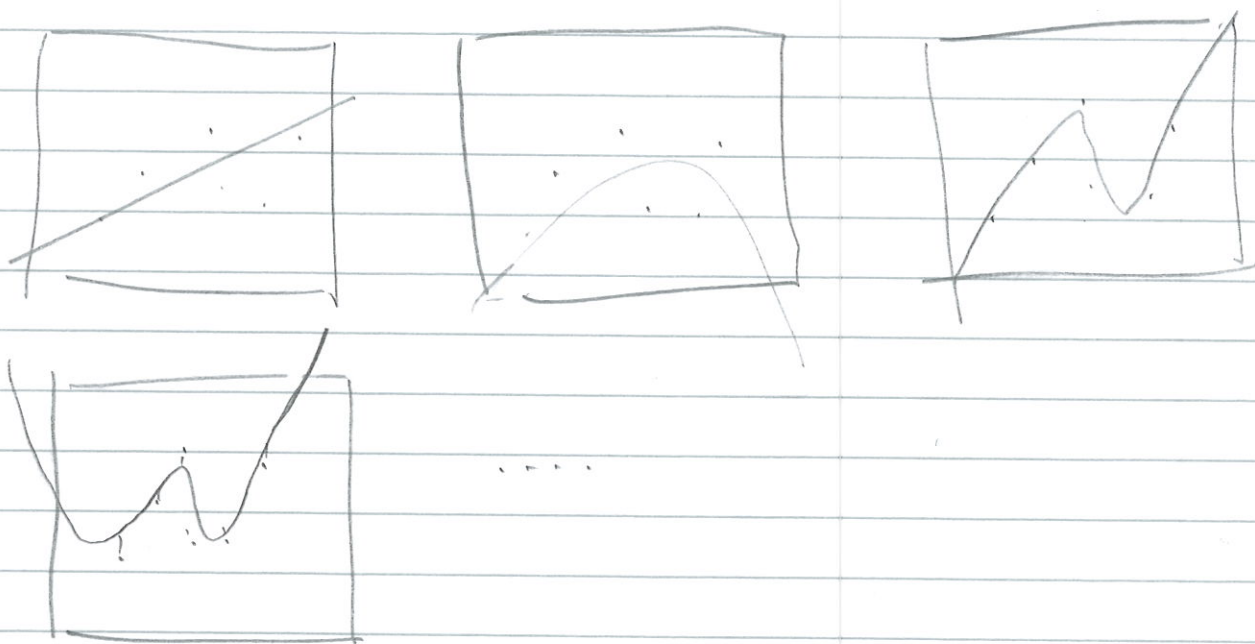


Regularization - a Bayesian Linear Regression (1)

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

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

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

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

This is quite interesting. It says that a Bayesian approach to linear regression introduces bias and, simultaneously, that solving for the MAP (maximum a posteriori) \vec{w} is clearly equivalent to solving the regularized least squares problem.

What is more, we can, here, analytically derive the full posterior distribution. If we know (and understand) the following about the M/V/N

μ, b, γ vectors
 A, L, Λ, Σ matrices

Bishop 2.113 \rightarrow 2.117

$$\Rightarrow \begin{cases} p(x) = \mathcal{N}(x | \mu, \Lambda^{-1}) \\ p(y|x) = \mathcal{N}(y | Ax + b, L^{-1}) \\ p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^T) \\ p(x|y) = \mathcal{N}(x | \Sigma \{ A^T L (y - b) + \Lambda \mu \}, \Sigma) \end{cases}$$

where $\Sigma = (\Lambda + A^T L A)^{-1}$

This is effectively Bayes rule for Gaussians and with it we can immediately derive

$$p(\vec{w} | \dots) = \mathcal{N}(\vec{w} | \Sigma \{ \Phi^T \frac{1}{\beta} \mathbf{I} \left(\frac{\vec{t}}{\beta} \right) \}, \Sigma)$$

$$= \mathcal{N}(\vec{w} | \frac{1}{\beta} \Sigma \Phi^T \vec{t}, \Sigma)$$

$k \times 1$
 $k \times n$
 $n \times 1$
 $k \times k$

and

$$\Sigma = \left(\frac{1}{\lambda} \mathbf{I} + \frac{1}{\beta} \Phi^T \Phi \right)^{-1}$$

which is the posterior distribution of the weight vector. We already know its mode and mean: having the full distribution allows ease of model combination and propagation of uncertainty throughout inference and computation.

Leaving that for more advanced treatments, we usually, in practice, are interested in making predictions about an output t for new data x .

This can be written as

$$p(t | x, \vec{\omega}, \Phi, \lambda, \beta) = \int p(t | \vec{\omega}, x, \beta) \underbrace{p(\vec{\omega} | \vec{t}, \Phi, \lambda, \beta)}_{\text{posterior of } \vec{\omega}} d\vec{\omega}$$

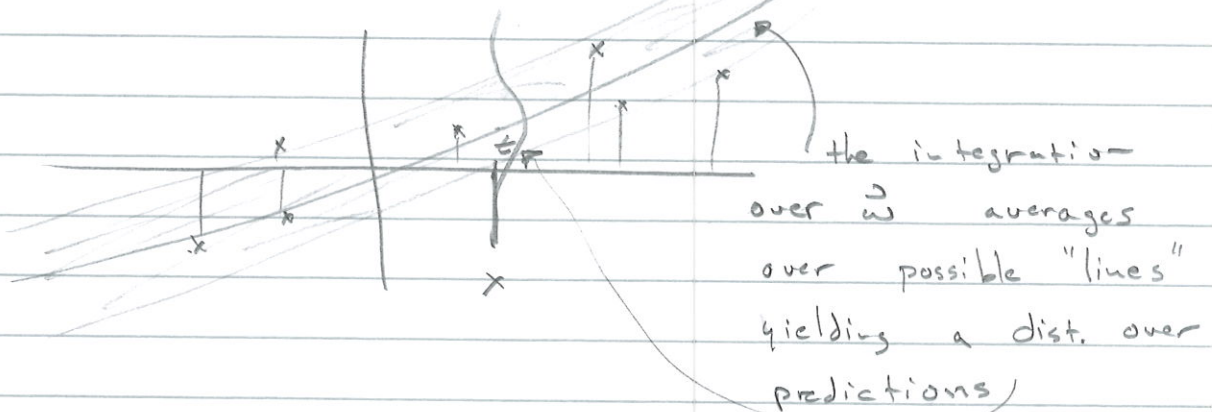
$\underbrace{p(t | \vec{\omega}, x, \beta)}_{\text{model of noise}}$

$$\Rightarrow p(t | \dots) = N\left(t; \vec{x}^T \vec{\omega}, \beta\right) N\left(\vec{\omega} | \frac{1}{\beta} \Sigma \Phi^T \vec{t}, \Sigma\right)$$

$\Sigma = \left(\frac{1}{\lambda} \mathbf{I} + \frac{1}{\beta} \Phi^T \Phi\right)^{-1}$

$$= N\left(t; \vec{x}^T \Sigma \Phi^T \vec{t}, \frac{1}{\beta} + \vec{x}^T \left(\frac{1}{\lambda} \mathbf{I} + \frac{1}{\beta} \Phi^T \Phi\right)^{-1} \vec{x}\right)$$

pictorially this looks like



As we can see here even the predictive dist. is Gaussian in the Bayesian linear regression setting. What does this mean though? By propagating uncertainty about the value of the weights throughout the computation we get error bars on our predictions. Analyses using these are constrained by the model family but otherwise are safer, in general, than using point estimates.