

WORKSHEET 2

1

PROBLEM 1

Find optimal parameter value, $\hat{\underline{w}}$, for

$$\mathcal{L} = \sum_{n=1}^N (t_n - \underline{w}^T \underline{x}_n)^2$$

Note: t_n is a scalar quantity

$$\underline{w} = (w_0, \dots, w_D)^T$$

$$\underline{x} = (x_0, \dots, x_D)^T$$

To find $\hat{\underline{w}}$, set $\frac{d\mathcal{L}}{d\underline{w}} = 0$

Find $\frac{d\mathcal{L}}{d\underline{w}}$ by application of the chain rule:

$$\frac{d(f(g(x)))}{dx} = \frac{df(g(x))}{dg(x)} \cdot \frac{dg(x)}{dx}$$

or by substitution:

$$\text{let } u_n = t_n - \underline{w}^T \underline{x}_n$$

$$\Rightarrow \frac{du_n}{d\underline{w}} = -\underline{x}_n$$

$$\Rightarrow \mathcal{L} = \sum_{n=1}^N u_n^2$$

$$\Rightarrow \frac{d\mathcal{L}}{du} = 2 \sum_{n=1}^N u_n = 2 \sum_{n=1}^N (t_n - \underline{w}^T \underline{x}_n)$$

Then
$$\frac{dL}{d\mathbf{w}} = \frac{dL}{du} \cdot \frac{du}{d\mathbf{w}} = 2 \sum_{n=1}^N (t_n - \mathbf{w}^T \mathbf{x}_n) (-\mathbf{x}_n)$$

Note that
$$\mathbf{w}^T \mathbf{x}_n = \mathbf{x}_n^T \mathbf{w}$$

and $(t_n - \mathbf{w}^T \mathbf{x}_n)$ is a scalar

$$\begin{aligned} \Rightarrow \frac{dL}{d\mathbf{w}} &= 2 \sum_{n=1}^N (t_n - \mathbf{x}_n^T \mathbf{w}) (-\mathbf{x}_n) \\ &= -2 \sum_{n=1}^N (\mathbf{x}_n) (t_n - \mathbf{x}_n^T \mathbf{w}) \\ &= -2 \left(\left(\sum_{n=1}^N \mathbf{x}_n t_n \right) - \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \mathbf{w} \right) \right) \end{aligned}$$

Setting
$$\frac{dL}{d\mathbf{w}} = 0$$

$$\begin{aligned} \Rightarrow \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \hat{\mathbf{w}} &= \sum_{n=1}^N \mathbf{x}_n t_n \\ \hat{\mathbf{w}} &= \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right)^{-1} \sum_{n=1}^N \mathbf{x}_n t_n \end{aligned}$$

Note that in matrix/vector notation

$$L = (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$

where
$$\mathbf{t} = (t_1, \dots, t_N)^T \quad \mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$$

Then
$$\frac{dL}{d\mathbf{w}} = 2(-\mathbf{X})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$

Setting
$$\frac{dL}{d\mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{t}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

as above.

\hat{w} provides a critical point for L .

Use the second derivative test to check whether this critical point corresponds to a maximum or a minimum point

$$\begin{aligned} \frac{\partial^2 L}{\partial w^2} &= -2 \left(- \sum_{n=1}^N x_n x_n^T \right) = 2 \sum_{n=1}^N x_n x_n^T \\ &= 2 X^T X \end{aligned}$$

So, if $X^T X$ is positive definite

(which as a sum of squares it will be)

then \hat{w} corresponds to a minimum when X is full rank, $N \geq d$.

and the squared loss function L is minimized at \hat{w} .

Scaling L by a constant will have no impact on the optimal parameter \hat{w} , thus the \hat{w} remains unchanged for both L and λL .

PROBLEM 2

Consider the weighted average loss

$$L = \frac{1}{N} \sum_{n=1}^N \alpha_n (t_n - \underline{w}^T \underline{x}_n)^2$$

where $t_n, \underline{w}, \underline{x}_n$ are defined as in Problem 1, and the weight / importance of the n^{th} data point, $\{t_n, \underline{x}_n\}$ is given by α_n

Find the optimal least squares parameter estimate by setting $\frac{dL}{d\underline{w}} = 0$

$$\frac{dL}{d\underline{w}} = \frac{1}{N} \sum_{n=1}^N 2\alpha_n (t_n - \underline{w}^T \underline{x}_n) (-\underline{x}_n)$$

which as in problem 1

$$= -\frac{2}{N} \left(\sum_{n=1}^N \alpha_n \underline{x}_n t_n - \sum_{n=1}^N \alpha_n \underline{x}_n \underline{x}_n^T \underline{w} \right)$$

setting $\frac{dL}{d\underline{w}} = 0 \Rightarrow \hat{\underline{w}} = \frac{\sum_{n=1}^N \alpha_n \underline{x}_n t_n}{\sum_{n=1}^N \alpha_n \underline{x}_n \underline{x}_n^T}$

or in matrix notation

$$\begin{aligned} L &= \frac{1}{N} \sum_{n=1}^N \alpha_n (t_n - w^T x_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N (t_n - w^T x_n) \alpha_n (t_n - w^T x_n) \\ &= \frac{1}{N} \sum_{n=1}^N (\underline{t} - Xw)^T A (\underline{t} - Xw) \end{aligned}$$

where $A = \begin{pmatrix} \alpha_1 & & & \\ & \alpha_2 & & 0 \\ & & \ddots & \\ 0 & & & \alpha_N \end{pmatrix}$ is the diagonal matrix of observation weights.

$$\text{then } \frac{dL}{dw} = \frac{1}{N} \cdot 2 \cdot (-X)^T \cdot A \cdot (\underline{t} - Xw)$$

$$= -\frac{2}{N} (X^T A \underline{t} - X^T A X w)$$

$$\Rightarrow \hat{w} = (X^T A X)^{-1} X^T A \underline{t}$$

as above.

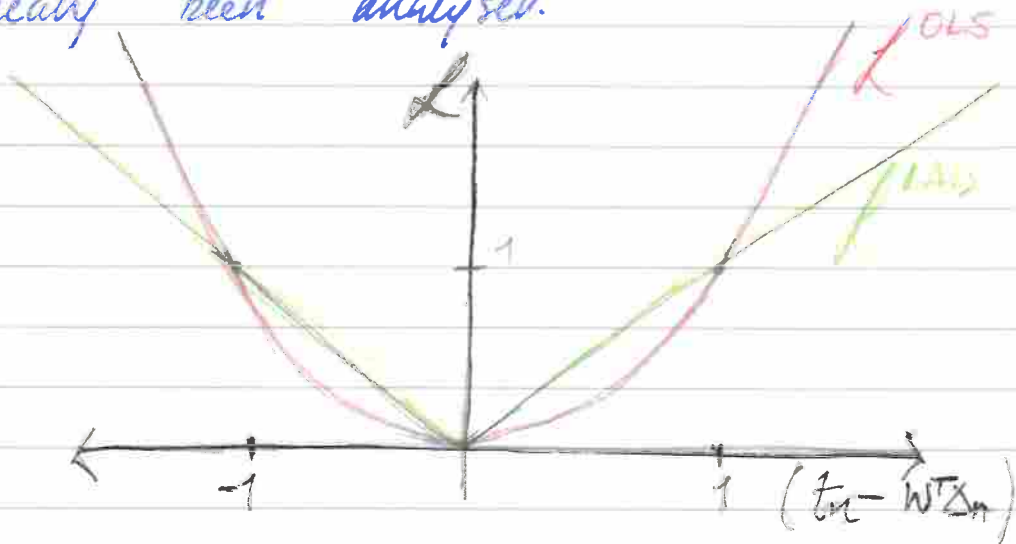
Note: For a guide to matrix/vector calculus see the Matrix lookbook, for which there is a link on the module webpage.

PROBLEM 3

Consider the loss function

$$L^{LAD} = \sum_{n=1}^N |t_n - w^T x_n|$$

versus the $L^{OLS} = \sum_{n=1}^N (t_n - w^T x_n)^2$ which has already been analysed.



Examining the plot above illustrates that minimising L^{LAD} places equal weight on all observations, while minimising L^{OLS} places a greater weight on those observations for which there are large errors.

Thus L^{OLS} is more heavily influenced by outliers when estimating the optimal parameter values, making L^{LAD} more robust to these outliers.

To find the optimal parameter values for L^{LAD} set $\frac{dL^{LAD}}{dW} = 0$

For $t_n - W^T \underline{x}_n > 0$

$$L^{LAD} = t_n - W^T \underline{x}_n$$

$$\frac{dL^{LAD}}{dW} = -\underline{x}_n \neq 0$$

For $t_n - W^T \underline{x}_n < 0$

$$L^{LAD} = -(t_n - W^T \underline{x}_n)$$

$$\frac{dL^{LAD}}{dW} = \underline{x}_n \neq 0$$

For $t_n - W^T \underline{x}_n = 0$

$$L^{LAD} = \pm (t_n - W^T \underline{x}_n)$$

$$\frac{dL^{LAD}}{dW} = \mp \underline{x}_n \neq 0$$

Thus L^{LAD} has no critical points and there is no analytic solution for optimal parameter values.