

Polynomial Curve Fitting

We start by solving a simple problem in ML. Suppose we have a *training set* of N values, that is, a vector $\mathbf{x} = \{x_1, \dots, x_N\}$ corresponding to $\mathbf{t} = \{t_1, \dots, t_N\}$. Now given a new value x , we want to predict its corresponding value t .

Given a set of coefficients $\mathbf{w} = \{w_0, \dots, w_M\}$ for some $M \in \mathbb{N}$, we have a polynomial of degree M expressed as $y(x, \mathbf{w}) = \sum_{i=0}^M w_i x^i$. Now we want to approximate the relation between \mathbf{x} and \mathbf{t} using a polynomial. We want to minimise an *error function* that tells us how good the approximation is given coefficients \mathbf{w} . The error function evaluates to 0 if and only if it passes through every $t \in \mathbf{t}$. Below is a simple way to do it.

Definition. The *sum of squares* error function E is given by

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2.$$

Deriving this gives us a linear map with a trivial kernel, so there is a unique \mathbf{w}^* such that $E(\mathbf{w}^*)$ is minimal. We still need to choose M , if it's too large then we have the problem of over-fitting, if too little, we don't have enough flexibility to fit accurately to the training set.

We could rigorously test whether over-fitting is a problem by using our function $y(x, \mathbf{w}^*)$ against a set with way more datapoints. Hence, we require the following definition that generalises N .

Definition. The *root-means-square* error function is defined as follows.

$$E_{\text{RMS}}(\mathbf{w}) = \sqrt{\frac{1}{N} E(\mathbf{w})}.$$

Curve Fitting Revisited

(Please read probability theory and fundamentals of Gaussian distribution first.)

We now assume the value we are trying to predict, t , has a Gaussian distribution with $\mu = y(x, \mathbf{w})$, and that

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(y(x, \mathbf{w}), \beta^{-1}).$$

Recall that $\beta = \frac{1}{\sigma^2}$. We assume the datapoints are independent, and we maximise the likelihood function below

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}).$$

This can be done by maximising its log.

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, \mathbf{w}))^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi$$

If we simplify the process of maximising \mathbf{w} by deleting the addition of the constant terms, then let $\beta = 1$, and minimise instead the negative log of $p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)$, it follows now that this maximisation process is exactly the minimisation of the sum-of-squares error function. Also note that

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N (y(x_n, \mathbf{w}_{\text{ML}}) - t_n)^2.$$

We may now predict a distribution given a new (x, t) . This is called a *predictive distribution*.

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1}).$$

Now we wish to introduce a more Bayesian approach. Given a precision α , we define a prior.

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M+1}{2}} \exp\left(-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right).$$

A parameter like α is called a *hyper-parameter*. The motivation is to define a probability distribution over \mathbb{R}^{M+1} around the origin with a precision α . Then, by Bayes' theorem,

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha).$$

We then use *maximum posterior* or *MAP* by minimising the negative log. We find that the maximum of the posterior is given by the minimum of the following

$$\frac{\beta}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}))^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}.$$

So maximising the posterior is to minimise the sum of squares error function with regularisation parameter $\lambda = \alpha/\beta$.

For a full Bayesian approach, we repeatedly apply the sum and product rules.

Ultimately, we wish to find $p(\mathbf{t}|\mathbf{x}, \mathbf{x}, \mathbf{t})$. Hence, we write it in the following form.

$$p(t|x, \boldsymbol{x}, \boldsymbol{t}) = \int p(t|x, \boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{x}, \boldsymbol{t}) \, d\boldsymbol{w}.$$

Recall that $p(t|x, \boldsymbol{w}) = \mathcal{N}(t|y(x, \boldsymbol{w}), \beta^{-1})$.

We will see that

$$p(t|x, \boldsymbol{x}, \boldsymbol{t}) = \mathcal{N}(t|m(x), s^2(x)).$$

where the mean and variance are given by

$$\begin{aligned} m(x) &= \beta \boldsymbol{\phi}(x)^T \boldsymbol{S} \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n \\ s^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \boldsymbol{S} \boldsymbol{\phi}(x) \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{S}^{-1} &= \alpha \boldsymbol{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \\ \boldsymbol{\phi}(x) &= (x, x^2, x^3, \dots, x^M)^T. \end{aligned}$$