

Introduction to Machine Learning

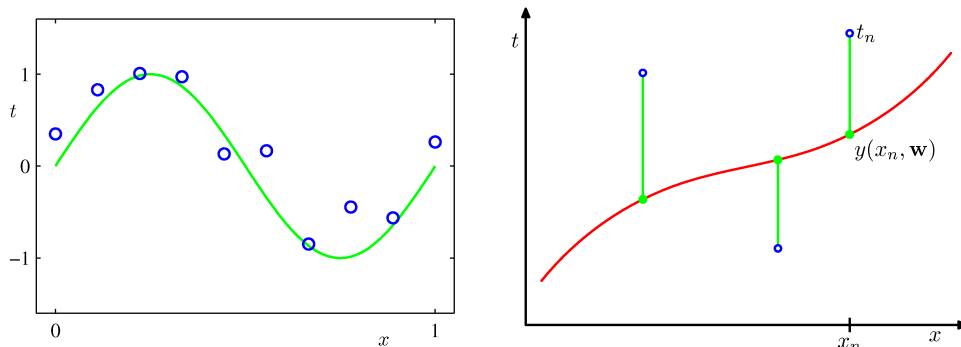
Estimating Parameters, Regularization

Maheśan Niranjana
University of Southampton

DiscNet Summer School

June 2025

Polynomial Curve Fitting



- One dimensional inputs and targets: $\mathbf{x} = [x_1, x_2, \dots, x_N]$; $\mathbf{t} = [t_1, t_2, \dots, t_N]$
- Fit a function (Eqn 1.1):

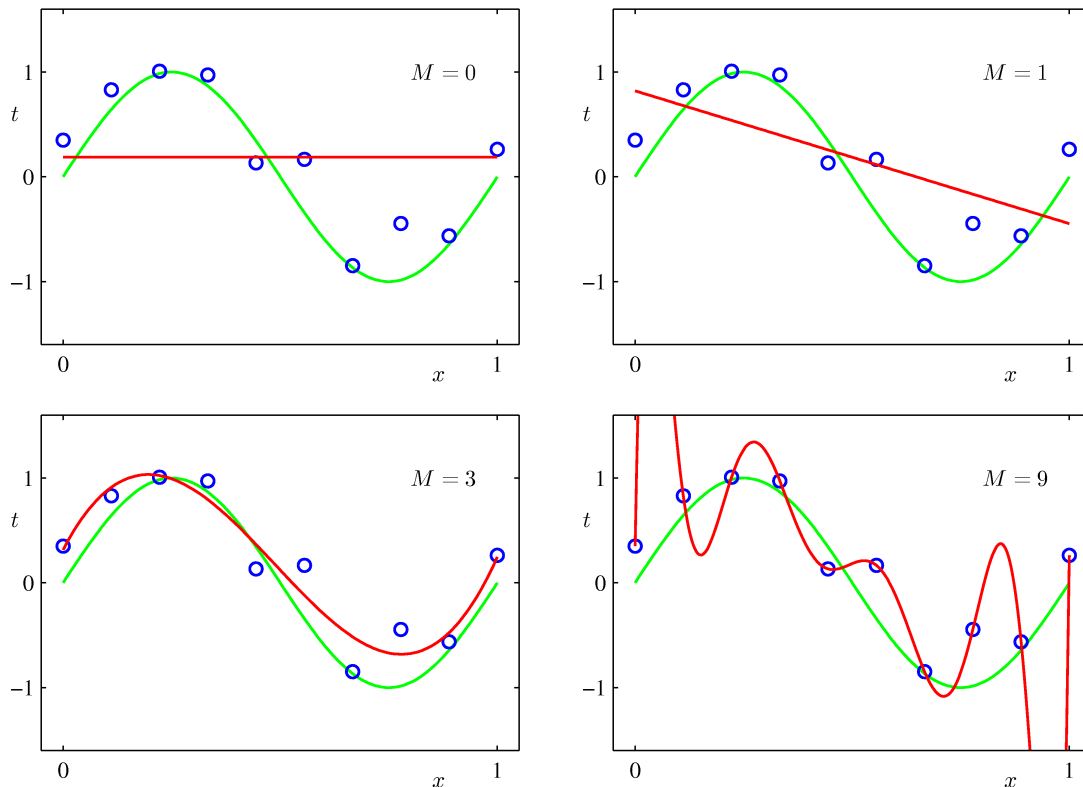
$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_j x^j$$

- Minimize error:

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

Polynomial Curve Fitting

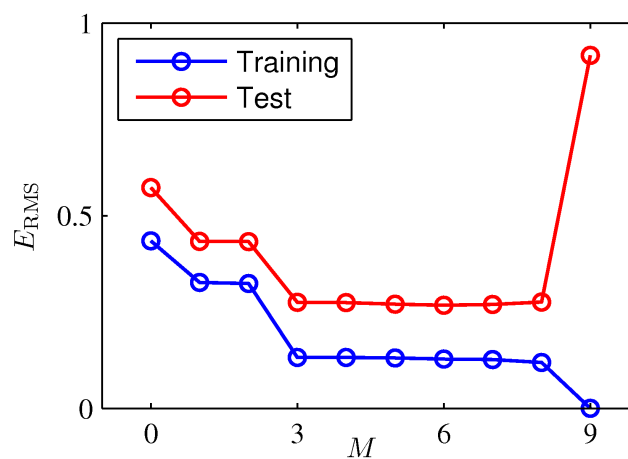
Illustrating Model Complexity



Error with Polynomial Order

- Root mean squared error (RMS)

$$E_{RMS} = \sqrt{2E(\mathbf{w}^*)/N}$$



Regularized Least Squares

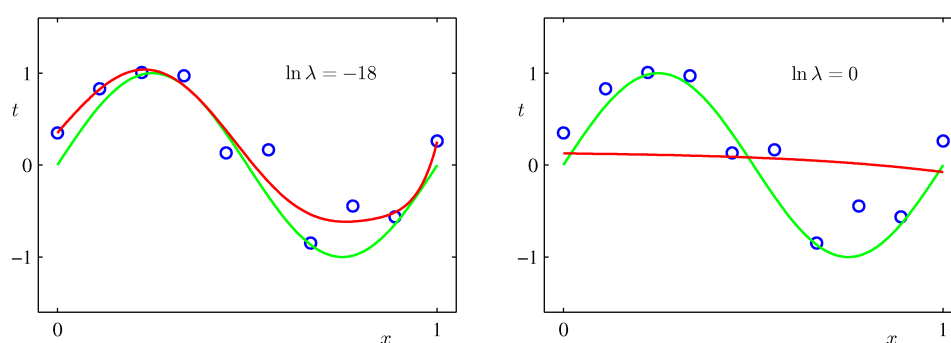
- As M increases, weight magnitudes get large
- Large positive and large negative weights cancelling out and fitting the data

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

Regularized Least Squares

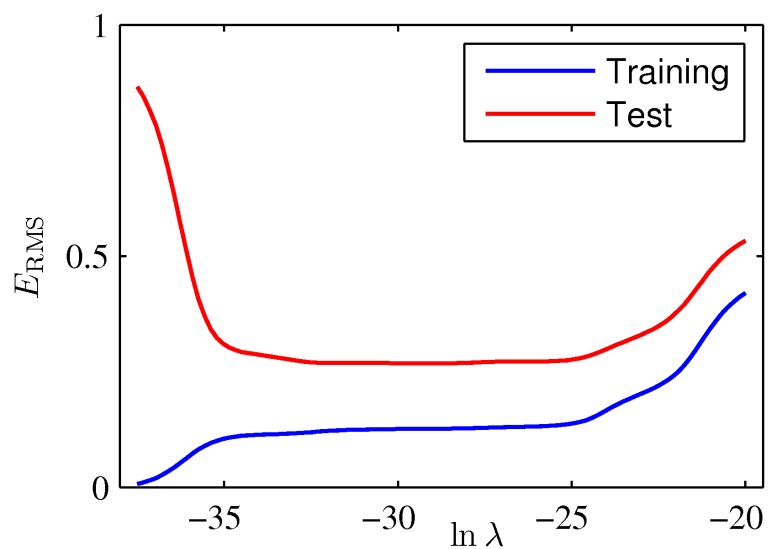
- The way to deal with this is *regularization*, a very important tool in machine learning.

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

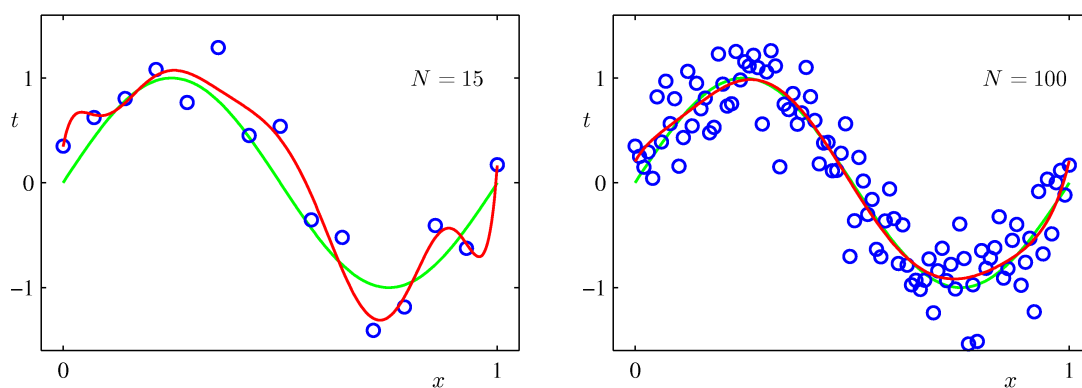


- Note the regularization weight in the figures is $\ln \lambda$
- Left hand figure is just about correct, right hand figure is over-doing it.

Choice of Regularization Parameter



Effect of dataset size



- Having a lot of data is good
- Lot of excitement around “*big data*”
- But is it always the case?

Probabilistic Inference

- In polynomial curve fitting earlier, we estimated the coefficients \mathbf{w} by minimising the squared error.
- We will formulate it in a probabilistic setting.
- First we note *Bayesian* versus *Frequentist* arguments

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

- Posterior \propto Likelihood \times Prior
- Denominator is integral over numerator:

$$p(\mathcal{D}) = \int p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})$$

- We will see this in classification, regression and estimation
- Bayesian vs Frequentist paradigms [Section 1.2.3]

Bayesian	Just one dataset, uncertainty through $p(\mathbf{w}) \rightarrow p(\mathbf{w} \mathcal{D})$
Frequentist	\mathcal{D} is just one realization of a process

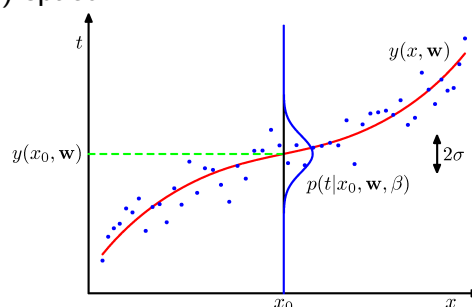
Polynomial Regression in a Probabilistic Setting

Section 1.2.5

- Data consists of Inputs $\mathbf{x} = (x_1, x_2, \dots, x_N)^T$ and corresponding targets $\mathbf{t} = (t_1, t_2, \dots, t_N)^T$
- Our probabilistic model is to assume that given x , the target is generated by a function $y(x; \mathbf{w})$ and is corrupted by zero mean Gaussian noise:

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

- β is precision parameter, inverse of variance.
- We assume (and usually this is the case) the data are independent, identically distributed (IID) samples from the (x, t) space.



- Likelihood of the data:

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

Probabilistic setting (cont'd)

- Likelihood (again) is a function of the parameters \mathbf{w} and β

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

- What values of the unknowns maximises the likelihood?
- It is convenient to work with log likelihood
 - Product turns into summation
 - Exponent can be avoided
- Log likelihood:

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

- When we maximize this respect to \mathbf{w} , we see we are minimising the mean squared error, giving \mathbf{w}_{ML} (ML Estimation)
- How exactly to solve for \mathbf{w}_{ML} , we will see later.
- When we maximize with respect to β , we get

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

Probabilistic setting (cont'd)

- With a probabilistic model, we now think of a *predictive distribution*

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

- We can also take a Bayesian inference perspective
- Assume a prior distribution over the parameters – an isotropic Gaussian

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

- Make sure you understand the above to be a special case of the general expression: $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ seen earlier!
- We can now write a posterior distribution over \mathbf{w}

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

- We can now maximize the posterior distribution (MAP Estimation); the maximum of this is the same as the minimum of

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

- Maximising the posterior is the same as minimising the *regularized* sum of squares.