

Assignment #1: Version Spaces, k-NN, Decision Trees

Instructor: Nika Haghtalab, Thorsten Joachims

Name: Kevin Klaben, Netid: kek228

Course Policy: Read all the instructions below carefully before you start working on the assignment, and before you make a submission.

- Please include your name and NetIDs on the first page. We recommend typesetting your submission in L^AT_EX, and an Overleaf template is linked on the canvas module. **When submitting, remember to mark which page has which response.**
- Assignments are due at 5 pm on the due date in PDF form on Gradescope.
- Late assignments can be submitted on Gradescope up to Sunday, Sept 20 at 5pm EST. This is also when the solutions will be released.
- You can do this assignment in groups of 1-2. Please submit no more than one submission per group. Collaboration across groups is not permitted.
- All sources of material outside the course must be cited. The University Academic Code of Conduct will be strictly enforced.

Problem 1

(a)

$$(a) \{ \text{mild, medium, hot} \} \times \{ \text{no meat, chicken, beef} \} \times \{ \text{yes, no} \}$$
$$3 \times 3 \times 2 = 18$$

Instance space of size 18

(b)

$$(b) 2^{18} = 262144 \text{ is the size of } H$$

(c)

(c) Clearly, this hypothesis is not consistent with the table. By definition h is consistent if and only if $h(x) = y$ for all (x, y) in the table. By Bob's hypothesis, $h(\text{any entry with hot})$ should be yes, however the entry for lentil curry contradicts this as the dish is hot but Alice does not like it. Thus Bob's hypothesis is not consistent with table T.

(d)

(d) With 5 of the 18 entries determined, the size of the version space is 2^{13} or $VS_{(H,T)} = 8192$.

(e)

(e) 3 possible combinations

$S+M$ | $M+C$

$$h_1(S,M) = \begin{cases} \text{yes, if } S+M > 0 \\ \text{no, otherwise} \end{cases} \quad h_2(M,C) = \begin{cases} \text{yes, if } M+C > 0 \\ \text{no, otherwise} \end{cases}$$

$$h'_1(S,M) = \begin{cases} \text{no, if } S+M > 0 \\ \text{yes, otherwise} \end{cases} \quad h'_2(M,C) = \begin{cases} \text{no, if } M+C > 0 \\ \text{yes, otherwise} \end{cases}$$

$S+C$ $h_3(S,C) = \begin{cases} \text{yes, if } S+C > 0 \\ \text{no, otherwise} \end{cases}$ $h'_3(S,C) = \begin{cases} \text{no, if } S+C > 0 \\ \text{yes, otherwise} \end{cases}$

$H' = \{h_1, h'_1, h_2, h'_2, h_3, h'_3\}$, thus H' has size 6.

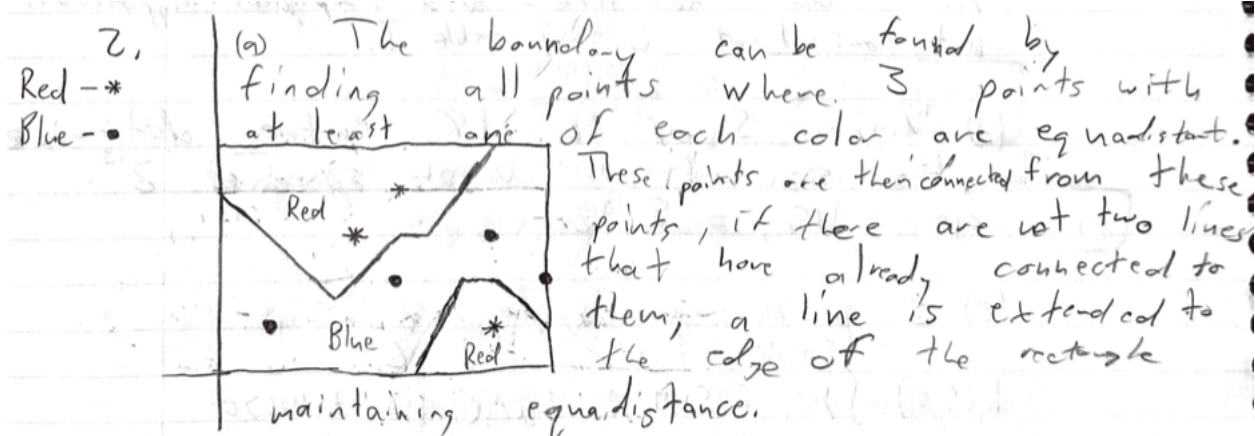
(f)

(f) Dish	S	M	C	h_1	h_2	h_3	h'_1	h'_2	h'_3
spicy chicken	1	0	1	yes	yes	yes	no	no	no
pepper steak	0	1	-1	yes	no	no	no	yes	yes
beet stew	-1	1	-1	no	no	no	yes	yes	yes
lentil curry	1	-1	1	no	no	yes	yes	yes	no
chicken soup	0	0	-1	no	no	no	yes	yes	yes

Clearly h_1 is the only function of $VS_{(H',T')}$ as it is consistent for T' . h_2, h_3 , and h'_1 are inconsistent due to the incorrect output for pepper steak whereas h'_2 and h'_3 are inconsistent due to incorrect output for spicy chicken entry. Thus, the size of $VS_{(H',T')}$ is 1.

Problem 2

(a)



(b)

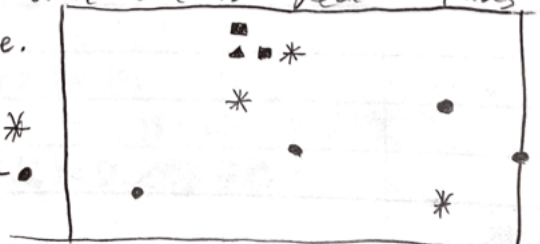


When working with real data this effect becomes a problem as points in the original un-scaled boundaries and data may be classified differently than when they are scaled. For example, the point (0,5) would be classified as Red, however following scaling, and the boundaries moving, this now scaled point (0,5) would be blue. One solution to mitigate these effects would be to simply scale both axis equally.

(c)

(c) Under the original set of points, $(3,4)$ would be classified as Red under 3-NN as the three nearest points are $(3,3)$, $(4,4)$, and $(4,2)$ of which 2 are Red. To change the classification of $(3,4)$ under 3-NN, a minimum of 2 points must be added such that they are now the closest points to $(3,4)$ such as $(3.5,4)$ and $(3,4.5)$ both Blue. Thus, the nearest 3 points are $(3.5,4)$, $(3,4.5)$ and a tie for third between $(3,3)$ and $(4,4)$ and two are Blue and one is red thus classification becomes Blue.

test point $(3,4)$ - \blacktriangle Red - $*$
 $(3.5,4)$ and $(3,4.5)$ - \blacksquare Blue - \bullet



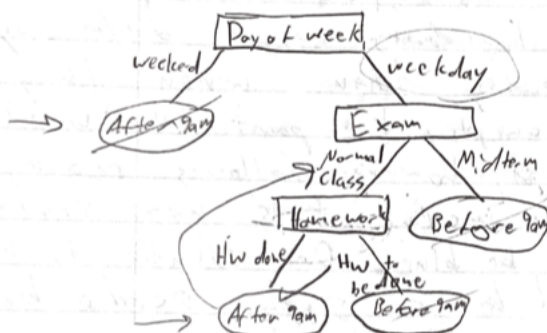
Problem 3

(a)

3.) (a) • weekend, normal class, homework to be finished, then student A will get up after 9am

• weekday, midterm day, and homework is done, then student A will get up before 9am

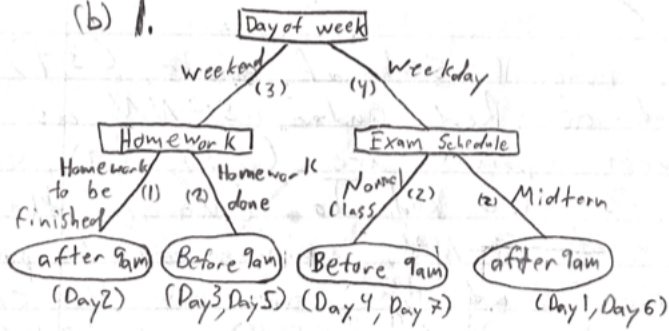
• weekday, normal class, homework is done, then student A will get up after 9am



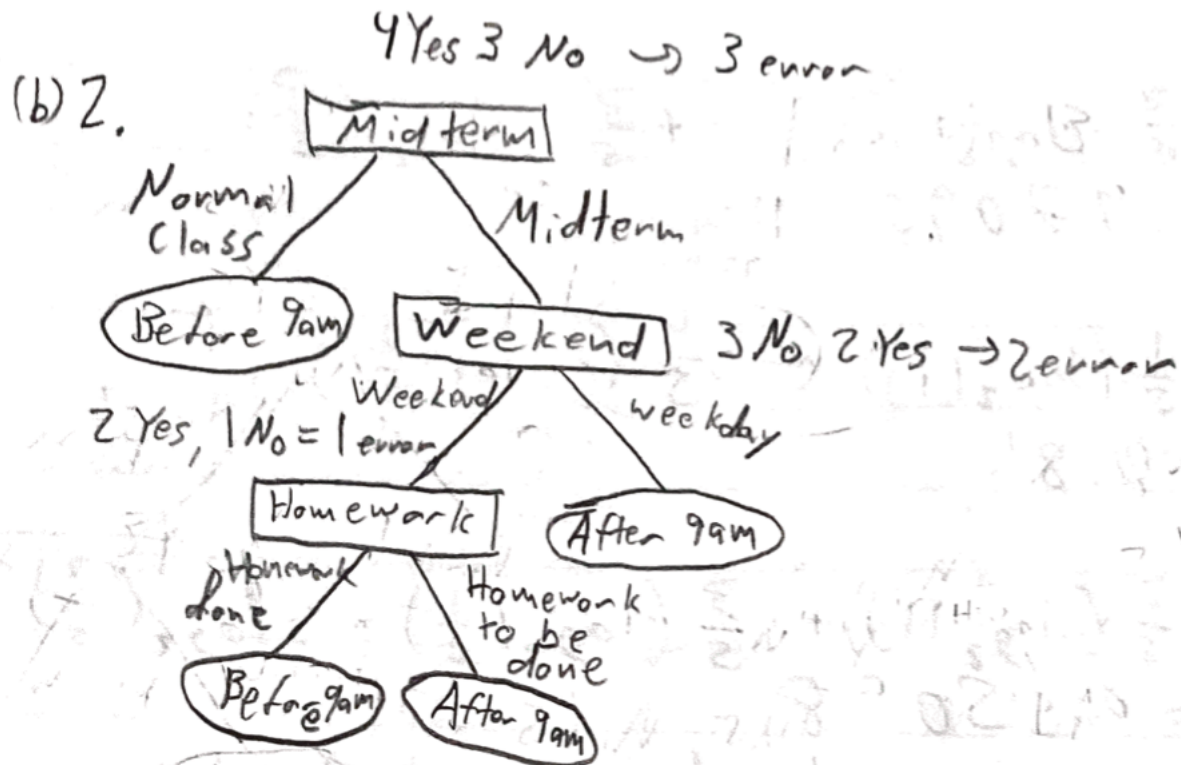
If student A woke up after 9am on a weekday, then it was not a midterm day.

(b)1.

(b) 1.



(b) 2.



(b) 3.

(b) 3.
$$H(S) = - \sum_{y \in Y} P_S(y) \log_2(P_S(y))$$

$$= - \left(\left(\frac{3}{7} \right) \log_2 \left(\frac{3}{7} \right) + \left(\frac{4}{7} \right) \log_2 \left(\frac{4}{7} \right) \right)$$

$$= 0.985228$$

$$H(S | \text{weekend}) = \sum_{a_i} P_S(S_i) H(S_i)$$

$$= \frac{4}{7} H(S_{\text{yes}}) + \frac{3}{7} H(S_{\text{no}})$$

$$= \frac{4}{7} \left(- \left(\left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) + \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) \right) \right) +$$

$$\frac{3}{7} \left(- \left(\left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right) + \left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) \right) \right)$$

$$= \frac{4}{7} \left(- \log_2 \left(\frac{2}{4} \right) \right)$$

$$= 0.96498$$

$$H(S | \text{Homework}) = \frac{2}{7} \left(- \left(\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) \right)$$

$$+ \frac{5}{7} \left(- \left(\left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) + \left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) \right) \right)$$

$$= 0.97925$$

$$H(S | \text{Midterm}) = \frac{2}{7} \left(- \left(1 \log_2(1) \right) \right)$$

$$+ \frac{5}{7} \left(- \left(\left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) + \left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) \right) \right)$$

$$= 0.69354$$

Maximizing $IG(S) = H(S) - H(S|A)$

when using midterm $IG(S) = H(S) - H(S | \text{midterm})$

$$= 0.985228 - 0.69354$$

$$IG(S) = 0.29270$$

$$H = \frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right) \\ \approx 0.97095$$

$$H_w = \frac{1}{5} (\log_2(1)) + \frac{4}{5} \left(\log_2\left(\frac{2}{5}\right) \right) \\ = 0.8$$

$$\text{weekend} = \frac{2}{5} (\log_2(1)) + \frac{3}{5} \left(\log_2\left(\frac{1}{3}\right) + \frac{2}{3} \log_2\left(\frac{2}{3}\right) \right) \\ = 0.55098$$

Max is when weekend is next feature

