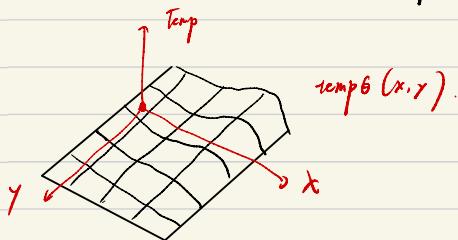


Note for Sep 8th

Data Science.

	M		H

(x, y) record of a temp



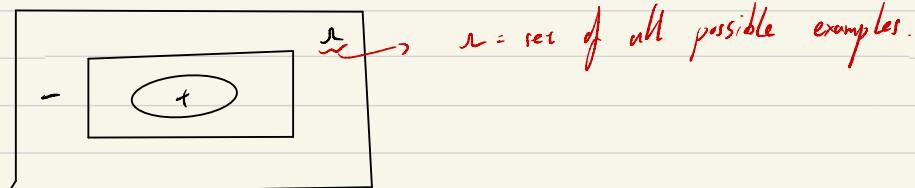
temp6(x, y).

Find a function $f(x, y, v) \rightarrow \text{temp}$

↓
found by testing known information, and explain hindsight or predicted future.

Positive examples: Expect positive results under hypothesis.

Negative ----; - - negative --- - - - .



The venn diagram shows the benefit of negative examples.

Note for Sep 13th.

Minimal progress loss → Backup work + commits to avoid loss

Collab process : - each collaborator creates a new branch.

- provide process for merging their branch into master. → PR (pull request)

key step: ① Fork CSS36 public repo

② Make all changes to your repo ★ keep copy up-to-date

③ request via pull request.

method: Using git upstream

0. git checkout master
2. git fetch upstream master
3. git rebase upstream/master

Work-flow:



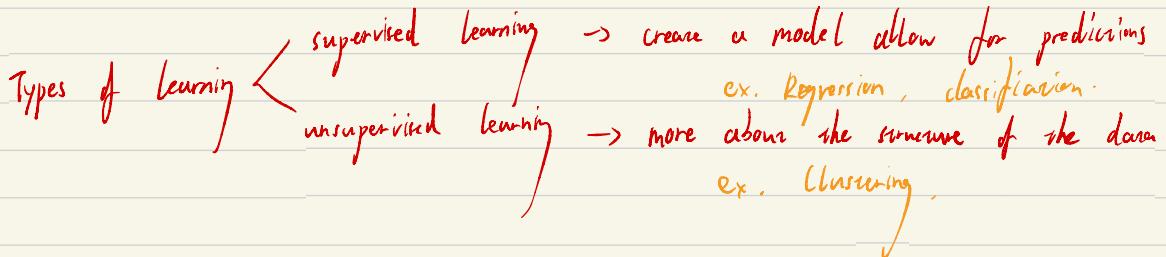
Note for Sep 15th.

max path sum

$$\begin{matrix} & 3 & & 10 & 7 \\ & 7 & 4 & 2 & 4 & 5 \\ 2 & 4 & 3 \end{matrix} \quad 12 \rightarrow \max(10+4, 7+4)$$

Data Representations

- Records (m dimensional vector as tuple array , holds different attributes)
- graphs (nodes connected with edges)
- image (matrix of RGB pixels).
- Text (list of words).
- String (list of characters)
- Time series (list of data at specific intervals of time).



Note for Sep 20th

Feature Space - set of all possible data points

Cosine similarity $s(x, y) = \cos \theta$ where θ is angle b/w x, y

using dissimilarity function : $d(x, y) = \frac{1}{s(x, y)}$

As $d(x, y) = k - s(x, y)$ for some k .

$$\underline{d(x, y)} = 1 - s(x, y).$$

* We use cosine over euclidean distance when direction matters more than magnitude.

Jaccard similarity : ex. take zero and each word is an attribute, and values are whether it is in the given doc.

$$J_{sim}(x, y) = \frac{|x \cap y|}{|x \cup y|} \quad J_{dist}(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$

Note Sep 22th:

Clustering → create groupings / assignments of objects s.t. 1. similar to each other
2. dissimilar to other clusters.

Clustering problem: • similar data points in same clusters
• dissimilar data points in different clusters.

- Types:
- ① Partitional
 - ② Hierarchical
 - ③ Density-based
 - ④ soft clustering

K-means

$$\sum_i \sum_{x \in C_i} \|x - \mu_i\|^2 \Rightarrow \text{Find } k \text{ points that minimize the cost function}$$

Note: easy if $k=1$ or $k=n$.

- 1 [Steps]:
- 1. Randomly pick k centers
 - 2. Assign each point in dataset to its closest center
 - 3. Compute the new centers as means of each cluster
 - 4. Repeat 2, 3 until converge

- Limitations:
- ① splits large number
 - ② dislike nonlobular cluster shape.
 - ③ dislike varying density

Optimal k selection: Elbow method:



Note Sep 27th

Hierarchical clustering :

- output: tree of clusters showing closest related points (We can cut at a special threshold to generate different)

Agglomerative Clustering Algorithm:

- ① each point in dataset be in its own cluster
- ② Compute distance b/w all pairs of clusters
- ③ merge 2 closest clusters
- ④ Repeat 2, 3 until converge.

Single link distance: $D_{SL}(C_1, C_2) = \max\{d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2\}$.

Avg Link distance: $D_{AL}(C_1, C_2) = \frac{1}{|C_1||C_2|} \cdot \left(d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2 \right)$

Centroid Distance: $D_C(C_1, C_2) = d(\mu_1, \mu_2)$

Ward's distance: $D_{WD}(C_1, C_2) = \sum_{p \in C_1} d(p, \mu_1) - \sum_{p \in C_2} d(p, \mu_1) - \sum_{p \in C_2} d(p, \mu_2)$

DB scan Algorithms: ① Find ϵ neighbourhood of each point

② label core if \geq min-points

③ Non core, but in neighborhood

④ else noise

⑤ For each core, assign to same cluster

⑥ Assign border points to nearby clusters.