

Ques 1:

- a). Model was trained on data from stable economic conditions and tested on data from economic downturn. Training data distribution no longer matches testing distribution.
- b). High quality, diverse and representative data enables the model to capture general patterns. A simple model can generalize well with good data, while a complex model will fail with biased data.
- c). Large volume of data is normal and unlabeled. Very few labeled fraud cases exist. Clustering based anomaly detection algorithm can be used.
- d). Outcome is binary with a probability p of success. Bernoulli distribution can be used.
-

Ques 2:

- a). Initial Parameters: $w=0.1$, $b=0$, $\alpha=0.05$

Applicant Data:

Applicant	Income(x)	Default(y)
A	$2 \rightarrow x^1$	$0 \rightarrow y^1$
B	$5 \rightarrow x^2$	$1 \rightarrow y^2$

Predicted Probability:

Applicant A, $z^1 = wx^1 + b = 0.1 \times 2 + 0 = 0.2$

$$\sigma(z^1) = \frac{1}{1 + e^{-z^1}} = \frac{1}{1 + e^{-0.2}} = \frac{1}{1 + 0.8187} = 0.5498$$

Applicant B, $z^2 = wx^2 + b = 0.1 \times 5 + 0 = 0.5$

$$\sigma(z^2) = \frac{1}{1 + e^{-z^2}} = \frac{1}{1 + e^{-0.5}} = \frac{1}{1 + 0.6065} = 0.6225$$

$$\hat{y}^1 = 0.55$$

$$\hat{y}^2 = 0.62$$

Batch Gradient:

$$\begin{aligned}\frac{\partial L}{\partial w} &= \frac{-1}{m} \left(\sum_{i=1}^m y^{(i)} - \hat{y}^{(i)} \right) x^{(i)} \\ &= \frac{1}{m} \left(\sum_{i=1}^m \hat{y}^{(i)} - y^{(i)} \right) x^{(i)}\end{aligned}$$

Applicant A, $(\hat{y}^1 - y^1) x^1 = (0.5498 - 0)(2) = 1.0996$

Applicant B, $(\hat{y}^2 - y^2) x^2 = (0.6225 - 1)(5) = -1.8875$

$$\frac{\partial L}{\partial w} = \frac{1}{2} (1.0996 - 1.8875) = \frac{-0.7879}{2} = -0.394$$

$$\frac{\partial L}{\partial w} = -0.394$$

b). SVM decision boundary $f(x) = 0.8x_1 - 1.5x_2 - 50 = 0$

Substitute new Stock values:

$$\begin{aligned}f(x) &= 0.8 * 72 - 1.5 * 15 - 50 \\ &= -14.9\end{aligned}$$

Buy class support vector $f(x) = 0.8 * 75 - 1.5 * 10 - 50$
 $= -5$

Sell class support vector $f(x) = 0.8 * 60 - 1.5 * 20 - 50$
 $= -32$

Importance of Support vectors:

SV are data points closest to decision boundary. They are the critical points that constrain the optimal separating hyperplane.

Why removing non-support vectors does not affect the boundary:

Non-support vectors lie outside the margin. They do not contribute to the optimization constraints.

Ques 3:

a).

Avg Intra Distance measures how compact the cluster is.
(lower is better)

Closest Inter Distance measures how well the cluster is separated from other clusters (higher is better).

Among C_1, C_2 and C_3 , C_2 has highest intra distance and lowest inter distance. C_2 shows the strongest indication of cluster assignment ambiguity.

b). Cumulative Variance:

$$PC1 + PC2 : 58 + 30 = 88\%$$

$$PC1 + PC2 + PC3 : 58 + 30 + 9 = 97\%$$

$$PC1 + PC2 + PC3 + PC4 : 58 + 30 + 9 + 3 = 100\%$$

c).

Precision: among all predicted as +ve, how many are actually +ve.

Recall: among those which are actually +ve, how many were predicted +ve.

Security team goal is to increase Recall
and

decrease unnecessary alerts

→ decrease false positives → Increase Precision

Threshold	Precision	Recall
0.40	0.72	0.88
0.75	0.93	0.46

→ 88% Recall and 28% FP
→ 46% Recall and 7% FP

} Threshold 0.8 preferrable

Ques 4:

a). Validation accuracy represents generalization to unseen data.

Method	Batch Size	Training Accuracy	Validation Accuracy
Full Batch GD	120,000	72%	70%
Mini Batch GD	256	89%	87%
Stochastic GD	1	97%	75%

Mini Batch GD performs best.

	Validation Accuracy
No Regularization	87%
L2 Regularization	90%
Dropout ($p=0.4$)	92%

Mini Batch with dropout is best choice.

b).

$$\text{leaky ReLU}(a_i) = \begin{cases} z_i & z_i > 0 \\ 0.1 z_i & z_i \leq 0 \end{cases}$$

$$z = \begin{bmatrix} 3 \\ -4 \\ 2 \end{bmatrix}$$

$$\frac{\partial L}{\partial a} = \begin{bmatrix} 5 \\ -1 \\ 4 \end{bmatrix}$$

local gradient:

$$\frac{\partial a_i}{\partial z_i} = \begin{cases} 1 & z_i > 0 \\ 0.1 & z_i \leq 0 \end{cases}$$

$$z_1 = 3 > 0 \Rightarrow \frac{\partial a_1}{\partial z_1} = 1$$

$$z_3 = 2 > 0 \Rightarrow \frac{\partial a_3}{\partial z_3} = 1$$

$$z_2 = -4 < 0 \Rightarrow \frac{\partial a_2}{\partial z_2} = 0.1$$

$$\text{Jacobian matrix} = J_{\text{ReLU}}(z) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Downstream Gradient

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \begin{bmatrix} 5 \times 1 \\ -1 \times 0.1 \\ 4 \times 1 \end{bmatrix} = \begin{bmatrix} 5 \\ -0.1 \\ 4 \end{bmatrix}$$

c). $v(\{ \}) = 50$ $v(\{C\}) = 60$
 $v(\{E\}) = 70$ $v(\{E, C\}) = 90$

Case	Subjects without C	Contribution of C
1.	$\{ \}$	$v(\{C\}) - v(\{ \}) = 60 - 50 = 10$
2.	$\{E\}$	$v(\{E, C\}) - v(\{E\}) = 90 - 70 = 20$

$$\text{Shapley value for feature } C = \sum_{z'} \frac{|z'|! (M - |z'| - 1)!}{M!} [v(z' \cup C) - v(z')]$$

\uparrow # features in subset \downarrow Subset without C \uparrow total # feature

For case 1, $|z'| = 0$, $M = 2$

$$\frac{|z'|! (M - |z'| - 1)!}{M!} = \frac{0! (2 - 0 - 1)!}{2!} = \frac{1}{2}$$

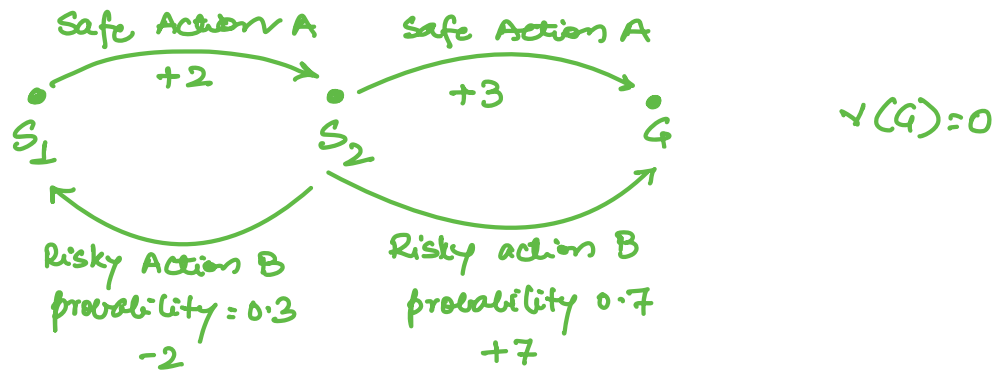
For case 2, $|z'| = 1$, $M = 2$

$$\frac{|z'|! (M - |z'| - 1)!}{M!} = \frac{1! (2 - 1 - 1)!}{2!} = \frac{1}{2}$$

$$\phi_C = \frac{1}{2} \times 10 + \frac{1}{2} \times 20 = 5 + 10 = 15$$

Ques 5:

a).



$$\gamma = 0.9$$

Policy π always chooses safe (A) in both S_1 and S_2 .

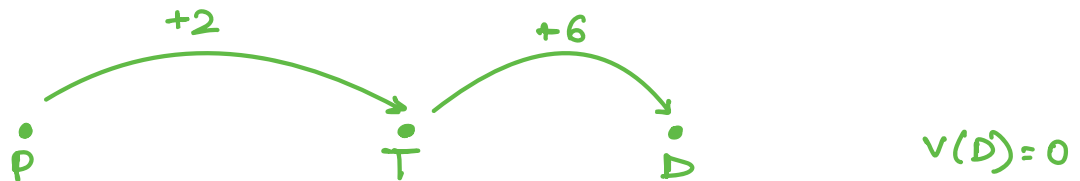
$$v^*(S_2) = 3 + \gamma v^*(G) = 3 + 0.9 * 0 = 3$$

$$v^*(S_1) = 2 + \gamma v^*(S_2) = 2 + 0.9 * 3 = 4.7$$

Risky (B) in S_2

$$\begin{aligned} Q^*(S_2, B) &= \sum_{s'} P(s' | S_2, B) [\pi(S_2, B, s') + \gamma v^*(s')] \\ &= 0.7 [7 + 0.9 * 0] + 0.3 [-2 + 0.9 * 4.7] \\ &= 0.7 * 7 + 0.3 * 2.23 = 4.9 + 0.669 \\ &= 5.569 \end{aligned}$$

b).



Discount factor $\gamma = 0.5$
Learning rate $\alpha = 0.5$

After first episode, $V_1(P) = 1$, $V_1(T) = 3$

$$\begin{aligned}
 v_2(P) &= v_1(P) + \alpha [r + \gamma v_1(T) - v_1(P)] \\
 &= 1 + 0.5 [2 + 0.5 * 3 - 1] \\
 &= 1 + 0.5 [2.5] = 1 + 1.25 = 2.25
 \end{aligned}$$

$$\begin{aligned}
 v_2(T) &= v_1(T) + \alpha [r + \gamma v_1(D) - v_1(T)] \\
 &= 3 + 0.5 [6 + 0.5 * 0 - 3] \\
 &= 3 + 0.5 [3] = 3 + 1.5 = 4.5
 \end{aligned}$$
