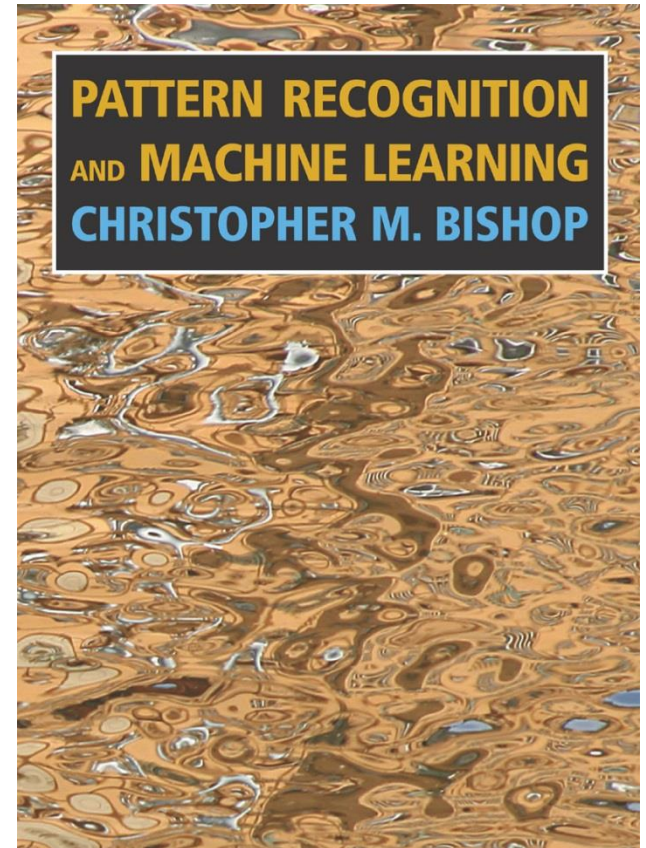
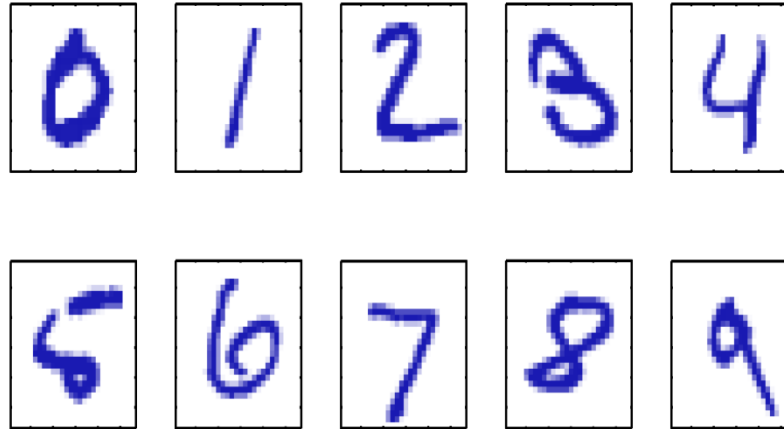


Machine-learning Crash Course

- This is not a course, people are expected to contribute (read, discuss, propose). Please bear in mind that we can be wrong!
- The initial idea is to follow Bishop's book, but there are other important references (e.g. Deep-learning book by Goodfellow).
- We thought of meeting twice a week 1 hour each session.
- Apart from discussing the concepts we (all) deem interesting, we could: do exercises, program, discuss relevant papers... Suggestions?



MNIST dataset



Automatically finding **regularities** in a dataset through the use of computer algorithms

Supervised learning: the training data comprises examples of the input vectors along with their corresponding target:

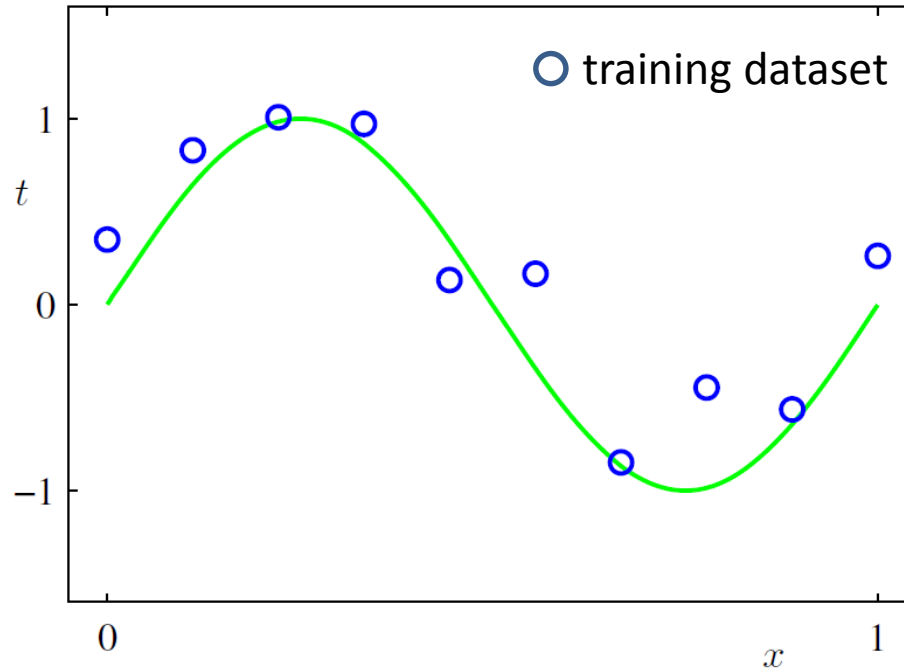
- Digit recognition (classification).
- Predicting price from size of a house (regression).

Unsupervised learning: there is not target (no ground truth), which makes the performance harder to evaluate:

- Clustering.
- Density estimation.
- Visualization.

Reinforcement learning: find the optimal policy (set of actions) to maximize the reward. Not cover in the book (and, probably, not cover by this course).

Example: polynomial curve fitting



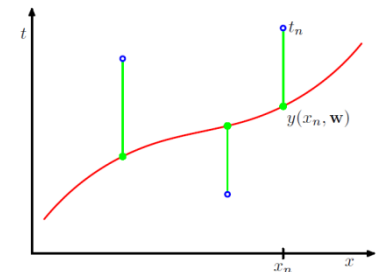
Our (linear) model:

$$y(x, \mathbf{w}) = \underbrace{w_0}_{\text{parameters}} + \underbrace{w_1}_{\text{parameters}}x + \underbrace{w_2}_{\text{parameters}}x^2 + \dots + \underbrace{w_M}_{\text{parameters}}x^{\underbrace{M}_{\text{hyperparameter(s)}}} = \sum_{j=0}^M w_j x^j$$

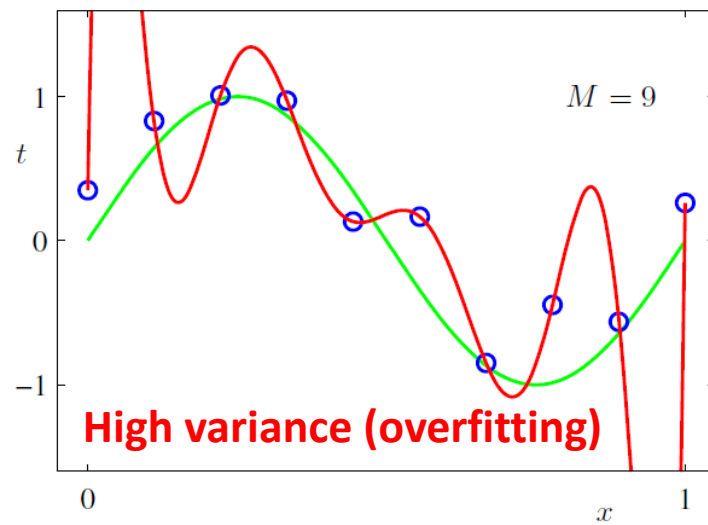
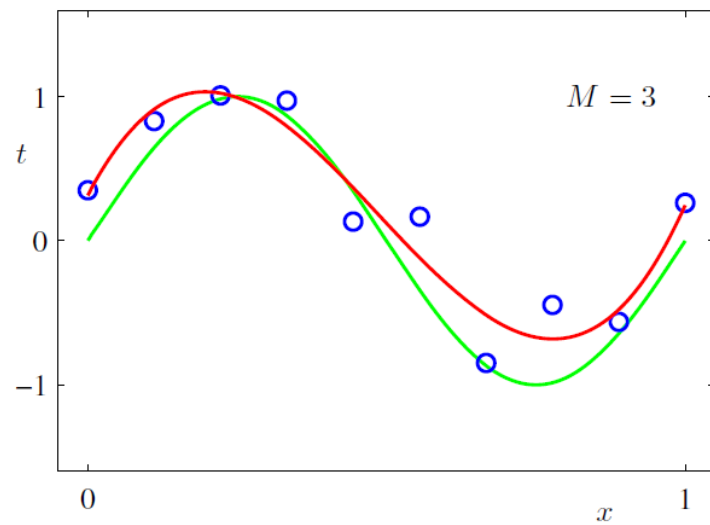
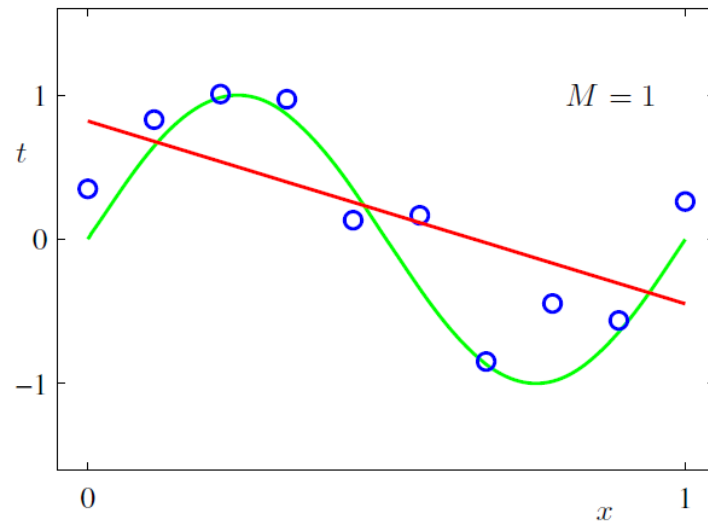
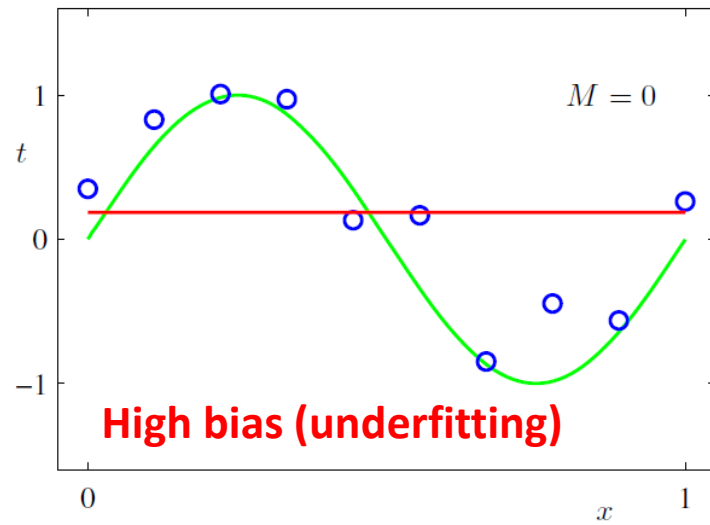
Our error function (sum of squares):

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

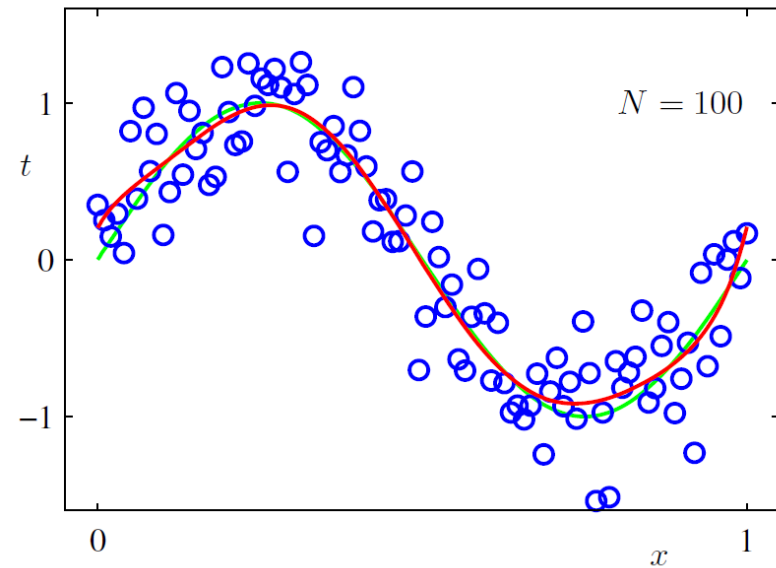
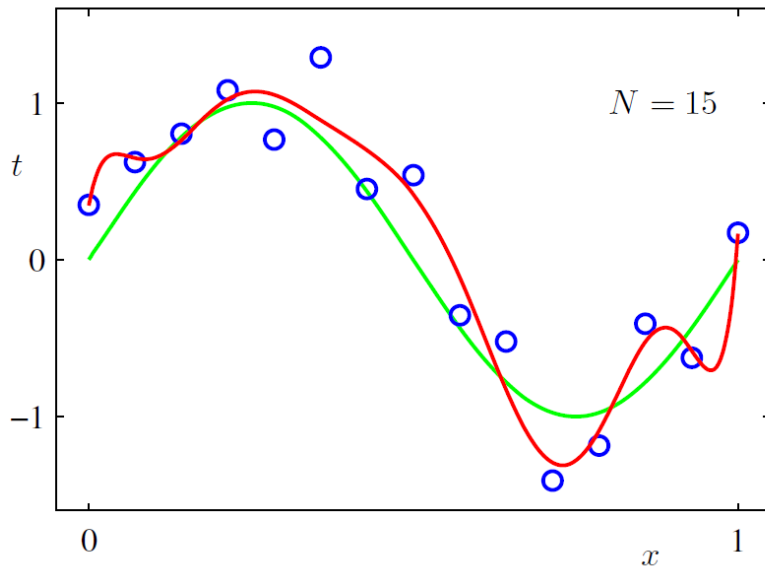
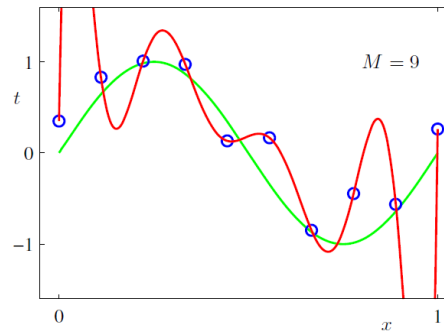
○ parameters
○ hyperparameter(s)



Model Selection (the bias-variance tradeoff)



The bias-variance tradeoff depends (a lot) on the amount of data.

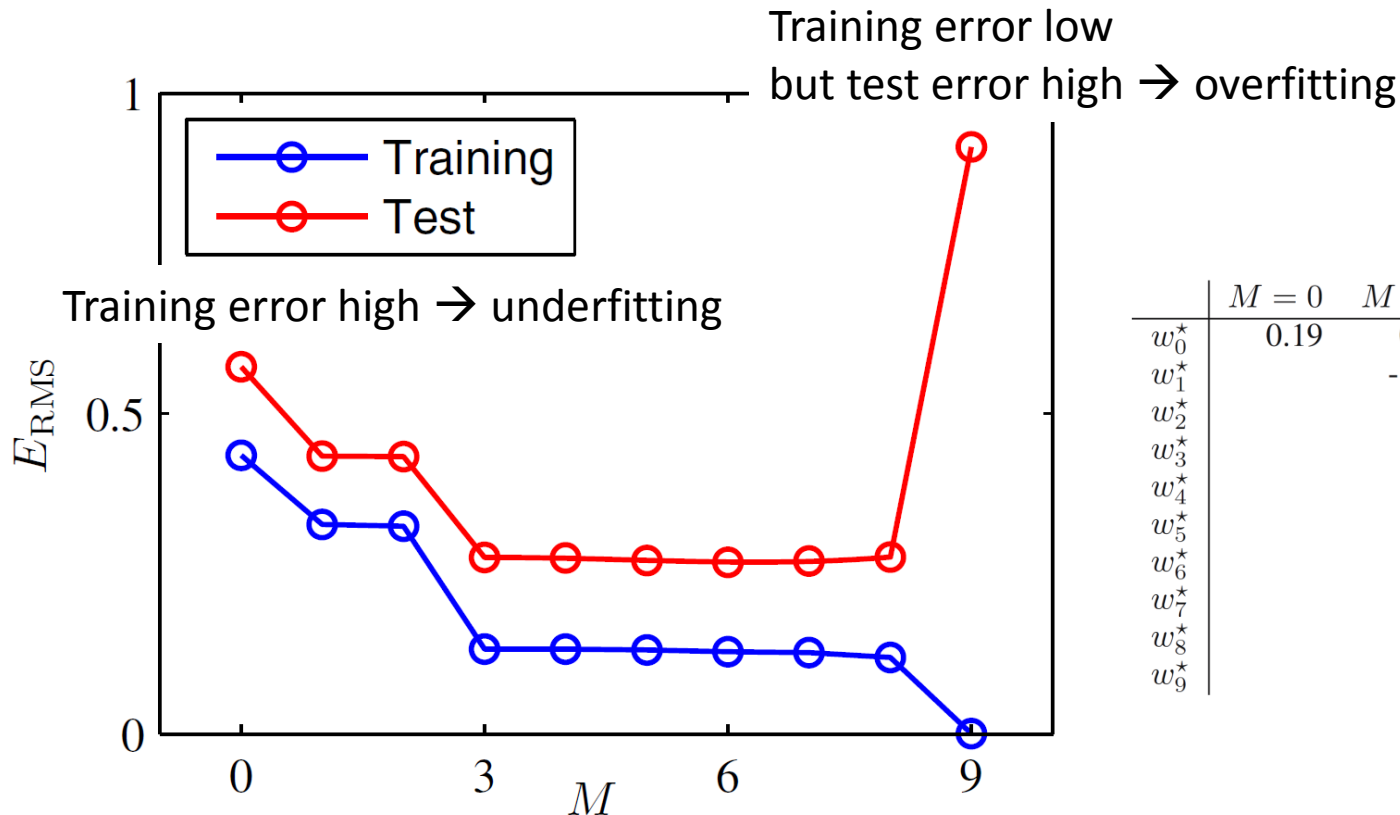


The larger the training dataset, the more flexible the model can be without risking a high variance issue

A way to control for overfitting:

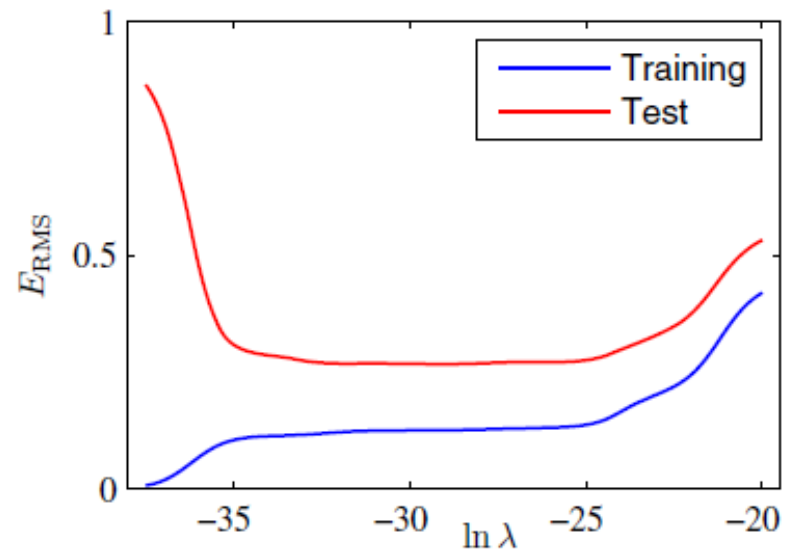
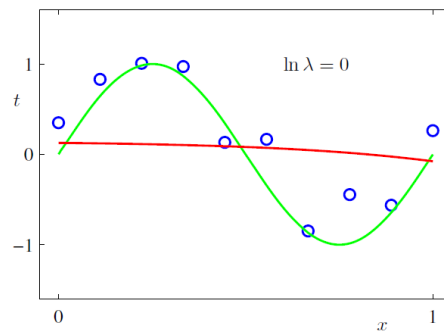
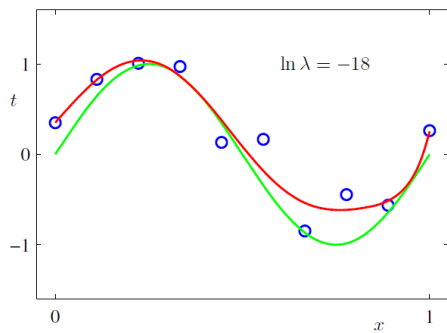
1. Divide dataset into *training* and *test* sets.
2. Train the model on the training dataset.
3. Test the model on the test dataset.

(note that this limits the amount of data available to train the model)



A way to avoid overfitting: regularization.

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



Probability Theory

The Rules of Probability

sum rule

$$p(X) = \sum_Y p(X, Y)$$

product rule

$$p(X, Y) = p(Y|X)p(X).$$

Bayesian approach

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

we would like to address and quantify the **uncertainty** that surrounds the appropriate choice for the model parameters **w**

posterior

$$p(\mathbf{w}|\mathcal{D}) = \frac{\overset{\text{likelihood}}{p(\mathcal{D}|\mathbf{w})} \overset{\text{prior}}{p(\mathbf{w})}}{p(\mathcal{D})}$$
$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) d\mathbf{w}$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

posterior

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

likelihood prior

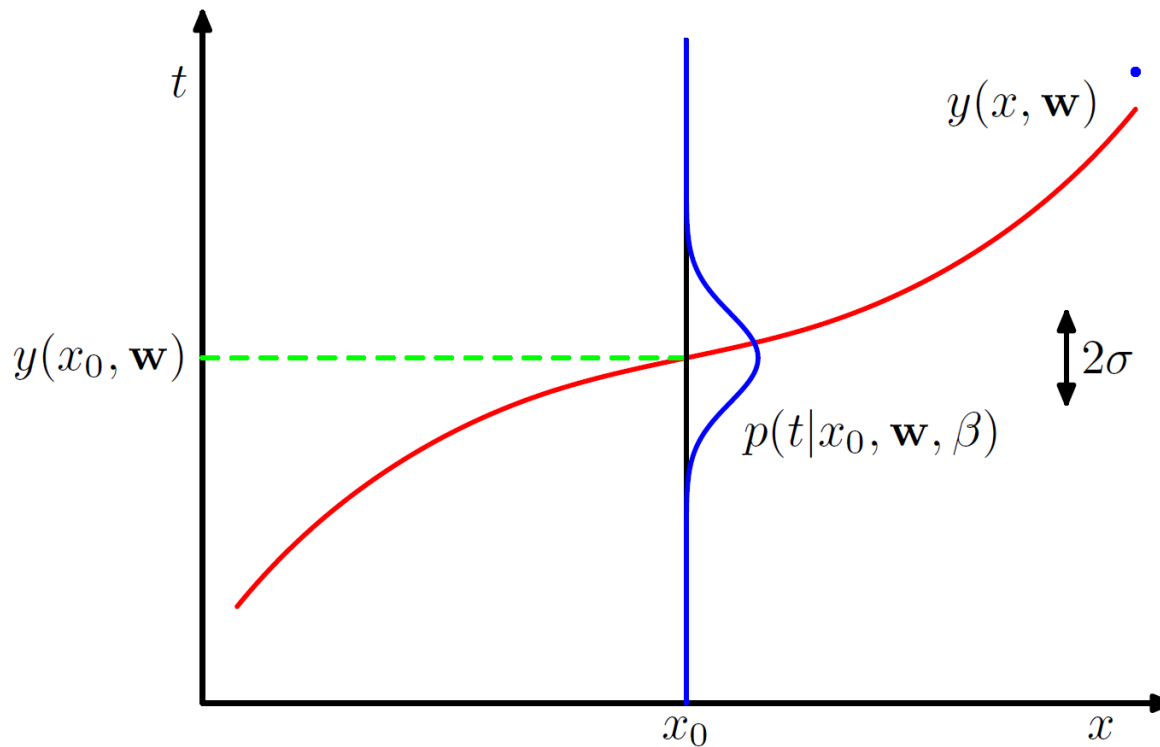
Maximum likelihood (frequentist approach):

- \mathbf{w} is set to the value that maximizes the likelihood function $p(\mathcal{D}|\mathbf{w})$.
- This corresponds to choosing the value of \mathbf{w} for which the probability of the observed data set is maximized.
- In the machine learning literature, the negative log of the likelihood function is called an *error function*.

A more probabilistic approach to curve fitting:

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

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



training data $\{\mathbf{x}, \mathbf{t}\}$

Likelihood function:

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

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

The sum-of-squares error function arises as a consequence of **maximizing likelihood under the assumption of a Gaussian noise distribution**.

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

predictive distribution that gives the probability distribution over t , rather than simply a point estimate

Maximum a Posteriori (MAP):

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

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

hyperparameter

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

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

maximizing the posterior distribution is equivalent to minimizing the regularized sum-of-squares error function

Full Bayesian approach:

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

