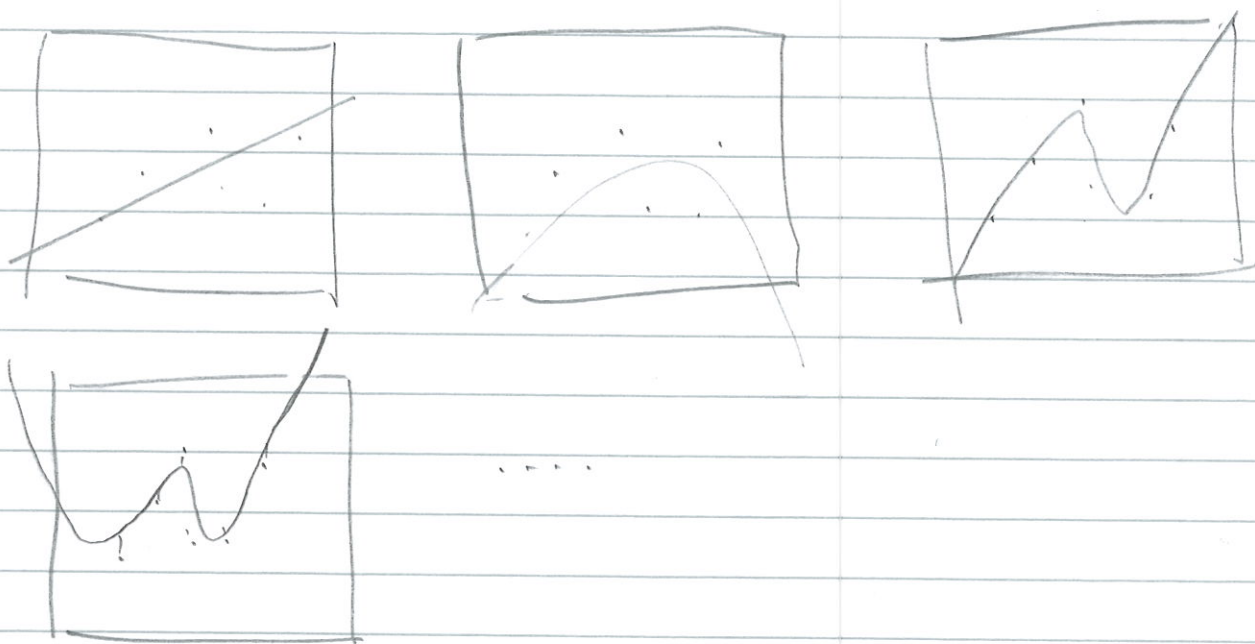


Regularization - a Bayesian Linear Regression (1)

As model complexity increases the ability for a model to fit the data becomes higher.

Polynomial regression serves as an excellent pedagogical tool for explaining what happens when a model begins to overfit:



As the polynomial order increases the ability of the model to exactly reproduce the target values increases; however, out-of-sample performance typically degrades.

To combat this we typically regularize or bias models towards reasonable solutions. In this case we might like the polynomial weights to be small.

One way to do this is to directly penalize large weights, i.e. modifying our ML objective

$$E_D(\omega) = \frac{1}{2} \sum_{n=1}^N \{t_n - \omega^T \phi(x_n)\}^2$$

$$= \frac{1}{2} (\vec{t} - \Phi \vec{\omega})^T (\vec{t} - \Phi \vec{\omega})$$

where E is "energy", i.e. log. prob. and

$$\Phi = \begin{matrix} & \begin{matrix} k \\ \phi_1(x_1) \dots \phi_k(x_1) \\ \vdots \\ \phi_1(x_N) \dots \phi_k(x_N) \end{matrix} \\ \begin{matrix} n \\ \hline \end{matrix} & \begin{bmatrix} \phi_1(x_1) \dots \phi_k(x_1) \\ \vdots \\ \phi_1(x_N) \dots \phi_k(x_N) \end{bmatrix} \end{matrix}$$

to include a term like

$$E_W(\vec{\omega}) = \frac{\lambda}{2} \vec{\omega}^T \vec{\omega}$$

where λ is a smoothing tuning parameter.

Maximum probability is minimum energy so we might wish to find

$$\vec{\omega}^* = \underset{\vec{\omega}}{\operatorname{argmin}} E_D(\vec{\omega}) + E_W(\vec{\omega})$$

$$= \underset{\vec{\omega}}{\operatorname{argmin}} \frac{1}{2} (\vec{t} - \Phi \vec{\omega})^T (\vec{t} - \Phi \vec{\omega}) + \frac{\lambda}{2} \vec{\omega}^T \vec{\omega}$$

which arises at $\frac{\partial}{\partial \vec{\omega}} (E_D(\vec{\omega}) + E_W(\vec{\omega})) = 0$

$$0 = \frac{\partial}{\partial \vec{\omega}} \left(\frac{1}{2} (\vec{t} - \Phi \vec{\omega})^T (\vec{t} - \Phi \vec{\omega}) + \frac{\lambda}{2} \vec{\omega}^T \vec{\omega} \right)$$

$$= -\Phi^T (\vec{t} - \Phi \vec{\omega}) + \lambda \vec{\omega}$$

$$= -\Phi^T \vec{t} + (\Phi^T \Phi + \lambda \mathbf{I}) \vec{\omega}$$

$$\Rightarrow \vec{\omega} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \vec{t}$$

$$u \begin{bmatrix} k \\ \end{bmatrix} \textcircled{2}$$

What does this mean and do?

1) Inverting $\Phi^T \Phi$ becomes unstable or impossible if Φ is rank deficient. Adding positive elements to the diagonal ensures that $(\Phi^T \Phi + \lambda \mathbf{I})$ is full rank and therefore invertible.

2) $\frac{\lambda}{2} \vec{\omega}^T \vec{\omega}$ looks a lot like a norm, i.e. $\frac{1}{2} (\vec{\omega} - \vec{0})^T \lambda \mathbf{I} (\vec{\omega} - \vec{0})$ which in energy terms will be minimized when $\vec{\omega}$ is close to $\vec{0}$.

3) This is a form of BIAS that is helpful!

To interpret this further we appeal to Bayesian reasoning, e.g. Bayesian linear regression.

Let's consider Bayes rule in the context of linear regression.

$$\begin{aligned} p(\vec{\omega} | \vec{z}, \Phi, \beta, \lambda) \\ = \frac{p(\vec{z} | \Phi, \vec{\omega}, \beta) p(\vec{\omega} | \lambda)}{p(\vec{z} | \Phi, \beta, \lambda)} \end{aligned}$$

this is a little confusing so let's drop all the constants

$$\begin{aligned} p(\vec{\omega} | \vec{z}) &= \frac{p(\vec{z} | \vec{\omega}) p(\vec{\omega})}{p(\vec{z})} \\ &\propto p(\vec{z} | \vec{\omega}) p(\vec{\omega}) \\ &= \mathcal{N}(\vec{z}; \Phi \vec{\omega}, \beta \mathbf{I}) \mathcal{N}(\vec{\omega}; \vec{0}, \lambda \mathbf{I}) \end{aligned}$$