

Machine Learning I

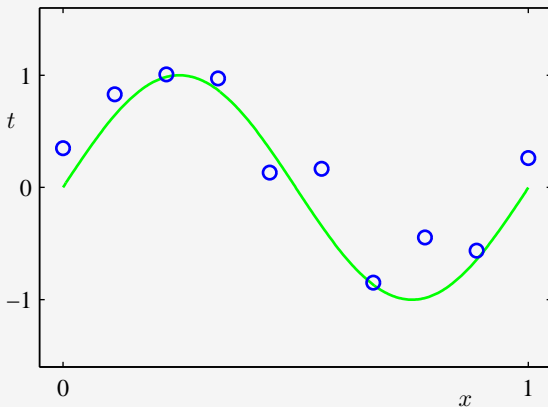
Week 2: Regression, Model Selection, Regularization and Validation

Martin Felder, Christian Osendorfer

Technische Universität München

29/30. October 2009

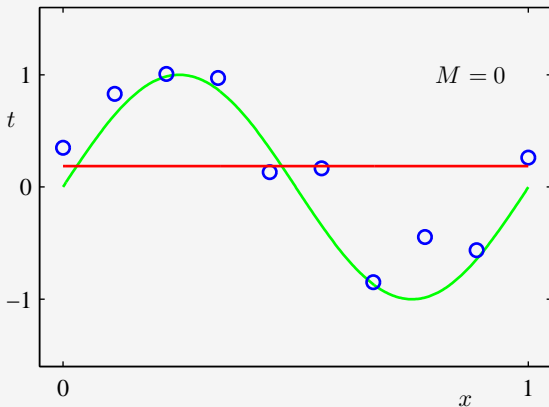
Consider noisy real-valued functions



$$\text{inputs: } \mathbf{X} = (x_1, \dots, x_N)^\top \quad (1)$$

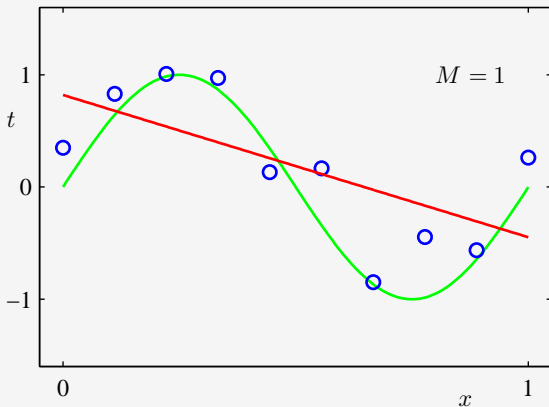
$$\text{targets: } \mathbf{T} = (t_1, \dots, t_N)^\top, \quad t_i = h(x_i) + \epsilon = \sin(2\pi x_i) + \epsilon \quad (2)$$

Model: 0th order polynomial



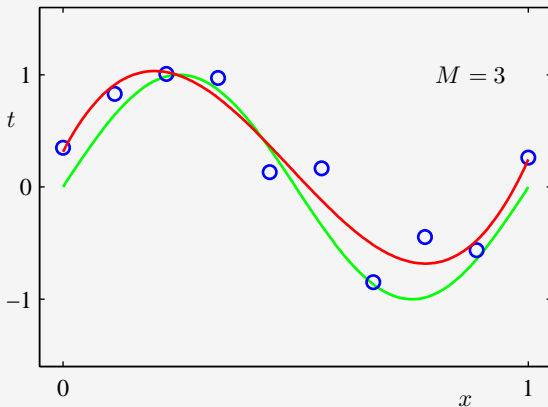
$$y(x, \mathbf{w}) = w_0$$

Model: 1st order polynomial



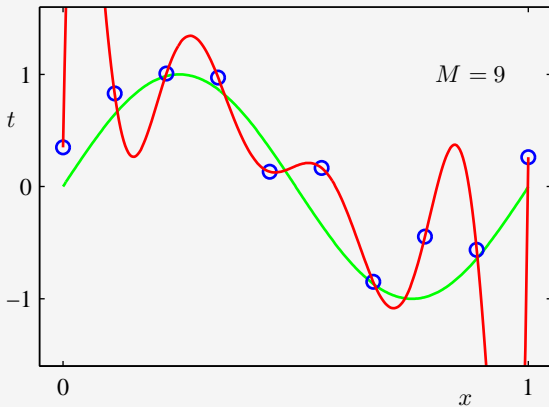
$$y(x, \mathbf{w}) = w_0 + w_1 x$$

Model: 3rd order polynomial



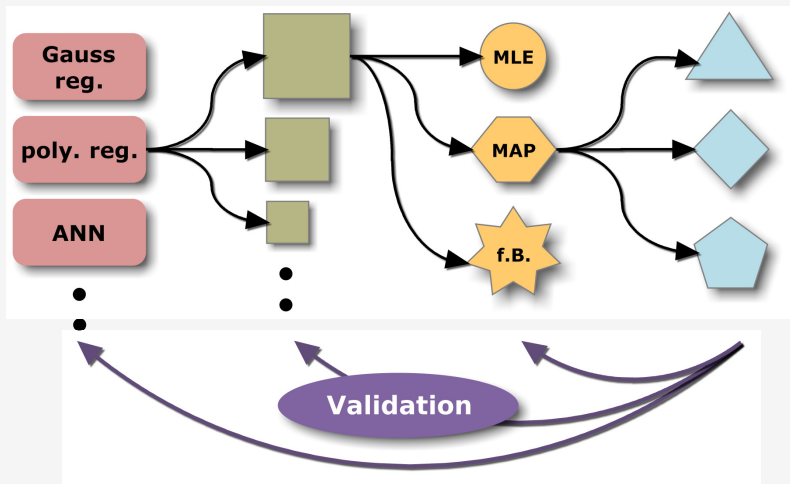
$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3$$

Model: 9th order polynomial



$$y(x, \mathbf{w}) = \sum_{j=0}^M w_j x^j$$

Model class Complexity Principle Optimizer



Problem Definition

Have: Input vector \mathbf{x} , model y with parameters \mathbf{w}

Question: What is output value t ?

\implies construct model from M functions $\phi(\mathbf{x})$:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad (3)$$

where

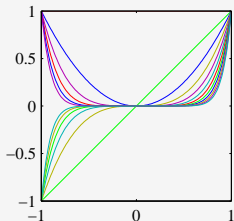
ϕ **basis function** – many choices, can be nonlinear

w_0 **bias** – equivalent to defining $\phi_0 \equiv 1$

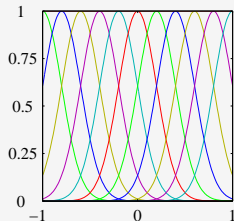
\implies still **linear** in \mathbf{w}

\implies compare to Taylor expansion, Fourier transform, wavelets...

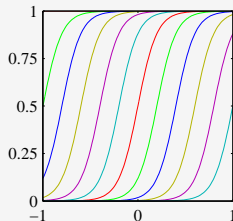
Typical Basis Functions



polynomials



Gaussians



“sigmoids”
(=S-shaped curves)

Now assume we measure random variable t as

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad [\epsilon: \text{Gaussian, zero mean}] \quad (4)$$

$$\Rightarrow \mathcal{E}[t|\mathbf{x}] = y(\mathbf{x}, \mathbf{w}) \quad \text{for least squares loss fct.}$$

Let dataset \mathcal{D} consist of a series of observations $\mathbf{T} = (t_1, t_2, \dots, t_N)$ and corresponding input vectors $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$. Then the **likelihood function** is

$$p(\mathbf{T}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad [\beta^{-1}: \text{precision}] \quad (5)$$

As usual, it's easier to deal with logarithmized probabilities (leave out \mathbf{X} for brevity):

$$\ln p(\mathbf{T}|\mathbf{w}, \beta) = \frac{N}{2} \ln \beta - \underbrace{\frac{N}{2} \ln(2\pi)}_{const.} - \beta \underbrace{\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{E_{\mathcal{D}}(\mathbf{w})} \quad (6)$$

Maximum Likelihood Solution

Now find max likelihood by setting gradient w.r.t. \mathbf{w} to zero:

$$\nabla_{\mathbf{w}} \ln p(\mathbf{T} | \mathbf{w}, \beta) = \sum^n (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) \quad (7)$$

$$\Rightarrow \mathbf{w}_{\text{ML}} = \underbrace{(\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}}_{=\Phi^\dagger} \quad \leftarrow \text{normal equations of least sq. prob.} \quad (8)$$

Φ^\dagger is called **Moore-Penrose pseudo-inverse** of Φ (because for a square matrix, $\Phi^\dagger = \Phi^{-1}$). Furthermore,

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & & \vdots \\ \vdots & \vdots & \ddots & \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} = \text{design matrix of } \phi.$$

What else can we say about the ML solution?

1.) Usually, $\phi_0(\mathbf{x}) \equiv 1$ (bias!) \hookrightarrow pull w_0 out of the sum \hookrightarrow solve $\nabla \ln p = 0$ again \hookrightarrow

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \quad (9)$$

\Rightarrow bias compensates for difference between (training set) **averages** of targets and weighted linear reconstruction.

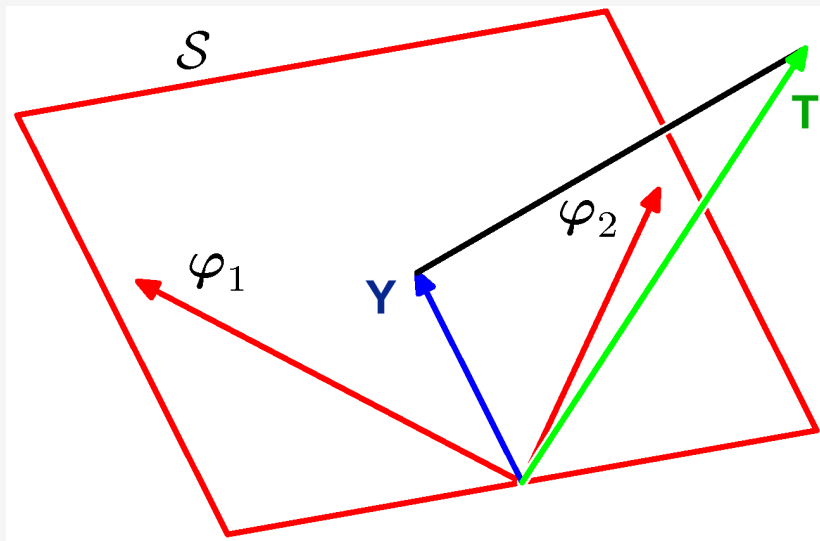
2.) Maximize likelihood w.r.t. $\beta \hookrightarrow$

$$\beta_{\text{ML}}^{-1} = \frac{1}{N} \sum^n (t_n - \mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_n))^2 \quad (10)$$

= inverse variance for **residual scatter** of weighted linear reconstruction around targets.

\Rightarrow if $\mathbf{w}_{\text{ML}}^T \phi$ is perfect fit to generating fct. $y(\mathbf{x}, \mathbf{w})$, then $\beta_{\text{ML}} \rightarrow$ inverse variance of error ϵ .

Geometrical Interpretation



The Bias-Variance Dilemma

We will follow the derivation of Bishop (1995) here. Consider again the quadratic error function, and let $N \rightarrow \infty$:

$$E_{\mathcal{D}} = \lim_{N \rightarrow \infty} \frac{1}{2N} \sum_n (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 \quad (11)$$

$$= \frac{1}{2} \int (y(\mathbf{x}, \mathbf{w}) - t)^2 p(t, \mathbf{x}) dt d\mathbf{x} \quad (12)$$

Now define

$$\langle t | \mathbf{x} \rangle \equiv \int t p(t | \mathbf{x}) dt \quad (13)$$

$$\langle t^2 | \mathbf{x} \rangle \equiv \int t^2 p(t | \mathbf{x}) dt. \quad (14)$$

The Bias-Variance Dilemma

By using the product rule, and including a zero expressed as $+\langle t|x \rangle - \langle t|x \rangle$ into the quadratic term, it can be shown that

$$\begin{aligned} E_{\mathcal{D}} &= \frac{1}{2} \int (y(\mathbf{x}, \mathbf{w}) - \langle t|x \rangle)^2 p(x) dx \\ &\quad + \frac{1}{2} \int (\langle t^2|x \rangle - \langle t|x \rangle^2) p(x) dx \end{aligned} \quad (15)$$

Only the first integral is relevant for minimization of $E_{\mathcal{D}}$, the second one merely sets a lower limit to the achievable error. Note that the absolute minimum \mathbf{w}^* of $E_{\mathcal{D}}$ is given by

$$y(x, \mathbf{w}^*) = \langle t|x \rangle, \quad (16)$$

i.e. **the average model output is equal to the average of the training data**. In practice, this state will be reached first where the training data density $p(x)$ is highest.

The Bias-Variance Dilemma

However, the size of $(y(\mathbf{x}, \mathbf{w}) - \langle t|x \rangle)^2$ depends on the exact choice of training data \mathcal{D} – there are infinitely many realizations! To keep the notation uncluttered, define

$$y_{\mathcal{D}} \equiv y(\mathbf{x}, \mathbf{w}|\mathcal{D}). \quad (17)$$

This dependence can be eliminated by considering the **ensemble average** $\mathcal{E}[\cdot] = \mathbb{E}_{\mathcal{D}}[\cdot]$ over many different training data sets \mathcal{D} :

$$\mathcal{E} [(y_{\mathcal{D}} - \langle t|x \rangle)^2] \quad (18)$$

Again, ideally this would be zero, but in practice no model $y_{\mathcal{D}}$ will perform equally well on all different training data sets, since it does not contain the full information of the infinite generator.

The Bias-Variance Dilemma

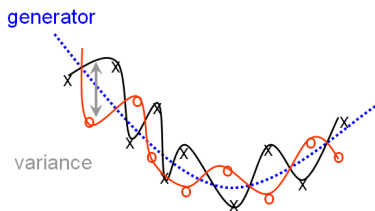
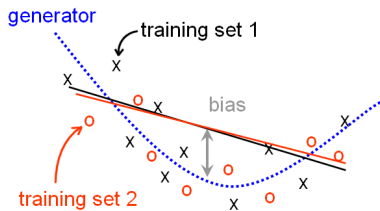
By introducing another zero made from $+\mathcal{E}[y_D] - \mathcal{E}[y_D]$, one obtains

$$\mathcal{E} [(y_D - \langle t|x \rangle)^2] = \underbrace{(\mathcal{E}[y_D] - \langle t|x \rangle)^2}_{\text{bias}^2} + \underbrace{\mathcal{E} [(y_D - \mathcal{E}[y_D])^2]}_{\text{variance}} \quad (19)$$

Bias = model pays not enough heed to the data, or is insufficiently complex.

Variance = model too strongly data-driven, cannot detect underlying generator.

$$\underbrace{(\mathcal{E}[y_D] - \langle t|x \rangle)^2}_{\text{bias}^2} + \underbrace{\mathcal{E}[(y_D - \mathcal{E}[y_D])^2]}_{\text{variance}}$$



→ **Regularization** ←
mediates between these extremes.

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
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

Regularization

MLE often suffers from overfitting \hookrightarrow use **regularization**:

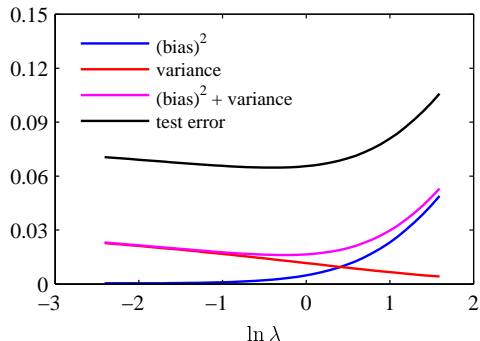
$$E_{\mathcal{D}} = \frac{1}{2} \sum^n (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum^M |w_j|^q \quad (20)$$

\Rightarrow this is like a Lagrange term specifying an additional constraint:

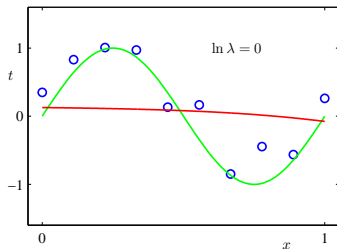
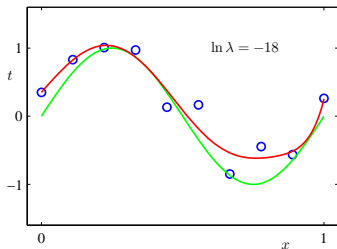
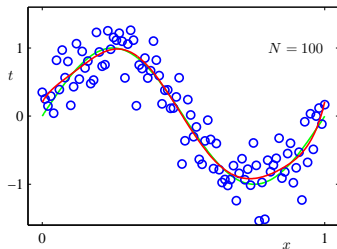
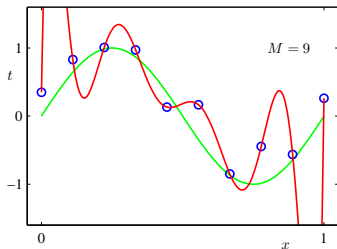
$$\sum^M |w_j|^q \leq \varepsilon$$

\hookrightarrow most often, use quadratic regularizer: $q = 2$, i.e.

$$\sum^M |w_j|^q = \mathbf{w}^T \mathbf{w}$$



Regularization and Model Complexity



MAP Solution

Recap: For the coin flip experiment, we introduced prior information to prevent overfitting. By analogy:

	train data	likelihood	prior	posterior
coin:	$\mathcal{D} = \mathbf{X}$	$p(\mathcal{D} \theta)$	$p(\theta a, b)$	$p(\theta \mathcal{D})$
regr.:	$\mathcal{D} = \{\mathbf{X}, \mathbf{T}\}$	$p(\mathbf{T} \mathbf{X}, \mathbf{w}, \beta)$	$p(\mathbf{w} \cdot)$	$p(\mathbf{w} \mathbf{X}, \mathbf{T}, \cdot)$

Note: As mentioned earlier, \mathbf{X} is usually dropped for clarity.

Prior: How to find a good one?

- recall (from Eq. 5) that our likelihood function $p(\mathbf{T}|\mathbf{w}, \beta)$ is a Gaussian,
- treat precision $\beta = 1/\sigma^2$ as a known parameter (for now),
- know that the **conjugate prior** for a Gaussian with known variance is also a Gaussian.

MAP Solution

Hence introduce the following prior:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \quad [\mathbf{m}_0: \text{mean}, \mathbf{S}_0: \text{covariance}] \quad (21)$$

Often we don't know much about the prior distribution anyway. For a suitably designed model with independent parameters \mathbf{w} , the following prior is usually reasonable:

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0 = 0, \mathbf{S}_0 = \alpha^{-1} \mathbf{I}) \quad (22)$$

This results in the posterior:

$$p(\mathbf{w} | \mathbf{T}, \alpha, \beta) \propto p(\mathbf{T} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) \quad (23)$$

Applying the log yields again an error function E_{MAP} to minimize:

$$\ln p(\mathbf{w} | \mathbf{T}, \alpha, \beta) = -\frac{\beta}{2} \sum^n (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + \text{const.} \quad (24)$$

\Rightarrow just like squared error minimization with a quadratic regularizer weighted with $\lambda = \frac{\alpha}{\beta}$ (Eq. 20).

Posterior Distribution

We can actually find a closed expression for the posterior!

Posterior parameter distribution

$$p(\mathbf{w}|\mathbf{T}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (25)$$

with

$$\mathbf{m}_N = \mathbf{S}_N \left(\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{T} \right) \quad (26)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi \quad (27)$$

Properties of the posterior:

- Since we again have a Gaussian, the maximum posterior solution equals the mode: $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$.
- In the limit of an infinitely broad prior, $\mathbf{S}_0^{-1} \rightarrow 0$, therefore $\mathbf{w}_{\text{MAP}} \rightarrow \mathbf{w}_{\text{ML}} = \Phi^\dagger \mathbf{T}$ (Eq. 8).
- For $N = 0$, i.e. no data points, we get the prior back.

Posterior Distribution for a Simple Prior

If we look at our simplified case:

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0 = 0, \mathbf{S}_0 = \alpha^{-1} \mathbf{I}) \quad (28)$$

the posterior parameters simplify to:

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{T} \quad (29)$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi \quad (30)$$

A simple example

Bayesian regression for the target values

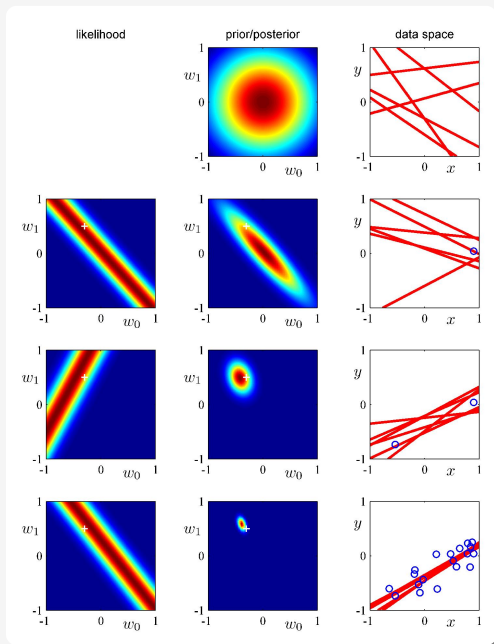
$$t_n = -0.3 + 0.5x_n + \epsilon$$

where ϵ is a Gaussian noise term ($\sigma = 0.2$).

To model this, we set $\phi(x) = \begin{bmatrix} 1 \\ x \end{bmatrix}$ and thus

$$y(x, \mathbf{w}) = w_0 + w_1 x$$

Sequential Estimation: The demo shows how the posterior's breadth gets smaller as more and more points t are taken into account, and how its mode converges to the optimum (=correct) values of the weights (white cross).



Predictive Distribution

Usually, we want to know output t for new values of \mathbf{x} – the model parameters \mathbf{w} are just a means to achieve this. To predict t , evaluate

$$p(t|\mathbf{x}, \mathbf{T}, \alpha, \beta) = \int \underbrace{p(t|\mathbf{x}, \mathbf{w}, \beta)}_{\text{likelihood (5)}} \underbrace{p(\mathbf{w}|\mathbf{T}, \alpha, \beta)}_{\text{posterior (25)}} d\mathbf{w} \quad (31)$$

(coin flip analogy: $p(x|\mathcal{D}, a, b) = \int_0^1 p(x|\theta)p(\theta|\mathcal{D}, a, b) d\theta$)

Predictive distribution

$$p(t|\mathbf{x}, \mathbf{T}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T\phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad (32)$$

with variance

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}). \quad (33)$$

Example (using 9 radial basis functions)

