

Ridge - Regu^l, exact:

We now minimize

$$\mathcal{L} = \frac{1}{N} \sum_n (w x_n - t_n)^2 + 2 \underbrace{\|w\|_2^2}_{= w^T \cdot w}$$

$$\nabla \mathcal{L} = \frac{2}{N} \sum_n (w x_n - t_n) x_n + 2\lambda \vec{w}$$

(We check $\nabla \|w\|^2 = \vec{w}$: $\nabla (w_1^2 + w_2^2 + \dots + w_D^2) = \begin{pmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_D \end{pmatrix}$)

$w^* = \text{argmin } \mathcal{L}$

$\Rightarrow w^*$ is such that $\nabla_{\vec{w}} \mathcal{L} = \vec{0}$

$$\Rightarrow \frac{2}{N} \underbrace{X^T (XW - T)}_{\text{shape } (D, N)(N, D)(D, 1) = (N, 1)} + 2\lambda \underbrace{w}_{\text{shape } (D, 1)} = 0$$

$$\Rightarrow (D, N)(N, 1) \rightarrow (D, 1)$$

$$\Rightarrow X^T (XW - T) + N\lambda W = 0$$

$$W = (X^T X + N\lambda I)^{-1} X^T T$$

(here is an implicit $I_{D,D}$ matrix)

Again, the hard step is inverting $X^T X$, which is just $O(D^3)$. (+ $O(ND^2)$ to new things)

When $D=1$, If $\text{Var}(x) \gg 2$, large w is permitted

$$W = \frac{X^T T}{N\lambda + X^T X} = \frac{\sum_{n=1}^N x_n t_n}{N\lambda + \sum_n x_n^2}$$

For large λ , $w \rightarrow 0$. This has to be compared with $\frac{1}{N} \sum x_n^2$, which is $\text{variance}(x_n)$, essentially.
 $(- \langle x \rangle)$