

# ARCHITECTING INTELLIGENCE

Assignment 1

Q1)

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_d x_{id}$$

$$X = \begin{bmatrix} 1 & x_{i1} & x_{i2} & \dots & x_{id} \\ 1 & x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$$

$$\hat{y} = X\beta. \quad J(\beta) = \frac{1}{2} \|y - X\beta\|^2$$

$$J(\beta) = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

normal eq: soln obtained by setting  $\frac{\partial J(\beta)}{\partial \beta} = 0$ ,

(a) what regression is doing:

→ gives us a hypothesis function to predict the value of the dependent variable  $y$  with regard to the independent variable  $x$ , while using parameters  $\beta$ .

why to minimize squared error:

→ squared error represents the difference between the predicted value of  $y$  and the actual value of  $y$ . minimizing this makes our model as accurate and as close to reality as possible.

$$(b) \cdot J(\beta) = \frac{1}{2} \|y - X\beta\|^2 = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\frac{\partial J}{\partial \beta} = \sum_{i=1}^n \left[ 0 - \frac{\partial \hat{y}_i}{\partial \beta} \right] \quad y_n x_1 \\ y^T x_n$$

$$J(\beta) = \frac{1}{2} (y - X\beta)^T \cdot (y - X\beta)$$

$$= y^T y - y^T X\beta - (X\beta)^T y + (X\beta)^T X\beta$$

$$= y^T y - y^T X\beta - \beta^T x^T y + \underbrace{(X\beta)^T X\beta}_{\text{all of these terms become scalars (w dimension } 1 \times 1)}$$

$$[y_n x_1, y^T y \equiv 1 \times 1; X_{n \times (n+1)}, \beta_{(n+1) \times 1}] \quad \text{terms become scalars (w dimension } 1 \times 1)$$

$$J(\beta) = y^T y - 2 y^T X\beta + \beta^T X^T X\beta$$

$$\frac{\partial J}{\partial \beta} = 2 y^T X \quad \nabla_{\beta} J = \begin{bmatrix} \frac{\partial J}{\partial \beta_1} \\ \vdots \\ \frac{\partial J}{\partial \beta_p} \end{bmatrix}$$

$$\text{OR} \quad J(\beta) = y^T y - 2 \beta^T X^T y + \beta^T X^T X\beta$$

$$\nabla_{\beta} J = \frac{1}{2} [-2 X^T y + 2 X^T X\beta]$$

$$= X^T X\beta - X^T y = 0$$

$$\therefore \beta = (X^T X)^{-1} X^T y$$

(c) directly inverting  $X^T X$  in the normal equations:

$$\beta = (X^T X)^{-1} X^T y.$$

is often problematic in high dimensional or large scale settings for both numerical and computational reasons.

Following are the issues:

(i) computational cost (time & memory):

- forming  $(X^T X)$ :  $O(n p^2)$ .
- inverting  $(X^T X)$ :  $O(p^3)$

$\Rightarrow$  in contrast, one gradient descent step costs only  $O(np)$ .

(ii) poor scalability to large datasets:

Normal equations require all data at once;  
- it must load the entire  $X$  into memory.

$\Rightarrow$  iterative methods allow mini-batch or stochastic gradient descent.

Q.2) (a) core idea behind backpropagation:

⇒ to efficiently compute how much each parameter (weight and bias) is responsible for the final prediction error, and then use this information to update the parameters so that the error is reduced.

Role of chain rules:

The neural network consists of:

- input transferred layer by layer through weighted sums.
  - activation functions until final output is produced.
  - loss function measuring how far the output is from the true target.
- ⇒ loss depends on output, output depends on earlier layers, which in turn depend upon weights and biases.
- ⇒ the chain rule makes it possible to calculate these dependencies systematically

$$(b) \quad z_1 = w_1 a_0 + b_1; \quad a_1 = \sigma(z_1). \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

$$z_2 = w_2 a_1 + b_2; \quad a_2 = \sigma(z_2) = \hat{y}$$

$$L = -[y \log a_2 + (1-y) \log (1-a_2)]$$

$$L = -[y \log a_2 + (1-y) \log(1-a_2)]$$

$$\bullet \frac{\partial L}{\partial w_2} = - \left[ \frac{y}{a_2} \cdot \frac{\partial a_2}{\partial w_2} + \frac{(1-y)}{(1-a_2)} \cdot \frac{\partial a_2}{\partial w_2} \right]$$

$$a_2 = \sigma(z_2) = \sigma(w_2 a_1 + b_2)$$

$$\sigma'(z) = -\frac{1}{(1+e^{-z})^2} \left[ -e^{-z} \right] = \frac{e^{-z}}{(1+e^{-z})^2} \cdot dz$$

$$\frac{\partial a_2}{\partial w_2} = \frac{\partial \sigma(z_2)}{\partial w_2} = \frac{\partial \sigma(z_2)}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} = a_1 \cdot \frac{\partial \sigma(z_2)}{\partial z_2} = \frac{\partial a_2}{\partial w_2}$$

$$\begin{aligned} \frac{\partial a_2}{\partial w_2} &= a_1 \left( \frac{1}{a_2^2} - \frac{1}{a_2^3} \right) \\ &= a_1 \frac{(a_2^2 - a_2^3)}{a_2^3} = a_1 \left( \frac{1 - a_2}{a_2^3} \right) \end{aligned}$$

$$\frac{\partial L}{\partial w_2} = +\frac{y a_1 (1-a_2)}{a_2^3} + \frac{(1-y)}{a_2^3} \cdot a_1$$

$$\bullet \frac{\partial L}{\partial b_2} = - \left[ \frac{y}{a_2} \cdot \frac{\partial a_2}{\partial b_2} + \frac{(1-y)}{(1-a_2)} \cdot \frac{\partial a_2}{\partial b_2} \right]$$

$$\frac{\partial a_2}{\partial b_2} = \frac{\partial \sigma(z_2)}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} = \frac{(1-a_2)}{a_2^3}$$

$$\frac{\partial L}{\partial b_2} = -\frac{y}{a_2^3} (1-a_2) + \frac{(1-y)}{a_2^3}$$

- ~~$\frac{\partial L}{\partial w_i}$~~  (2). for deriving  $\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial b_1}$ , we use similar methods.

(c). in gradient descent, for a parameter  $A$ ,

$$A_n = A_{n-1} - \eta \cdot \frac{\partial L}{\partial A} \quad L: \text{loss func.}$$

$\eta = \text{learning rate}$

$$w_1 \rightarrow w_1 - \eta \frac{\partial L}{\partial w_1}$$

$$w_2 \rightarrow w_2 - \eta \frac{\partial L}{\partial w_2}$$

$$b_1 \rightarrow b_1 - \eta \frac{\partial L}{\partial b_1}$$

$$b_2 \rightarrow b_2 - \eta \frac{\partial L}{\partial b_2}$$

$\Rightarrow$  role of learning rate  $\eta$ :

controls how big each step is.

- if  $\eta$  is too large: updates may overshoot. loss might oscillate or diverge.
- if  $\eta$  is too small: training remains very slow.

$\Rightarrow$  we need a proper learning rate for a steady decrease in loss and a stable convergence

- (Q3) (b) Simple RNNs struggle w. long term dependencies because gradients become very small or very large when propagated over many time steps.
- (c) LSTM gates control what info is stored, updated or removed from the cell state.
- (d) LSTMs address the vanishing gradient problem by maintaining a nearly constant error flow through the cell state across time steps.
- (e) ANNs are best for tasks like image classification, RNNs for short sequence prediction and LSTMs for long sequence tasks such as language translation.
- (Q4) (a) ~~The book the~~ In a sentence like "the book that the professor who the students admired wrote was famous." understanding that 'was' refers to the 'book' depends on long-range dependency across many words, and a standard RNN would struggle to model this relatively reliably due to vanishing gradients over long sequences.

(b) An LSTM's memory cell acts like a conveyor belt that carries important info forward, while gates what to keep, update or discard, while allowing relevant context to persist over long sequences. In machine translation, the forget gate will need to be close to 0 when moving from one sentence or class to another so that outdated context from the previous class does not interfere with translating the next one correctly.

Q3) (a) An ANN processes each input independently with no memory, while RNN processes sequences by passing info from one step to the next through a hidden state.