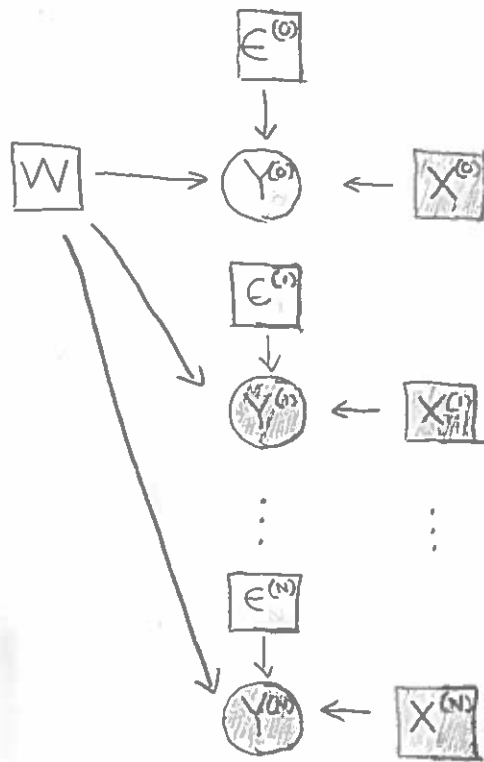


LINEAR REGRESSION: MLE

① Recall "ordinary linear regression":



where: $P_{\epsilon}(\epsilon^{(n)}) \sim \text{Normal}(0, \sigma^2) \quad \forall n \in \{0, \dots, N\}$
 $y^{(n)} \leftarrow w^T x^{(n)} + \epsilon^{(n)}$

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

LINEAR REGRESSION: MLE

③ The second step is not too bad:

$$\begin{aligned} P(y^{(0)} | \hat{w}, x^{(0)}) &= \int P(y^{(0)}, \epsilon^{(0)} | \hat{w}, x^{(0)}) d\epsilon^{(0)} && \text{[Total Probability]} \\ &= \int P(\epsilon^{(0)} | \hat{w}, x^{(0)}) P(y^{(0)} | \epsilon^{(0)}, \hat{w}, x^{(0)}) d\epsilon^{(0)} \\ &= \int P(\epsilon^{(0)}) P(y^{(0)} | \epsilon^{(0)}, \hat{w}, x^{(0)}) d\epsilon^{(0)} && \text{[Chain Rule]} \\ & && \text{[d-sep.]} \\ &= P(\epsilon^{(0)} = y^{(0)} - \hat{w}^T x^{(0)}) \\ &= P_{\epsilon}(y^{(0)} - \hat{w}^T x^{(0)}) \end{aligned}$$

[this is the only value of $\epsilon^{(0)}$ s.t. $P(y^{(0)} | \epsilon^{(0)}, \hat{w}, x^{(0)}) \neq 0$]

Therefore:

$$\begin{aligned} \hat{y}^{(0)} &= \operatorname{argmax}_{y^{(0)}} P(y^{(0)} | \hat{w}, x^{(0)}) \\ &= \operatorname{argmax}_{y^{(0)}} P_{\epsilon}(y^{(0)} - \hat{w}^T x^{(0)}) \end{aligned}$$

Since $P_{\epsilon} \sim \text{Normal}(0, \sigma^2)$, therefore $P_{\epsilon}(y^{(0)} - \hat{w}^T x^{(0)})$ is maximized when $y^{(0)} - \hat{w}^T x^{(0)} = 0$, thus:

$$\boxed{\hat{y}^{(0)} = \hat{w}^T x^{(0)}}$$

LINEAR REGRESSION: MLE

- ④ How do we compute the first step? First let's turn those annoying products into friendly sums:

$$\begin{aligned} & \operatorname{argmax}_w \prod_{n=1}^N P(y^{(n)} | w, x^{(n)}) \\ &= \operatorname{argmax}_w \log \prod_{n=1}^N P(y^{(n)} | w, x^{(n)}) \\ &= \operatorname{argmax}_w \underbrace{\sum_{n=1}^N \log P(y^{(n)} | w, x^{(n)})}_{\text{let's call this } \ell(w)} \end{aligned}$$

- ⑤ Next, let's do some manipulations of $\ell(w)$.

$$\begin{aligned} \ell(w) &= \sum_{n=1}^N \log P(y^{(n)} | w, x^{(n)}) \\ &= \sum_{n=1}^N \log P_e(y^{(n)} - w^T x^{(n)}) \quad [\text{from } \textcircled{3}] \\ &= \sum_{n=1}^N \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left(\frac{-1}{2\sigma^2} (y^{(n)} - w^T x^{(n)})^2 \right) \right] \\ &= \sum_{n=1}^N \log \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} + \log \exp \left(\frac{-1}{2\sigma^2} (y^{(n)} - w^T x^{(n)})^2 \right) \\ &= \sum_{n=1}^N \frac{-1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y^{(n)} - w^T x^{(n)})^2 \\ &= \left(\sum_{n=1}^N \frac{-1}{2} \log 2\pi\sigma^2 \right) - \left(\sum_{n=1}^N \frac{1}{2\sigma^2} (y^{(n)} - w^T x^{(n)})^2 \right) \\ &= \frac{-N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2 \end{aligned}$$

LINEAR REGRESSION: MLE

⑥ Thus:

$$\begin{aligned}\arg\max_w l(w) &= \arg\max_w \left(-\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2 \right) \\ &= \arg\max_w -\frac{1}{2\sigma^2} \sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2 \\ &= \arg\max_w -\sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2\end{aligned}$$

⑦ We can express $\sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2$ without the explicit summation by resorting to vector dot product:

$$\sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2 = \begin{bmatrix} y^{(1)} - w^T x^{(1)} \\ \vdots \\ y^{(N)} - w^T x^{(N)} \end{bmatrix}^T \begin{bmatrix} y^{(1)} - w^T x^{(1)} \\ \vdots \\ y^{(N)} - w^T x^{(N)} \end{bmatrix}$$

and noticing that:

$$\begin{aligned}\begin{bmatrix} y^{(1)} - w^T x^{(1)} \\ \vdots \\ y^{(N)} - w^T x^{(N)} \end{bmatrix} &= \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} - \begin{bmatrix} w^T x^{(1)} \\ \vdots \\ w^T x^{(N)} \end{bmatrix} \\ &= y - \begin{bmatrix} w_1 x_1^{(1)} + \dots + w_D x_D^{(1)} \\ \vdots \\ w_1 x_1^{(N)} + \dots + w_D x_D^{(N)} \end{bmatrix} \\ &= y - \begin{bmatrix} x_1^{(1)} & \dots & x_D^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_D \end{bmatrix} \\ &= y - Xw\end{aligned}$$

LINEAR REGRESSION: MLE

$$\begin{aligned}\textcircled{8} \text{ So } \operatorname{argmax}_w \ell(w) &= \operatorname{argmax}_w -(y - Xw)^T (y - Xw) \\ &= \operatorname{argmax}_w -\cancel{y^T y} + y^T Xw + (Xw)^T y - (Xw)^T Xw \\ &= \operatorname{argmax}_w ((Xw)^T y)^T + (w^T X^T) y - (w^T X^T) Xw \\ &\quad \quad \quad [\text{since } (AB)^T = B^T A^T] \\ &= \operatorname{argmax}_w (w^T X^T y)^T + w^T X^T y - w^T X^T Xw \\ &\quad \quad \quad [\text{since } (AB)^T = B^T A^T]\end{aligned}$$

Notice that $w^T X^T y$ is a 1×1 matrix (i.e. $(1 \times D) \cdot (D \times N) \cdot (N \times 1)$), so $(w^T X^T y)^T = w^T X^T y$. That gives us:

$$\operatorname{argmax}_w \ell(w) = \operatorname{argmax}_w 2w^T X^T y - w^T X^T Xw$$

$\textcircled{9}$ At this point, we're pretty close. We've shown (over $\textcircled{4}-\textcircled{8}$) that the point estimate \hat{w} of our weight vector is:

$$\begin{aligned}\hat{w} &= \operatorname{argmax}_w \prod_{n=1}^N P(y^{(n)} | w, x^{(n)}) \\ &= \operatorname{argmax}_w 2w^T X^T y - w^T X^T Xw \\ &= \operatorname{argmin}_w w^T X^T Xw - 2w^T X^T y\end{aligned}$$

So the loss function for ordinary linear regression is:

$$L_{\text{lin}}(w) = w^T X^T Xw - 2w^T X^T y$$

LINEAR REGRESSION: MLE

⑩ We can use the identities:

$$\boxed{\frac{\partial}{\partial a} a^T b = \frac{\partial}{\partial a} b^T a = b} \quad \text{and} \quad \boxed{\frac{\partial}{\partial a} a^T X a = (X + X^T) a}$$

to compute the gradient of $L_{\text{lin}}(w)$:

$$\begin{aligned} & \frac{\partial}{\partial w} (-2w^T X^T y + w^T X^T X w) \\ &= \frac{\partial}{\partial w} -2w^T X^T y + \frac{\partial}{\partial w} w^T X^T X w \\ &= -2X^T y + \frac{\partial}{\partial w} w^T (X^T X) w \quad \left[\text{b/c } \frac{\partial}{\partial a} a^T b = b \right] \\ &= -2X^T y + (X^T X + (X^T X)^T) w \quad \left[\text{b/c } \frac{\partial}{\partial a} a^T X a = (X + X^T) a \right] \\ &= -2X^T y + 2X^T X w \quad \left[\text{b/c } (AB)^T = B^T A^T \right] \end{aligned}$$

⑪ We can then compute $\arg\max_w \prod_{n=1}^N P(y_n | w, x_n)$ by finding when the gradient equals zero:

$$\begin{aligned} & -2X^T y + 2X^T X w = 0 \\ \Rightarrow & X^T X w = X^T y \\ \Rightarrow & w = (X^T X)^{-1} X^T y \end{aligned}$$

So we have our answer:

$$\boxed{\arg\min_w L_{\text{lin}}(w) = (X^T X)^{-1} X^T y}$$

LINEAR REGRESSION: MLE

⑫ So we can now go back to ② and make our MLE algorithm more concrete:

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

(b) compute $\hat{y}_0 = \underset{y_0}{\operatorname{argmax}} P(y_0 | \hat{w}, x_0)$

i.e.

$$\hat{w} \leftarrow (X^T X)^{-1} X^T y$$

$$\hat{y}_0 \leftarrow \hat{w}^T x_0$$