30         # explored and added to queue
31         # If its a border point, it gets assignment
# but it doesn't get its neighborhood added to the
queue

```

```

32
33         while queue:
34
35             nextP = queue.pop(0)
36
37             if assignments[nextP] != 0:
38                 continue
39
40             assignments[nextP] = assignment # We
41             shouldn't be overwriting will go on forever otherwise
42
43             neighborhood_of_nextP = self.
44             _get_epsilon_neighborhood(nextP)
45             if len(neighborhood_of_nextP) >= self.
46             min_pts:
47                 # Core point
48                 queue += neighborhood_of_nextP
49
50             return assignments
51
52     def dbSCAN(self):
53         """
54             returns a list of assignments. The index
55             of the
56             assignment should match the index of the
57             data point
58             in the dataset.
59         """
60
61         # assignment vs assignments is not great
62         # naming
63         assignments = [0 for _ in range(len(self.
64         dataset))]
65         assignment = 1
66         # Iterate over points in dataset,
67         # If you iterate over points themselves
68         # You need to assignment of the point to match
69         # the index of the point
70         # So this can get messy
71         # So iterate over indexes instead
72         # Or use enumerate to get index and the point
73
74

```

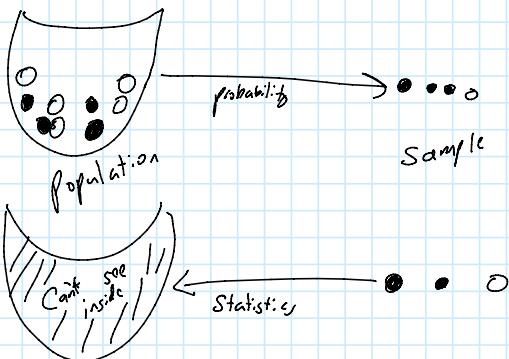
```
65     for i, x in enumerate(self.dataset):
66         if assignments[i] != 0:
67             continue # Skip the rest of this
68             iteration of the loop
69             # Need to determine whether or not X is a
70             core point
71             # This statement doesn't make sense,
72             because you want to get the actual neighborhood from
73             them helper:
74             # if self._get_epsilon_neighborhood(x
75             ) >= self.min_pts:
76                 neighborhood = self.
77                 _get_epsilon_neighborhood(i)
78                 if len(neighborhood) >= self.min_pts:
79                     # Core point
80                     # Explore x neighborhood for core
81                     points and give them the same assignment
82                     # Instead of doing that now, we are
83                     going to save:
84                     # Helper function
85                     assignments[i] = assignment
86                     assignments = self.
87                     _make_bfs_assignments(neighborhood, assignments,
88                     assignments)
89                     assignment += 1
90                     # We need to account for stumbling across
91                     a core point a second time
92                     # Without a core point, you can't really
93                     do anything
94                     # Everything is initially labeled noise
95                     # You know it is a border point if it's
96                     in a core point's neighborhood, but isn't a core
97                     point
98                     # Will iterate over queue and pop off
99                     return assignments
100
101 # Try to do on your own before next time
102 # Someone please upload this code to github
103 # My code isn't working properly
```

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import sklearn.datasets as datasets
4 from sklearn.preprocessing import StandardScaler
5
6 from cs506 import dbscan
7
8 centers = [[1, 1], [-1, -1], [1, -1]]
9 X, _ = datasets.make_blobs(n_samples=750, centers=
    centers, cluster_std=0.4,
10                           random_state=0)
11 plt.scatter(X[:, 0], X[:, 1], s=10, alpha=0.8)
12 plt.show()
13
14 clustering = dbscan.DBC(X, 3, .2).dbscan()
15 colors = np.array([x for x in
    'bgrcmykbgrcmykbgrcmykbgrcmyk'])
16 colors = np.hstack([colors] * 20)
17 plt.scatter(X[:, 0], X[:, 1], color=colors[clustering
    ].tolist(), s=10, alpha=0.8)
18 plt.show()
19
```

20210929 Probability

Wednesday, September 29, 2021 4:59 PM

Probability vs. statistics \rightarrow reasoning in opposite directions



Example: Coin Toss

$$\Omega = \{H, T\}$$

Omega: set of all possible outcomes

script f: all possible events

$$\mathcal{F} = \{\emptyset, \{H\}, \{T\}, \{H, T\}\}$$

$$P: \mathcal{F} \rightarrow [0, 1]$$

(Ω, \mathcal{F}, P) \rightarrow a probability space

where

- $P(\Omega) = 1$

- $P(A) \geq 0, A \in \mathcal{F}$

any possible outcome prob.
is @ 1st O

- $P(\bigcup_{i=1}^n A_i) = \sum_{j=1}^n P(A_j)$

if $A_i \cap A_j = \emptyset$

double check
subscripts!

assuming not

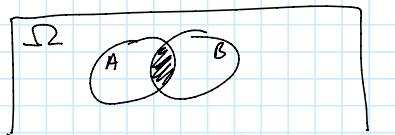
- $P(A^c) = 1 - P(A)$

- $P(A) \leq 1 \rightarrow P(A^c) \geq 0$

- $P(\emptyset) = 0$

- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

subtract intersection



- $P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B | A) P(A)}{P(B)}$

given

* A & B are independent if:

$$P(A \cap B) = P(A) \cdot P(B)$$

Bayes Rule

$$\left\{ \begin{array}{l} P(A \cap B) = P(A) \cdot P(B) \\ P(A|B) = P(A) \\ P(B|A) = P(B) \end{array} \right.$$

Focus back on coin toss

Distribution characterizes the population

If you list all possible outcomes & their probabilities, then you've characterized the population

Use Random Variable

A random variable $X: \Omega \rightarrow \mathbb{R}$
maps outcomes to real line

$$X(\omega) = \begin{cases} 1 & \text{if } \omega \in H \\ 0 & \text{if } \omega \in T \end{cases}$$

lowercase omega

$$P_x(1) = P(X=1)$$

$$\text{Using } P: \mathcal{F} \rightarrow [0, 1]$$

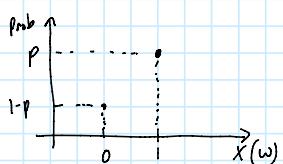
$$= P(\{\omega | X(\omega)=1\})$$

$$= P(H)$$

Extremely rigorous way of describing mapping

What is the distribution of our random var X

Distribution of our coin toss



May be unfair coin

Bernoulli Trial

$$X_1, X_2, \dots$$

1st coin flip 2nd

What is the probability that it will take k # of tosses to get first heads?

Probability that it takes k tosses to get first H.

$$P_{x_1}(0) P_{x_2}(0) P_{x_3}(0) \dots P_{x_k}(1)$$

Assuming independence of random variables

$$= (1-p)^{k-1} \cdot p$$

assuming independence or random variables

$$= (1-p)^{k-1} \cdot p$$

Geometric Distribution

What probability that you get k # of heads in first n coin flips?

Prob of getting k Hs in n tosses?

Binomial Distribution

$$P_n(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

On average, you are getting ~60 heads per hour.

Flip coin over large period of time, on average we get 60 Hs/hour.

Probability of k heads in a minute?

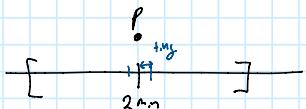
↳ need a different distribution function

Poisson Dist

$$P(k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$\lambda = \mu$ \downarrow down to units
Maximum when $k = \lambda$? \checkmark double check

What is probability that b/w 2 successes is exactly 2 min ?
↳ ∞ # of values that are very close to 2 min value



No way to sum ∞ # of values & have it equal one

↳ So this? doesn't make sense for probability

B/w 1 min & 2 min makes more sense in probability

Time domain is continuous, but we have been talking about discrete outcomes.

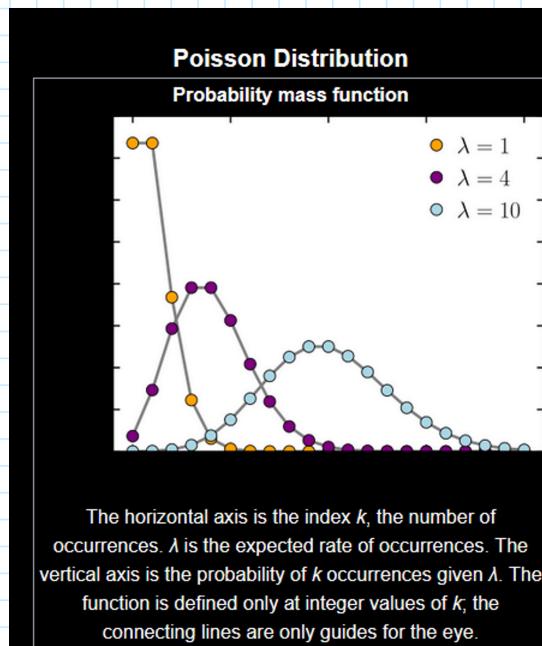
Can't talk things w/ continuum of values.

With a continuum of possible outcomes, we want to look @
 $P(X \leq t)$ ← Cumulative Distribution

Exponential Distribution

$$P(\text{time b/w consecutive Hs}) = e^{-\lambda t}$$

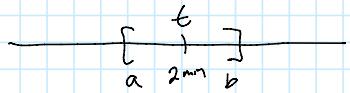
https://en.wikipedia.org/wiki/Poisson_distribution



The horizontal axis is the index k , the number of occurrences. λ is the expected rate of occurrences. The vertical axis is the probability of k occurrences given λ . The function is defined only at integer values of k ; the connecting lines are only guides for the eye.

Instantaneous Prob. dist.

Instantaneous change of prob @ point



as $\Delta t \rightarrow 0$, what happens in this infinitesimally small range

$$f(t) = \mu e^{-\mu t} \quad (\text{is a derivative})$$

Probability Density Function

$$P(t \in [t_1, t_2]) = \int_{t_1}^{t_2} \mu e^{-\mu t} dt = P(t \geq t_1, t \leq t_2)$$

If you were to average the time b/w successive heads?

↳ What is this probability distribution

↳ normal distribution

↳ Central limit theorem

All these different distributions can be covered
just by flipping a coin