

# 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.

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