

1. Multiple Choice Questions:

Lasso: Question A:

True

Bagging & Bootstrap Sampling

Question B:

Option B.

Convolutional Filters

Question C:

Fig 2 (a) : K_2

Fig 2 (b) : K_1

Fig 2 (c) : K_2

Multiclass SVMs:

Question D:

$\hat{y} = z$ is the maximizer of $||\cdot||$

Feature Maps:

Question E:

Option A

Hard-Margin SVMs:

Question F:

False

Question G:

Option A.

AdaBoost:

Question H:

False.

Tensor Model Training

Question I:

True

Bias - Variance Decomposition:

Question J:

False

HMM EM Learning:

Question K:

Option A

Non-Negative Matrix Factorization:

Question L:

True

Decision Tree:

Question M:

False

Overfitting:

Question N:

True

2. Naïve Bayes:

Question 1:

using maximum likelihood estimate formula:

$$\begin{aligned} \text{eg: } P(\text{Grade} = A \mid \text{Happy?} = \text{Yes}) &= \frac{1 + \sum_{(x,y)} 1_{[x=\text{Grade}=A \wedge y=\text{Yes}]} }{2 + \sum_{(x,y)} 1_{[y=\text{Yes}]} } \\ &= \frac{1+3}{2+4} = \frac{2}{3} \end{aligned}$$

we can calculate:

$$P(\text{Grade} = A \mid \text{Happy?} = \text{No}) = \frac{1}{3}$$

$$P(\text{Grade} = C \mid \text{Happy?} = \text{Yes}) = \frac{1}{3}$$

$$P(\text{Grade} = C \mid \text{Happy?} = \text{No}) = \frac{2}{3}$$

therefore, final probability table for:

(table 1): $P(\text{Grade} \mid \text{Happy?})$	Grade = A	Grade = C
Happy? = Yes	$\frac{2}{3}$	$\frac{1}{3}$
Happy? = No	$\frac{1}{3}$	$\frac{2}{3}$

Similarly, we can calculate:

$$P(\text{Year} = \text{Freshman} \mid \text{Happy?} = \text{Yes}) = \frac{1}{3}$$

$$P(\text{Year} = \text{Freshman} \mid \text{Happy?} = \text{No}) = \frac{1}{2}$$

$$P(\text{Year} = \text{Senior} \mid \text{Happy?} = \text{Yes}) = \frac{2}{3}$$

$$P(\text{Year} = \text{Senior} \mid \text{Happy?} = \text{No}) = \frac{1}{2}$$

therefore, final probability table for:

(table 2): $P(\text{Year} \mid \text{Happy?})$	Year = Freshman	Year = Senior
Happy? = Yes	$\frac{1}{3}$	$\frac{2}{3}$
Happy? = No	$\frac{1}{2}$	$\frac{1}{2}$

Question 2:

$$P(\text{Year} = \text{Freshman}, \text{Grade} = C, \text{Happy?} = \text{No})$$

$$= P(\text{Happy?} = \text{No}) P(\text{Grade} = C \mid \text{Happy?} = \text{No}) P(\text{Year} = \text{Freshman} \mid \text{Happy?} = \text{No})$$

$$= \frac{1}{2} \cdot \frac{2}{3} \cdot \frac{1}{2} = \frac{1}{6}$$

Question 3:

Pseudocode for drawing a sample.

// first sample y (Happy?)

$R_y = \text{random}()$

Set Happy? = Yes if $R_y < P(\text{Happy?} = \text{Yes})$

else Set Happy? = No

// then sample each x^d (Grade & Year)

$R_{x_1} = \text{random}()$

Set Grade = A if $R_{x_1} < P(\text{Grade} = A | y)$

else Set Grade = C

$R_{x_2} = \text{random}()$

Set Year = Freshman if $R_{x_2} < P(\text{Year} = \text{Freshman} | y)$

else Set Year = Senior

3. Data Transformation:

Question 1:

We know:

$$\tilde{x} = Ax$$

$$\omega^T x = \tilde{\omega}^T \tilde{x} = \tilde{\omega}^T Ax$$

therefore,

$$\omega^T = \tilde{\omega}^T A$$

thus,

$$\omega = A^T \tilde{\omega} = A \tilde{\omega}$$

Question 2:

$$\operatorname{argmin}_{\tilde{\omega}} \frac{\lambda}{2} \|\tilde{\omega}\|^2 + \sum_i (y_i - \tilde{\omega}^T \tilde{x}_i)^2$$

$$= \operatorname{argmin}_{\omega} \frac{\lambda}{2} \|A^T \omega\|^2 + \sum_i (y_i - \omega^T x_i)^2$$

Question 3:

when x is rescaled by A ,

the regulation sector in loss function is changed from $\frac{\lambda}{2} \|\omega\|^2$ to $\frac{\lambda}{2} \|A^T \omega\|^2$, then the optimal ω would be different

4. Latent Markov Embedding

Question 1:

as U, V, X are optimal choices which maximize $P(s)$,
single-data model is more restricted to dual-point model
because it requires that " U " and " V " equals to each other,

$$P(s) = \prod_{p \in s} \prod_{i=1}^{M_p} \frac{e^{-\|U(p^{i1}) - V(p^{i1-1})\|_2^2}}{Z(p^{i1-1})} \quad (17)$$

$$P(s) = \prod_{p \in s} \prod_{i=1}^{M_p} \frac{e^{-\|X(p^{i1}) - X(p^{i1-1})\|_2^2}}{Z(p^{i1-1})} \quad (18)$$

therefore, $P(s)_{\text{dual}} \geq P(s)_{\text{single}}$

Question 2:

if (16) is equal to (18) for every pair of songs s and s' .

this case is satisfied when the optimal dual point model returns a solution U, V where both are equal to each other

$$U = V.$$

if this is the case, then the dual-point model changes into single-point model, as U, V, X are all optimal choices already, what we can interpret is that.

$$U = V = X.$$

5. Neural Net Backprop Gradient Derivation

Question 1:

$$\frac{\partial}{\partial w_{11}} L(y, f(x)) = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{11}}$$

we expand each, $\frac{\partial L}{\partial f} = -2(y - f)$

$$\frac{\partial f}{\partial h_1} = \frac{\partial \sigma(u_1 h_1 + u_2 h_2)}{\partial h_1} = \sigma(u_1 h_1 + u_2 h_2) (1 - \sigma(u_1 h_1 + u_2 h_2)) u_1$$

$$\frac{\partial h_1}{\partial w_{11}} = \frac{\partial \sigma(w_{11} x_1 + w_{21} x_2)}{\partial w_{11}} = \sigma(w_{11} x_1 + w_{21} x_2) (1 - \sigma(w_{11} x_1 + w_{21} x_2)) x_1$$

Question 2:

using formula in Question 1 and data provided.

we have,

$$\begin{aligned} \frac{\partial}{\partial w_{11}} L(y - f(x)) &= \frac{\partial}{\partial w_{11}} (y - f(x))^2 \\ &= \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{11}} \end{aligned}$$

$$w_{11} x_1 + w_{21} x_2 = 0.25 \times 0.1 + 0.05 \times 0.5 = 0.05$$

$$\frac{\partial h_1}{\partial w_{11}} = \sigma(0.05) (1 - \sigma(0.05)) \times 0.1 = 0.025$$

$$\begin{aligned} u_1 h_1 + u_2 h_2 &= u_1 (\sigma(w_{11} x_1 + w_{21} x_2)) + u_2 (\sigma(w_{12} x_1 + w_{22} x_2)) \\ &= 0.5 \times 0.5125 + (-0.1) \times 0.4717 = 0.2091 \end{aligned}$$

$$\frac{\partial f}{\partial h_1} = \sigma(0.2091) (1 - \sigma(0.2091)) \times 0.5 = 0.1236$$

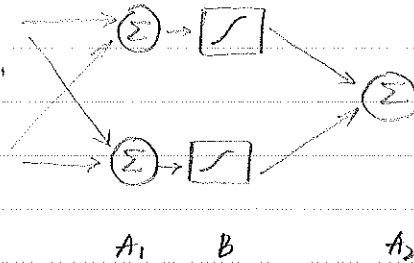
$$f = \sigma(u_1 h_1 + u_2 h_2) = \sigma(0.2091) = 0.5121$$

$$\frac{\partial L}{\partial f} = -2(y - f) = -2(0.25 - 0.5121) = -0.3958$$

$$\text{therefore, } \frac{\partial}{\partial w_{11}} L(y - f(x)) = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{11}} = (-0.3958) \cdot 0.1236 \cdot 0.025 = -0.0012$$

Question 3:

$\frac{\partial h_1}{\partial w_{11}}$ that is the ...



$\frac{\partial h_1}{\partial w_{11}}$ is always between 0 and 1, and when neural networks with more layers, they multiply together and become really small and gradient descent process becomes very inefficient and results in vanishing gradient problem. More Layers means more value $\in (0, 1)$ to be multiplied, so the problem is exacerbated exponentially with layer number.

we can solve the problem by changing the sigmoid function to Rectilinear function.