

Exam 2

Evan Nuss

0	0	-2
1	10	

$$Q((0,0), \text{down}) = 10$$

$$Q((1,1), \text{left}) = 10$$

$$Q((0,1), \text{left}) = 0 + .7(10) = 7$$

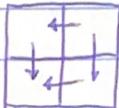
$$Q((0,1), \text{down}) = 0 + .7(10) = 7$$

$$Q((0,0), \text{right}) = -2 + .7(7) = 2.9$$

$$Q((1,1), \text{up}) = -2 + .7(7) = 2.9$$

0	1
10	2.9
	-2.9

optimal policy:



~~Value Function~~

~~a. Pick random action~~

~~reward = state reward + discount factor * highest rewards for staying in new states~~

b. state = initial_state

Do until optimal policy determined

state. pick_random_action

$Q_{\text{reward}}(\text{state}, \text{action}) = \text{reward}_{\text{new-state}} + \text{discount-factor} * \text{highest } Q\text{-value}$ for leaving new-state

state = new-state

2a. centroid 1 = (1, 1)
centroid 2 = (1, 2)

$$dist_{AC} = |1-3| + |2-5| = 6 \quad dist_{AD} = |1-5| + |2-6.5| = 9.5$$

$$dist_{BC} = |2-3| + |2-5| = 5 \quad dist_{BD} = |1-5| + |2-6.5| = 8.5$$

cluster1 = {A, 3}
cluster2 = {B, C, D}

centroid 1 = (1, 1)

$$\frac{1+3+5}{3} = 3 \quad \frac{2+5+6.5}{3} = 4.5$$

centroid 2 = (3, 4.5)

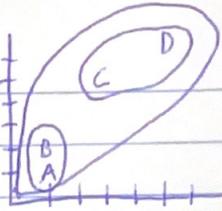
dist_{C2-B} = |3-1| + |4.5-2| = 4.5 dist_{C2-B} = |1-1| + |1-2| = 1

$$dist_{C2-C} = |3-3| + |4.5-5| = 0.5$$

$$dist_{C2-D} = |3-5| + |4.5-6.5| = 4$$

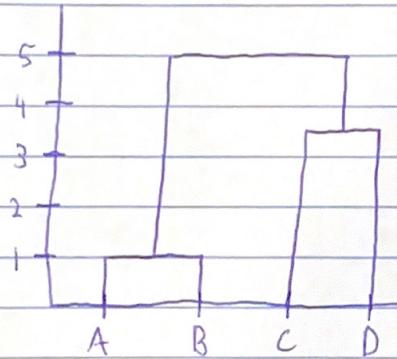
cluster 1 = {A, B}	centroid 1 = (1, 1.5)
cluster 2 = {C, D}	centroid 2 = (4, 5.75)

b. $\text{dist}_{AB} = 1$ $\text{dist}_{AC} = 6$ $\text{dist}_{AD} = 9.5$
 $\text{dist}_{BC} = 5$ $\text{dist}_{BD} = 8.5$
 $\text{dist}_{CD} = |3-5| + |5-6.5| = 3.5$



$\text{dist}_{BC} = 5$ $\text{dist}_{AD} = 9.5$ $\text{dist}_{CD} = 3.5$

$\text{dist}_{BC} = 5$



3a. Genetic Algorithm

$t = 0$

Initialize population $P(t)$

Evaluate $P(t)$

Do until given stopping criteria is met E

$t += 1$

Parent-selection $P(t)$

recombine $P(t)$

Mutate $P(t)$

Evaluate $P(t)$

Survive $P(t)$

b. Assumptions:

- The data is represented such that the attributes of the genome are their own chunk within the genome, meaning any crossover & mutation must not affect only a part of the attribute.
- I will use ~~single~~ 2 point crossover & any mutation will be single point. That point & its new value will be picked at random.
- Parent selection will be based on the pizzas with the best fitness score. However, to maintain diversity, a lower performing pizza will be selected at times.
- Survival will be determined by fitness score. The pizzas with the best fitness scores will survive, as well as some with lower scores in order to maintain diversity.

crust	amount of cheese	meat/veggie ratio	# of pepperoni
P_1 = [thin, high]	,	4:1	20
P_2 = [stuffed, medium]	,	1:6	2
P_3 = [thick, high]	,	2:1	10
P_4 = [thin, low]	,	1:9	0

- fitness score is on a scale 0-10

$$F_1 = 7 \quad F_2 = 5 \quad F_3 = 9 \quad F_4 = 3$$

- P_2 & P_3 will be parents. We will do 2 point crossover & single point mutation.

$$P_5 = [\text{stuffed, high, } 1:6, 10] \Rightarrow \boxed{\text{stuffed, high, } 1:3, 10}$$

$$P_6 = [\text{thick, medium, } 2:1, 2] \Rightarrow [\text{thick, medium, } 2:1, 10]$$

$$- F_5 = 9 \quad F_6 = 8$$

- Kill off P_4 & P_1 to maintain population of 4

- P_2 & P_3 will be parents

$$P_7 = [\text{stuffed, high, } 1:6, 10] \Rightarrow [\text{thin, high, } 1:6, 10]$$

$$P_8 = [\text{stuffed, medium, } 1:3, 2] \Rightarrow [\text{stuffed, low, } 1:3, 2]$$

$$F_7 = 8 \quad \text{kill off } P_3 \& P_2$$

$$F_8 = 2$$

$$\boxed{\text{population} = \{P_3, P_5, P_6, P_7\}}$$

4a. prior-p = 4/7 prior-n = 3/7
L \Rightarrow likelihood

$$\begin{aligned} L_{\text{sev or } p} &= 2/4 & L_{\text{sof or } p} &= 1/4 & L_{\text{fresh } p} &= 1/4 \\ L_{\text{med } p} &= 3/4 & L_{\text{hi } p} &= 1/4 & L_{\text{lo } p} &= 0/4 \\ L_{\text{true } p} &= 3/4 & L_{\text{false } p} &= 1/4 \end{aligned}$$

$$L_{\text{sev or } n} = 2/3 \quad L_{\text{med } n} = 1/3 \quad L_{\text{false } n} = 2/3$$

posterior Find posterior for each output class based on new point

$$\begin{aligned} \text{posterior-}p &= \text{prior-}p (L_{\text{sev or } p}) (L_{\text{med } p}) (L_{\text{false } p}) \\ &= 4/7 (2/4) (3/4) (1/4) \\ &= .0536 \end{aligned}$$

$$\begin{aligned} \text{posterior-}n &= \text{prior-}n (L_{\text{sev or } n}) (L_{\text{med } n}) (L_{\text{false } n}) \\ &= 3/7 (2/3) (1/3) (2/3) \\ &= .0635 \end{aligned}$$

$$\text{probability-}p = \frac{.0536}{.0536 + .0635} = \boxed{\text{probability-}p = .458}$$

$$\text{probability-}n = \frac{.0635}{.0536 + .0635} = \boxed{\text{probability-}n = .542}$$

b. The assumption is that using only past data will be enough to get ~~more~~ accurate probabilities given a new datapoint. It is "naive" of the current data.

5a. Un-normalized attributes will cause points to be significantly further apart, esp if there is a feature with a very wide range of values. This will mean distances between points are often much higher, which will decrease the effect that feature has on the predicted class, which could be very detrimental to the model's accuracy, especially if using inverse distance weighting.

b. An irrelevant feature that is one of the K nearest neighbors will ~~never~~ still have a say in what the output value ends up being. It could have as much as, or even more influence on the final prediction, as a very relevant feature, depending on how close it is to the point to be classified. This influence will be very detrimental to the model's accuracy.

6a. Bagging:

$b = \# \text{ of models}$

for $x \in \text{range}(b)$:

$\text{training_sets}[x] = \text{Uniform randomly sampled subset of dataset}$

for $i \in \text{range}(b)$

$\text{models}[i].\text{train}(\text{training_sets}[i])$

for $\text{model} \in \text{models}$:

$\text{outputs.append(model.\text{predict(whole_dataset, datapoint)})}$

$\text{output} = \text{class that appears most in outputs}$

Boosting:

X = dataset

D = distribution of how likely each datapoint is selected.

For each model:

subset = sampling of X using D

model. train(subset)

Update D based on points model incorrectly predicted

For each model:

vote = model.predict * model.training_accuracy

votes[i] += vote

~~class~~ output-class = argmax(votes)

b. High variance comes from letting noise affect the model's predictions too much. By using bagging, the many models each have an equal vote in the output class, which means a model with high variance doesn't have a disproportionate influence & if its vote is incorrect, that will become negligible since the majority of the other models will not have ^{the same} predicted value. The results with high variance are unlikely to make it through the voting system.

c. If enough of the models end up overfitting & having really good accuracy on the training set, they will have a much stronger vote in the output class & could cause the whole ensemble to overfit. Whereas, the models that don't overfit might have worse training accuracy, & therefore not as much say in the final output class decision. This is b/c in boosting, votes are weighted by the model's training accuracy.

KNN Manhattan distance: $|x[0] - y[0]| + |x[1] - y[1]| + \dots |KBF|$

- Calculate distances between new point & old ones.
- Choose the K closest points (smallest distance).
- Whichever class has most points in list of K neighbors wins.

Distance Weighting: $\text{outputs} = \text{rbf-results} * \text{weights}$. Whichever output value is greatest, that's the class for that point.

For every class:

$$\text{votes}[class] += \frac{1}{\text{neighbor_distance}^2}$$

Regression:

For each distance-label pair:

$$\text{weight} = 1/\text{distance}^2$$

weights.append(weight)

values.append(label+weight)

output = $\frac{\text{sum of all values}}{\text{sum of all weights}}$

- irrelevant features have big effect. Dimensionality & distances tend toward infinity.

- use PCA, filters/wrapers, attribute weighing, decision tree to limit variance of features.

Bagging: uniformly randomly sample dataset b times, b = # of models. Train each model independently on its dataset. Then for predicting, each model gets whole data set & sets unweighted vote in class.

Naive Bayes / Bayesian learning: $\text{prob}_{\text{class}} = \frac{\text{prior}_{\text{class}} * \prod_{i=1}^n \text{posterior}_{i, \text{class}}}{\text{total}}$

HAC: HAC = Manhattan dist between each point. Shortest distance. If two points form cluster - repeat depending on single link / complete link.

Silhouette: $s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$, where $a(i)$ = avg distance between i & all other points in cluster, $b(i)$ = avg dist b/w i & points in other cluster. $s(i) \in [-1, 1]$. Silhouette having $s(i) > 0$ means i is well clustered.

Un-labeled data use clustering: K-means : $\text{Q}(\text{state action}) = \text{reward of state you be in next} + \gamma \text{discount factor} * \text{highest Q leaving the square you're going to.}$ Genetic algorithm: $E = \frac{1}{n} \sum_{i=1}^n P(x_i)$, $E' = \frac{1}{n} \sum_{i=1}^n P'(x_i)$. $E' < E$: better selection, recombine, mutate, evolution, and survive.

Boosting: X decision, D determines how likely each train a model, given point is to be closer to class by points off be voted by next model. Each model votes on class. If a vote is skewed by its accuracy for its class.

Neural networks: except when we have a bunch of nodes, then it becomes very expensive.