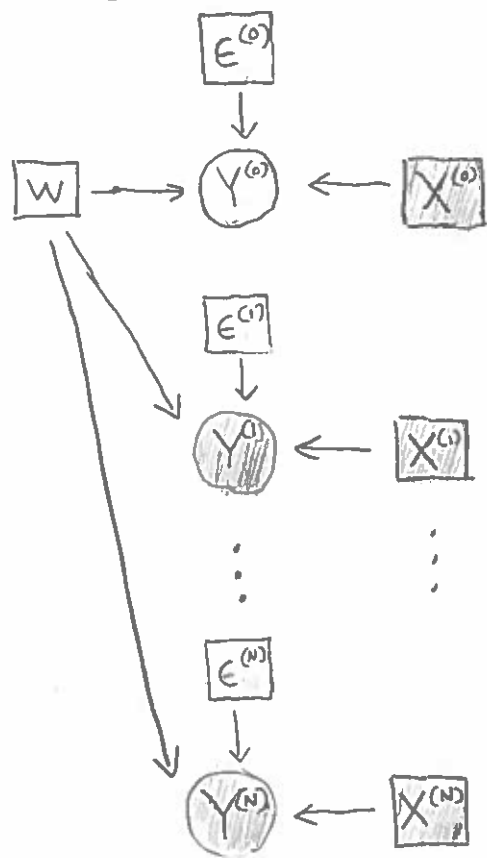


LOGISTIC REGRESSION: MLE

① Recall "logistic regression":



where: $P(\epsilon^{(n)}) \sim \text{Constant}(0, 1) \quad \forall n \in \{0, \dots, N\}$

$$y^{(n)} \leftarrow 1_{\epsilon^{(n)} < (1 + \exp(-w^T x^{(n)}))^{-1}}$$

② Also recall that one way to estimate the value of the unobserved response variable $Y^{(0)}$ is through maximum likelihood estimation (MLE):

(a) compute $\hat{w} = \underset{w}{\operatorname{argmax}} \prod_{n=1}^N P(y^{(n)} | w, x^{(n)})$

(b) compute $\hat{y}^{(0)} = \underset{y^{(0)}}{\operatorname{argmax}} P(y^{(0)} | \hat{w}, x^{(0)})$

LOGISTIC REGRESSION: MLE

③ To compute the second step, observe:

$$P(Y^{(n)} = 1 | w, x^{(n)})$$

$$= \int_0^1 P(Y^{(n)} = 1, \epsilon^{(n)} | w, x^{(n)}) d\epsilon^{(n)} \quad [\text{total probability}]$$

$$= \int_0^1 P(Y^{(n)} = 1 | w, x^{(n)}, \epsilon^{(n)}) P(\epsilon^{(n)} | w, x^{(n)}) d\epsilon^{(n)} \quad [\text{Chain Rule}]$$

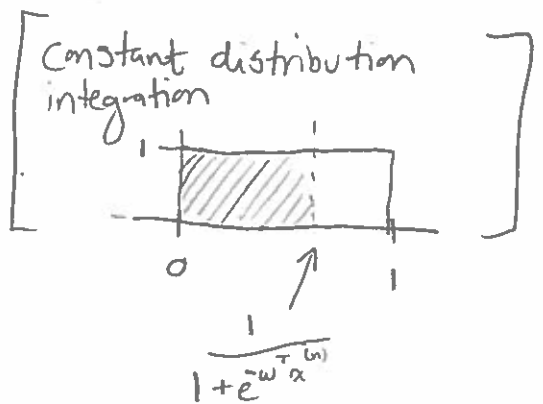
$$= \int_0^1 P(Y^{(n)} = 1 | w, x^{(n)}, \epsilon^{(n)}) P(\epsilon^{(n)}) d\epsilon^{(n)} \quad [\text{d-separation}]$$

$$= \int_0^1 P\left(\epsilon^{(n)} < \frac{1}{1 + e^{-w^T x^{(n)}}}\right) P(\epsilon^{(n)}) d\epsilon^{(n)} \quad \left[\begin{array}{l} \text{b/c} \\ y^{(n)} \leftarrow 1_{\epsilon^{(n)} < (1 + \exp(-w^T x^{(n)}))^{-1}} \end{array} \right]$$

$$= \int_0^{\frac{1}{1 + e^{-w^T x^{(n)}}}} P(\epsilon^{(n)}) d\epsilon^{(n)}$$

$$\left[\begin{array}{l} \text{everywhere else,} \\ P(\epsilon^{(n)} < \frac{1}{1 + e^{-w^T x^{(n)}}}) = 0 \end{array} \right]$$

$$= \frac{1}{1 + e^{-w^T x^{(n)}}}$$



Thus:

$$\begin{aligned} P(Y^{(n)} = 0 | w, x^{(n)}) &= 1 - \frac{1}{1 + e^{-w^T x^{(n)}}} = \frac{1 + e^{-w^T x^{(n)}}}{1 + e^{-w^T x^{(n)}}} - \frac{1}{1 + e^{-w^T x^{(n)}}} \\ &= \frac{e^{-w^T x^{(n)}}}{1 + e^{-w^T x^{(n)}}} \end{aligned}$$

LOGISTIC REGRESSION: MLE

④ Or, more compactly:

$$P(y^{(n)} | w, x^{(n)}) = \frac{e^{-(1-y^{(n)})w^T x^{(n)}}}{1 + e^{-w^T x^{(n)}}}$$
$$= \begin{cases} \frac{e^{-w^T x^{(n)}}}{1 + e^{-w^T x^{(n)}}} & \text{if } y^{(n)} = 0 \\ \frac{1}{1 + e^{w^T x^{(n)}}} & \text{if } y^{(n)} = 1 \end{cases}$$

⑤ This allows us to express ②(b) as:

$$\hat{y}^{(0)} = \underset{y^{(0)}}{\operatorname{argmax}} P(y^{(0)} | w, x^{(0)})$$

$$= \underset{y^{(0)} \in \{0,1\}}{\operatorname{argmax}} \frac{e^{-(1-y^{(0)})w^T x^{(0)}}}{1 + e^{-w^T x^{(0)}}}$$

LOGISTIC REGRESSION: MLE

⑥ To compute 2(a), we start with some simplifications:

$$\hat{w} = \arg\max_w \prod_{n=1}^N P(y^{(n)} | w, x^{(n)})$$

$$= \arg\max_w \log \prod_{n=1}^N P(y^{(n)} | w, x^{(n)})$$

$$= \arg\max_w \sum_{n=1}^N \log P(y^{(n)} | w, x^{(n)})$$

$$= \arg\max_w \sum_{n=1}^N \log \frac{e^{-(1-y^{(n)})w^T x^{(n)}}}{1 + e^{-w^T x^{(n)}}} \quad [\text{from } ④]$$

$$= \arg\max_w \sum_{n=1}^N \log e^{-(1-y^{(n)})w^T x^{(n)}} - \log (1 + e^{-w^T x^{(n)}})$$

$$= \arg\max_w \sum_{n=1}^N -(1-y^{(n)})w^T x^{(n)} + \log \frac{1}{1 + e^{-w^T x^{(n)}}}$$

$$= \arg\min_w \sum_{n=1}^N (1-y^{(n)})w^T x^{(n)} - \log \frac{1}{1 + e^{-w^T x^{(n)}}}$$

So for logistic regression, our loss function is:

$$L_{\text{logistic}}(w) = \sum_{n=1}^N (1-y^{(n)})w^T x^{(n)} - \log \frac{1}{1 + e^{-w^T x^{(n)}}}$$

LOGISTIC REGRESSION: MLE

⑦ To compute the gradient of $L_{\text{logistic}}(w)$, we'll first prove the following lemma:

Lemma: If $\sigma(a) = \frac{1}{1+e^{-a}}$, then:

$$\frac{d}{da} \sigma(a) = \sigma(a)(1 - \sigma(a))$$

Proof:

$$\begin{aligned} \frac{d}{da} \sigma(a) &= \frac{-1}{(1+e^{-a})^2} \cdot e^{-a} \cdot -1 \\ &= \frac{e^{-a}}{(1+e^{-a})^2} \\ &= \left(\frac{1}{1+e^{-a}} \right) \left(\frac{e^{-a}}{1+e^{-a}} \right) \\ &= \sigma(a) \left(\frac{1+e^{-a}-1}{1+e^{-a}} \right) \\ &= \sigma(a) \left(\frac{1+e^{-a}}{1+e^{-a}} - \frac{1}{1+e^{-a}} \right) \\ &= \sigma(a)(1 - \sigma(a)) \quad \square \end{aligned}$$

LOGISTIC REGRESSION: MLE

⑧ So the gradient of the loss function is:

$$\frac{d}{dw} L_{\text{logistic}}(w)$$

$$= \sum_{n=1}^N \frac{d}{dw} (1 - y^{(n)}) w^T x^{(n)} - \frac{d}{dw} \log \sigma(w^T x^{(n)})$$

$$= \sum_{n=1}^N (1 - y^{(n)}) \frac{d}{dw} w^T x^{(n)} - \frac{1}{\sigma(w^T x^{(n)})} \frac{d}{dw} \sigma(w^T x^{(n)})$$

$$= \sum_{n=1}^N (1 - y^{(n)}) \frac{d}{dw} (w^T x^{(n)}) - \frac{1}{\sigma(w^T x^{(n)})} \sigma(w^T x^{(n)}) (1 - \sigma(w^T x^{(n)})) \frac{d}{dw} w^T x^{(n)}$$

[from Lemma]

$$= \sum_{n=1}^N (1 - y^{(n)}) \frac{d}{dw} (w^T x^{(n)}) - (1 - \sigma(w^T x^{(n)})) \frac{d}{dw} (w^T x^{(n)})$$

$$= \sum_{n=1}^N (1 - y^{(n)}) x^{(n)} - (1 - \sigma(w^T x^{(n)})) x^{(n)}$$

[b/c $\frac{d}{dw} w^T x = x$]

$$= \sum_{n=1}^N (1 - y^{(n)} - 1 + \sigma(w^T x^{(n)})) x^{(n)}$$

$$= \sum_{n=1}^N (\sigma(w^T x^{(n)}) - y^{(n)}) x^{(n)}$$

LOGISTIC REGRESSION: MLE

⑨ This can be expressed even more compactly in terms of the evidence matrix X and response vector y :

$$\frac{d}{dw} L_{\text{logistic}}(w) = \sum_{n=1}^N (\sigma(w^T x^{(n)}) - y^{(n)}) x^{(n)}$$

$$= X^T (\sigma(Xw) - y)$$

Exercise: Show $X^T (\sigma(Xw) - y) = \sum_{n=1}^N (\sigma(w^T x^{(n)}) - y^{(n)}) x^{(n)}$

LOGISTIC REGRESSION: MLE

- ⑩ As usual, there isn't a known way to solve directly for $\frac{d}{dw} L_{\text{logistic}}(w) = 0$, however we are free to use gradient descent.

LOGISTIC REGRESSION $(X, y, x^{(0)})$:

- compute point estimate $\hat{w} = \text{GRADDESCENT}(L_{\text{logistic}})$
- compute prediction $\hat{y}^{(0)} = \underset{y^{(0)} \in \{0,1\}}{\operatorname{argmax}} \frac{e^{-(1-y^{(0)})w^T x^{(0)}}}{1 + e^{-w^T x^{(0)}}}$
- return $\hat{y}^{(0)}$