

Homework 1

Bryan Pitsker Data Mining

9/15/23

① Let $\vec{X} = [1, x_1, x_2, \dots, x_D]^T = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$

- So we have D features for each x vector.

Let $W = [w_0, w_1, w_2, \dots, w_D]^T = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$

- And we have D weights for each of the D independent features

$y_* = W^T \vec{X}_* = \begin{bmatrix} w_0, w_1, \dots, w_D \end{bmatrix} \cdot \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix} =$

prediction for y -output

$$y_* = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

Now if we let

$$X = \begin{bmatrix} \vec{X}_1^T \\ \vec{X}_2^T \\ \vdots \\ \vec{X}_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_D \\ 1 & x_1 & x_2 & \dots & x_D \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1 & x_2 & \dots & x_D \end{bmatrix}$$

Note

" $\vec{X}_1, \vec{X}_2, \dots, \vec{X}_N$
are all individual
vectors that we
initialized above in
the first line"

* So, we have D features/columns in our matrix, but there are N total rows, or observances, of each feature in the matrix.

$$Xw = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_D \\ 1 & x_1 & x_2 & \dots & x_D \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1 & x_2 & \dots & x_D \end{bmatrix} \cdot \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}$$

$$N \cdot D \quad \cdot \quad D \cdot 1$$

$$N \cdot 1$$

$$Xw = \begin{bmatrix} w_0 + x_1 w_1 + x_2 w_2 + \dots + x_D w_D \\ w_0 + x_1 w_1 + x_2 w_2 + \dots + x_D w_D \\ \vdots \\ w_0 + x_1 w_1 + x_2 w_2 + \dots + x_D w_D \end{bmatrix}$$

N rows, 1 column

- Now if the actual outputs y are:

$$y = [y_1, y_2, y_3, \dots, y_N]^T = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Then this means that

$$Xw - y = \begin{bmatrix} (w_0 + x_1 w_1 + \dots + x_D w_D) - y_1 \\ (w_0 + x_1 w_1 + \dots + x_D w_D) - y_2 \\ \vdots \\ (w_0 + x_1 w_1 + \dots + x_D w_D) - y_N \end{bmatrix}$$

Thus, the L2-norm squared is: $\|Xw - y\|^2$

$$\hookrightarrow \|Xw - y\|^2 =$$

$$= \sqrt{\left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_1 \right]^2 + \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_2 \right]^2 + \dots + \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_N \right]^2}$$

The square root and the squared cancel out, so all we are left with is:

$$\|Xw - y\|^2 = \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_1 \right]^2 + \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_2 \right]^2 + \dots + \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_N \right]^2$$

• Therefore $\frac{1}{N} \cdot \|Xw - y\|^2 =$

$$= \frac{1}{N} \cdot \left[\left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_1 \right]^2 + \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_2 \right]^2 + \dots + \left[(w_0 + x_1 w_1 + \dots + x_D w_D) - y_N \right]^2 \right]$$

BUT, this is the exact same thing as writing

$$\Rightarrow \frac{1}{N} \cdot \sum_{n=1}^N (w^T \vec{x}_n - y_n)^2 =$$

reminder: X is the vector here, NOT an observance or a point

$$= \frac{1}{N} \cdot \left[(w^T \vec{x}_1 - y_1)^2 + (w^T \vec{x}_2 - y_2)^2 + \dots \right. \\ \left. \dots + (w^T \vec{x}_N - y_N)^2 \right]$$

Note: $w^T \cdot \vec{x}_*$ we showed before to be

$$w^T \cdot \vec{x}_* = w_0 + w_1 x_1 + \dots + w_D x_D$$

So, \rightarrow the equation at the very top becomes

$$= \frac{1}{N} \cdot \left[[(w_0 + w_1 x_1 + \dots + w_D x_D) - y_1]^2 \right. \\ \left. + [(w_0 + w_1 x_1 + \dots + w_D x_D) - y_2]^2 + \dots \right. \\ \left. \dots + [(w_0 + w_1 x_1 + \dots + w_D x_D) - y_N]^2 \right]$$

which is what we got before!

Hence, we proved that the Mean Squared Error:

$$E(w) = \frac{1}{N} \cdot \sum_{n=1}^N (w^T x_n - y_n)^2 = \frac{1}{N} \cdot \| \underbrace{X}_{\text{matrix } X} w - y \|^2$$