

Exam 2 - Rajat Sethi

1.) ~~Ans~~ Use $\sum_{i=1}^n a_i x_i K(x_i, z) + b = f(z)$

$$2.5(1)(2z+1)^2 + 7.3(-1)(5z+1)^2 + 4.83(1)(6z+1)^2 + b = f(z)$$

$$2.5(4z^2 + 4z + 1) + -7.3(25z^2 + 10z + 1) + 4.83(36z^2 + 12z + 1) + b = f(z)$$

$$10z^2 + 10z - 183.3z^2 - 73.3z + 174z^2 + 58z + b = f(z) \quad (\text{All } x^0 \text{ terms are merged into } b)$$

$$f(z) = 0.6z^2 - 5.3z + b$$

Let $z=2$, ~~then~~ $f(z)=1$, $b=9$

$b=9$ for $z=5$ and 6 too

Discriminate Function:

$$f(z) = 0.666...z^2 - 5.333...z + 9$$

$$f(3) = -1$$

2.)

	a_1	a_2	c	a_1	a_2	c
h_1	10	01	1	11	10	0
h_2	01	11	0	10	01	0

Crossover for h_1 - $\langle 1, 3 \rangle$

Crossover for h_2 - $\langle 1, 8 \rangle$

	a_1	a_2	c	a_1	a_2	c	a_1	a_2	c
h_3	11	11	0	10	01	1	11	10	0
h_4	00	01	0						

3.) Lazy Learning - Simply storing the training set and not preprocessing or fixing it

Eager Learning - Preprocessing the dataset to ease the burden on future predictions.

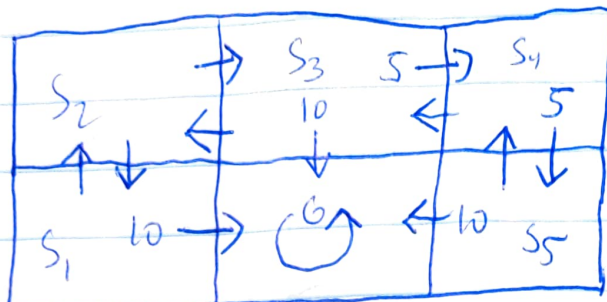
Advantages of IBL:

- Works well for complex datasets

Disadvantages of IBL:

- Computationally expensive

4a)



$$V^* = r + \gamma \cdot V_{n+1}^*$$

$$G = 0$$

$$S_1 = 10 + 0.8 \cdot 0 = 10$$

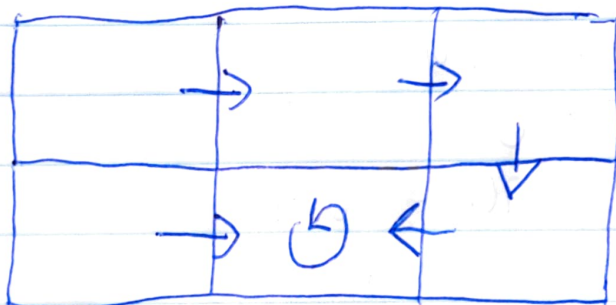
$$S_5 = 10 + 0.8 \cdot 0 = 10$$

$$S_4 = 5 + 0.8 \cdot 10 = 13$$

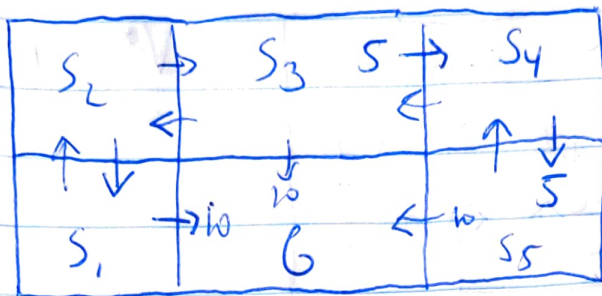
$$S_3 = 5 + 0.8 \cdot 13 = 15.4$$

$$S_2 = 0 + 0.8 \cdot 15.4 = 12.32$$

Optimal Policy:



4B.



Episode 1:

$$\hat{Q}(S_1, \text{right}) = 10$$

$$\hat{Q}(S_1, \text{up}) = 0 + 0 = 0$$

$$\hat{Q}(S_2, \text{down}) = 0 + 0.8(10) = 8$$

$$\hat{Q}(S_2, \text{right}) = 0 + 0 = 0$$

$$\hat{Q}(S_3, \text{left}) = 0 + 0.8(8) = 6.4$$

$$\hat{Q}(S_3, \text{down}) = 10 + 0 = 10$$

$$\hat{Q}(S_4, \text{right}) = 5 + 0 = 5$$

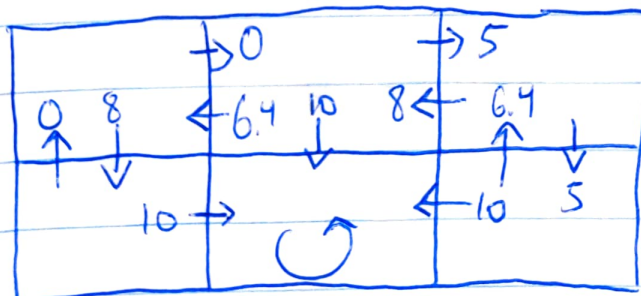
$$\hat{Q}(S_4, \text{left}) = 0 + 0.8(10) = 8$$

$$\hat{Q}(S_4, \text{down}) = 5 + 0 = 5$$

$$\hat{Q}(S_5, \text{up}) = 0 + 0.8(8) = 6.4$$

$$\hat{Q}(S_5, \text{left}) = 10 + 0 = 10$$

\hat{Q} Values after first episode



$$5.) C_2 = (C - (C_1 - \{L\})) \cup \{\sim L\}$$

G exists in C_1 but not C

$$\text{Let } L = G$$

$$C_1 - \{G\} = (A \vee B) \wedge \sim G$$

$$C - (C_1 - \{G\}) = (A \vee B) \wedge G$$

$$(C - (C_1 - \{G\})) \cup \{\sim G\} = \sim G \vee A \vee B$$

$$C_2 = \sim G \text{ or } \sim G \vee A \vee B$$

$$6. L_1 \theta_1 = \neg L_2 \theta_2 \therefore \text{then } L_2 = \neg L_1 \theta_1 \theta_2^{-1}$$

$$C = (C_1 - \{L_1\}) \theta_1 \cup (C_2 - \{L_2\}) \theta_2$$

$$C - (C_1 - \{L_1\}) \theta_1 = (C_2 - \{L_2\}) \theta_2$$

$$(C - (C_1 - \{L_1\}) \theta_1) \theta_2^{-1} = C_2 - \{L_2\}$$



$$C_2 = (C - (C_1 - \{L_1\}) \theta_1) \theta_2^{-1} \cup \{L_2\}$$

$$C_2 = (C - (C_1 - \{L_1\}) \theta_1) \theta_2^{-1} \cup \{\neg L_1 \theta_1 \theta_2^{-1}\}$$

7. For every state and action, imagine an infinite number of iterations/episodes.

As the episodes approach infinity, ~~then~~ the error between \hat{Q} and the true Q -value approaches the maximum error,

Specifically, the error will continuously ~~decrease~~ decrease by a factor of γ , the discount rate.

8.) Feature Selection Methods

Wrapper - A greedy, heuristic algorithm that searches for features with a local optimum. Combines features until the local optimum declines due to overfitting

Filter - Uses ~~an~~ a learning algorithm and statistical measurements to determine the features. ~~May miss~~ Might miss important features.

Embedded - Combines the Wrapper and Filter methods. Less prone to overfitting, reduces computational costs, and gets the useful features.

9.) PCR - Discards eigenvalues that are the smallest. PCR eliminates low variance and prioritizes high variance.

PLS - Shrinks low-variance and inflates high-variance components using least squares

Ridge - Shrinks all components, but shrinks the low-variance components more.

10.) Shrinkage Methods

Ridge Regression - ~~Deterministic~~ Reduces high variance by adding a little bias. Does so by penalizing the squares of coefficients. Ultimate goal is to shrink coefficients and reduce high variance

LASSO - Similar to Ridge Regression, ~~in that~~ in that it reduces high variance by adding a little bias. However, instead of squaring the coefficients, it takes the absolute value. Like Ridge, it shrinks coefficients (and can eliminate them) and reduces high variance.