

(1)

1.3 Vector / Matrix Notation

Add more parameters:

$$t = f(x, s_1, \dots, s_8; w_0, \dots, w_9) = \\ w_0 + w_1 x + w_2 s_1 + w_3 s_2 \dots$$

As you add more paras the number of equations and partial derivatives becomes unfeasibly large

the solution is to use vectors & matrices

x^T = vector or matrix transpose
 it is easier to read as a row so
 the transpose sign is used to tell
 us that a row is actually a col vect

$$x_n = [x_{n1}, \dots, x_{nq}]^T$$

Recall matrix is called rows x cols

(2)

Converting $t = w_0 + w_1 x$ to vector form

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad x_n = \begin{bmatrix} 1 \\ x_n \end{bmatrix}$$

Parameter

Data Vector

$$f(x_n; w_0, w_1) = w^T x_n = w_0 + w_1 x_n$$

matrix/vector multiplication

Any instance of $w_0 + w_1 x_n$ can be replaced with $w^T x_n$

$$\text{Sgloss} = \frac{1}{N} \sum (t_n - \underline{w^T x_n})^2$$

$\hookrightarrow \frac{1}{N} (t - X_w)^T (t - X_w)$

How does this notion work?

First collect all x_n & t_n into vectors

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \quad t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$$

(3)

Perform MM on X_w

$$X_w = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \times \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ \vdots \\ w_0 + w_1 x_n \end{bmatrix} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}$$

Note that this summarises the $t_n - w_0 - w_1 x_n$ component and also encapsulates the \sum part

(?) wrong needs matrix for this

to capture the squared and sum part of the loss function
MM $(X_w - t)$ with itself

$$t - X_w = t_i - w_0 - w_1 x_i$$

$$(t - X_w)^2 = (t_i - w_0 - w_1 x_i)^2$$

$$\hookrightarrow (t - X_w)^T (t - X_w)$$

creates a pairwise combination of $(t_n - X_w_n)$ against their which squares it and then sums it also

$$\text{Hence } L = \frac{1}{n} (t - X_w)^T (t - X_w)$$

(4)

Factorize this down:

$$\begin{aligned}
 h &= \frac{1}{N} (X_w - t)^T (X_w - t) \\
 &= ((X_w)^T - t^T)(X_w - t) \\
 &= \frac{1}{N} (X_w)^T X_w - \frac{1}{N} t^T X_w - \frac{1}{N} (X_w)^T t + \\
 &\quad \frac{1}{N} t^T t
 \end{aligned}$$

$$h = \frac{1}{N} w^T X^T X_w - \frac{2}{N} w^T X^T t + \frac{1}{N} t^T t$$

Differentiating loss in vector/matrix form

We want to work out the value of the w vector (w_0, w_1) at the turning point (min)

take partial Derivative of h with respect to w vector

take PD of h by each element of w & stack into a vector

$$\frac{\partial h}{\partial w} = \begin{bmatrix} \frac{\partial h}{\partial w_0} \\ \frac{\partial h}{\partial w_1} \end{bmatrix}$$

recall we calced both of these PD's in the non-matrix example

5

Start with the multiplied out expression from above

$$L = \frac{1}{n} (\omega^T X^T X \omega - 2\omega^T X^T E + E^T E)$$

We need multiply out these terms
which is a long & involved process
of transposing and matrix multiplication

note $t^T t$ term is ignore as doesn't include w
result:

$$\underbrace{w_0^2 + 2w_0 w_1 \frac{1}{N} (\sum x_n) + w_1^2 \frac{1}{N} (\sum x_n^2)}_{w^T X^T X w} + 2w_0 \frac{1}{N} (\sum x_{n \text{ not } t}) + 2w_1 \frac{1}{N} (\sum x_{n \text{ not } t}),$$

$w^T X^T X w$

$2w^T X^T t$

$$\text{reach } w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

note the Σ 's hold shorthand averages
 (that we went over in the previous
 section) that are relevant for calculating
 linear models, i.e. $\frac{1}{n} \sum x_n = \bar{x} = \arg \min$

This allows us to write the expression much more concise

Partial Diff this equation =

$$\frac{\partial L}{\partial w_0} = 2w_0 + 2w_1 \bar{n} - 2\bar{t}$$

$$\frac{\partial L}{\partial w_1} = 2w_0 \bar{n} + 2w_1 \bar{n}^2 - 2\bar{t}$$

Here instead of setting the PD's to 0 and solving to find the optimal w_0, w_1 , the author highlights how we can stay in the vectorized land using identities:

recall the vectorized loss function:

$$L = \frac{1}{n} (t - Xw)^T (t - Xw)$$

~~$$or L = \frac{1}{n} t^T t - 2t^T Xw + w^T X^T Xw$$~~

Expand this expression out using the properties of transposes & matrix multi (long process)

then you take the partial derivative in matrix form using the identities that the author provides

How To Partial Diff vectors/matrix

(7)

Expanded & Simplified loss function:

$$L = \frac{1}{N} * \left(t^T t - \underbrace{2w^T X^T t}_{t_1} + \underbrace{w^T X^T X w}_{t_2} \right) + \underbrace{\lambda w^T w}_{t_3}$$

useful identities when differentiating w/ resp to a vec

$f(w)$	$\frac{\partial f}{\partial w}$
$w^T X$	X
$X^T w$	X
$w^T w$	$2w$
$w^T C w$	$2Cw$

term 1: $t^T t$, no w vector so deriv = 0

term 2: $-2w^T X^T t$ rewrite as $2(X^T t)^T w$

- Treat $(X^T t)$ as a constant as it has no reference to $w = n$
- now, $-2w^T X$
- using identity rule, $w^T X = X$
why? in calculus deriv of $a^n = a$ where a is a constant (Power rule)
- Now, $-2X$

(8)

- However this result is not in the format we want. We are trying to optimise w to minimize the loss which depends on X & t
- Recall that in step 1 $x = X^T t$. we can sub this back in
- So, $-2x = -2X^T t$

term 3: $w^T X^T X w$

- $X^T X = C$ where $C = \text{Symmetric Matrix}$
~~involves growth of w~~
- $X^T X$ creates a square matrix as $X^T = p \times n$ and $X = n \times p$
- Symmetric matrix = its transpose = itself. the proof that something is symmetric is confusing so i wont write it here
- Now, $t^T \beta = w^T C w$
- Next identity says $w^T C w = 2 C w$
the reason why is do with the quality of the symmetric matrix which doubles up the summation due to the product rule: $uv = u'v + v'u'$
- To bring back the terms $2 C w = 2 X^T X w$

(9)

Putting this differentiating together:

$$\text{Vector loss: } \frac{1}{N} (t - Xw)^T (t - Xw)$$

$$\text{Expand & Simplify: } \frac{1}{N} (t^T t - 2w^T X^T t + w^T X^T Xw)$$

$$\text{Partial Diff } w: \frac{1}{N} (0 - 2X^T t + 2X^T Xw)$$

$$\text{Reorder & Simplify: } \frac{2}{N} X^T Xw - \frac{2}{N} X^T t$$

Next, set the PD to zero & solve for w

- Steps:
- Multiply both sides by $\frac{2}{N}$
 - Isolate $w = + X^T t$ both sides

$$① \frac{2}{N} X^T Xw - \frac{2}{N} X^T t = 0$$

$$② X^T Xw - X^T t = 0$$

$$③ X^T Xw = X^T t$$

Next use inverse matrix to remove $X^T X$.
note dividing isn't defined for matrices

$$(X^T X_w)^{-1} X^T X_w = \text{Identity Matrix} = I$$

$$④ Iw = (X^T X)^{-1} X^T t$$

$$\text{Opt value of } \hat{w} = (X^T X)^{-1} X^T t$$

(6)

Numerical Example using Small Data

n	x_n	t_n	$x_n t_n$	x_n^2
1	1	4.8	4.8	1
2	3	11.3	33.9	9
3	5	17.2	86	25
$\frac{1}{n} \sum_{n=1}^n$	3	11.1	41.57	11.67

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 5 \end{bmatrix} \quad t = \begin{bmatrix} 4.8 \\ 11.3 \\ 17.2 \end{bmatrix}$$

① calculate $X^T X$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 3 & 5 \end{bmatrix} \times \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 5 \end{bmatrix} = \begin{bmatrix} 3 & 9 \\ 9 & 35 \end{bmatrix}$$

② calculate this inverse $(X^T X)^{-1}$

$$\text{Formula} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \quad \frac{1}{24} \begin{bmatrix} 35 & -9 \\ -9 & 3 \end{bmatrix}$$

③ multiply inverse by ~~26~~ $X^T = (X^T X)^{-1} X^T$

$$\frac{1}{24} \begin{bmatrix} 35 & -9 \\ -9 & 3 \end{bmatrix} \times \begin{bmatrix} 1 & 1 & 1 \\ 1 & 3 & 5 \end{bmatrix} = \frac{1}{24} \begin{bmatrix} 26 & 8 & -10 \\ -6 & 0 & +6 \end{bmatrix}$$

(11)

④ finally multiply by $t: ((x^T x)^{-1} x^T) t$

$$\frac{1}{24} \begin{bmatrix} 26 & 8 & -10 \\ -6 & 0 & 6 \end{bmatrix} \times \begin{bmatrix} 4.8 \\ 11.3 \\ 17.2 \end{bmatrix} = \begin{bmatrix} 1.8 \\ 3.1 \end{bmatrix}$$

⑤ there the formula is:

$$f(x; w_0, w_1) = 1.8 + 3.1x$$

Note the textbook page 25
 Also walks through an example
 using only mathematical notation

Making predictions

$$t_{\text{new}} = \hat{w}^T X_{\text{new}} \quad \text{where } \hat{w} = \begin{bmatrix} 1.8 \\ 3.1 \end{bmatrix}$$