

Date _____						
M	T	W	T	F	S	S

ASSIGNMENT 2

Question 1:

(a) As M increases the model becomes more flexible and can fit more complex relationships. For high M , the polynomial can follow noise in the training data and produce large oscillations; this is overfitting. Overfitting occurs as high-degree polynomials have many parameters and can fit random fluctuations in the training set instead of the underlying pattern.

This problem can be addressed by including a penalty term to smoothen the curve, known as regularization. Increasing training data and using cross-validation also addresses the overfitting problem.

(b) Split the data into training and validation. Train polynomial models for different values of M . Evaluate performance on validation set. Select M that gives the best validation metric.

RC

Signature _____

No.

Date _____

M	T	W	T	F	S	S
---	---	---	---	---	---	---

(c) $\hat{y}(x_i, \omega) = \sum_{j=0}^M \omega_j x_i^j$

MSE :

$$J(\omega) = \frac{1}{2N} \sum_{i=1}^N (x_i \omega - y_i)^2 \Rightarrow \frac{1}{2N} \|x\omega - y\|^2$$

without regularization.

→ gradient descent can be used to reduce the cost function:

$$\frac{\partial J}{\partial \omega} = \frac{1}{N} x^T (x\omega - y) + \alpha \omega, \quad \omega_j = \omega_j - \alpha \frac{\partial J}{\partial \omega_j}$$

Question 2 :

- (a) Stochastic gradient descent is preferred over batch gradient descent when the dataset is large because stochastic is faster than batch gradient descent, it also improves generalization.
- (b) The learning rate (α) scales how far we move along the gradient in each update. A large value of α may lead to divergence of the algorithm. A small α value leads to extremely slow convergence. Moderate α with decay; large steps for fast progress, smaller steps later for convergence.

Page No. _____

RC

Signature _____

Date _____

M	T	W	T	F	S	S
---	---	---	---	---	---	---

- (c) The randomness in SGD stochastic gradient descent helps avoid getting stuck at saddle points and acts as regularization, however it needs tuning for stable convergence.

Q3:

$$(a) P(y=m|x) = h_{\theta_m}(x) = \frac{e^{\theta_m^T x}}{\sum_{k=0}^{m-1} e^{\theta_k^T x}}$$

$$\Rightarrow L(\theta) = - \sum_{i=1}^n \sum_{k=1}^K y_{(i,k)} \log(P_{(i,k)})$$

(b) Gradient:

$$\frac{\partial L}{\partial \theta_k} = \sum_{i=1}^n (h_{\theta}(x_i) - y_{(i,k)}) x_i$$

$$\Rightarrow \theta_k = \theta_k - \frac{\alpha \partial L}{\partial \theta_k}$$

- (c) Yes, cross entropy is exactly the negative log-likelihood used in logistic regression.

$$L_i = -[y_i \log p_i + (1-y_i) \log (1-p_i)]$$

- (d) Cross entropy handles multiple classes by comparing the true one-hot distribution to the predicted probability distribution, it uses the softmax output and checks for loss.

$$-\sum_{k=1}^K y_{(i,k)} \log(h_{\theta}(x_i)_k)$$

RC

Question 4:

$$(a) \begin{aligned} \beta_0 &= -3, -2, -1 \\ \beta_1 &= 0.4, 0.6, 0.9 \\ \beta_2 &= 0.05, 0.07, 0.1 \end{aligned} \quad \left. \begin{array}{l} \Rightarrow \text{Hours of study} = s \\ \Rightarrow \text{Previous score} = 70. \end{array} \right\}$$

$$\Rightarrow z_0 = -3 + 0.4(s) + 0.05(70) = 2.5$$

$$\Rightarrow z_1 = -2 + 0.6(s) + 0.07(70) = 5.9$$

$$\Rightarrow z_2 = -1 + 0.9(s) + 0.1(70) = 10.5$$

$$\Rightarrow e^{z_0} = e^{2.5} = 12.18$$

$$\Rightarrow e^{z_1} = e^{5.9} = 365.04.$$

$$\Rightarrow e^{z_2} = e^{10.5} = 36316.503$$

$$\Rightarrow \sum_{k=0}^{m-1} e^{0_k T_k} = 36692.72$$

(b)

$$t_0 = t_1 = 0$$

$$t_2 = 1$$

$$\Rightarrow L = -\log(0.9897) = 0.0104.$$

$$\Rightarrow P_0 = \frac{12.18}{36692.72} = 0.00033$$

$$P_1 = \frac{365.037}{36692.72} = 0.00996$$

$$P_2 = \frac{36316.503}{36692.72} = 0.9897$$