09/09/2024

Lecture: Linear Regression & Optimization

Pg 6.  $\mathcal{L}(y, t) = \frac{1}{2}(y-t)^2$

Cost $\mathcal{J}(w,b) =$ average of the sum of all $\mathcal{L}$'s.

loss calculated by = $\frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_i = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2}(y^i - t^i)^2$   $j =$ index of the feature
(regarding one sample)

the ground truth and
the prediction of the $= \frac{1}{2N} \sum_{i}^{N} (w^T x^i + b - t^i)^2$
$i^{th}$ vector

ground truth/target of the

$= \frac{1}{2N} \sum_{i}^{N} (\sum_{j}^{D} w_j x_j^i + b - t^i)^2$   $i$ th vector

Weight of the $j$th feature
within a vector/sample

Value of the $j$th feature in the $i^{th}$ sample/
vector in the training set

Pg 10.

$$X = \begin{bmatrix} 1 & [x^{(1)}]^T \\ 1 & [x^{(2)}]^T \\ \vdots & \vdots \\ 1 & [x^{(N)}]^T \end{bmatrix} \text{ and } w = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w^N \end{bmatrix} \in \mathbb{R}^{D+1}$$

why we do this: Incorporate the bias $b$ into the manipulation of matrix of all samples by essentially creating a new "dummy"/"placeholder" feature consisting of $1$'s. This allows the weight vector to incorporate $b$ into it, for better vectorization.

Essentially, $y = w X + b\mathbf{1} = Xw + b\mathbf{1}$ becomes $y = Xw$, to get rid of $b$.

We want to get rid of $b$ for better vectorization.

P.g. 11 "To show that $z^*$ minimizes $f(z)$, show that $\forall z, f(z) \geq f(z^*)$"

this line refers to the loss function ($f$ is a analogy to $\mathcal{R}$). To show that a particular $z^*$ (* means we are talking about a particular one / optimal one), in which $z^*$ is an analogy to parameters $w$ and $b$, minimizes $f/\mathcal{R}$, you need to show that for all combinations of $(w, b)$, or for all possible $z$, $f(z^*)$ is the smallest ($z^*$ minimizes the loss function).

P.g. 12. ① $\dfrac{\partial y}{\partial w_j} = \dfrac{\partial}{\partial w_j}\left(\sum_{j'}^{D} w_{j'} x_{j'} + b\right)$

$\longrightarrow$ looping through all dimensions

$= \dfrac{\partial}{\partial w_j}(w_0 x_0 + w_1 x_1 + \cdots + w_j x_j + \cdots + w_D x_D + b)$

$= \dfrac{\partial}{\partial w_j}(w_0 x_0)^0 + \cdots + \dfrac{\partial}{\partial w_j}(w_j x_j)^{1 \cdot x_j} + \cdots + \dfrac{\partial}{\partial w_j}(w_D x_D)^0 + \dfrac{\partial}{\partial w_j} b^0$

$= x_j.$

② $\dfrac{\partial y}{\partial b} = \dfrac{\partial}{\partial b}\left(\sum_{j'}^{D} w_{j'} x_{j'} + b\right) = \dfrac{\partial}{\partial b}(w_0 x_0)^0 + \cdots + \dfrac{\partial}{\partial b}(w_D x_D)^0 + \dfrac{\partial}{\partial b} b^1 = 1$

③ $\dfrac{\partial \mathcal{L}}{\partial w_j} = \dfrac{\partial \mathcal{L}}{\partial y} \cdot \boxed{\dfrac{\partial y}{\partial w_j}}^{x_j} = \dfrac{d}{dy}\left(\tfrac{1}{2}(y-t)^2\right) x_j = (y-t) x_j$

④ $\dfrac{\partial \mathcal{L}}{\partial b} = \dfrac{\partial \mathcal{L}}{\partial y} \cdot \dfrac{\partial y}{\partial b} = \dfrac{d}{dy}\left(\tfrac{1}{2}(y-t)^2\right) \cdot 1 = (y-t)$

**Remark:** we calculate EQ ① and EQ ② to calculate EQ ③ & EQ ④, as ③ & ④ relate to how $w_j$ & $b$ affect the loss.

Remark:
We apply the chain rule to get to the expanded form of ③ & ④ because we cannot do $\dfrac{\partial \mathcal{L}}{\partial w_j}$ & $\dfrac{\partial \mathcal{L}}{\partial b}$ directly — $\mathcal{L}$ is not defined (directly) using $w_j$ and $b$.