

Linear Models for Regression

Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea
seungjin@postech.ac.kr

Outline

- ▶ Regression
- ▶ Least Squares Method
- ▶ Ridge Regression
- ▶ Least Mean Squares
- ▶ Recursive (Sequential) least squares
- ▶ Bias-Variance Dilemma
- ▶ Bayesian Linear Regression

Regression

- ▶ **Regression** aims at modeling the dependence of a **response** Y on a **covariate** X . In other words, the goal of regression is to predict the value of one or more continuous target variables y given the value of input vector \mathbf{x} .
- ▶ The regression model is described by

$$y = f(\mathbf{x}) + \epsilon.$$

- ▶ Terminology:
 - ▶ \mathbf{x} : **input, independent variable, predictor, regressor, covariate**
 - ▶ y : **output, dependent variable, response**
- ▶ The dependence of a response on a covariate is captured via a conditional probability distribution, $p(y|\mathbf{x})$.
- ▶ Depending on $f(\mathbf{x})$,
 - ▶ **Linear regression**: $f(\mathbf{x}) = \sum_{j=1}^M w_j \phi_j(\mathbf{x}) + w_0$.
 - ▶ **Kernel regression**: $f(\mathbf{x}) = \sum_{i=1}^N w_i k(\mathbf{x}, \mathbf{x}_i) + w_0$.

Regression Function: Conditional Mean

We consider the mean squared error and find the MMSE estimate:

$$\begin{aligned}\mathcal{E}(f) &= \langle \|y - f(x)\|^2 \rangle \\ &= \int \int \|y - f(x)\|^2 p(x, y) dx dy \\ &= \int \int \|y - f(x)\|^2 p(x) p(y|x) dx dy \\ &= \int p(x) \left[\underbrace{\int \|y - f(x)\|^2 p(y|x) dy}_{\text{to be minimized}} \right] dx\end{aligned}$$

$$\frac{\partial}{\partial f(x)} \left[\int \|y - f(x)\|^2 p(y|x) dy \right] = 0 \Rightarrow \boxed{f(x) = \int y p(y|x) dy = \langle y|x \rangle}.$$

Linear Regression

Linear regression refers to a model in which the conditional mean of y given the value of \mathbf{x} is an **affine function** of $\phi(\mathbf{x})$

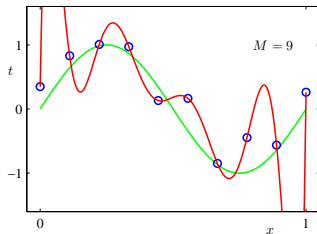
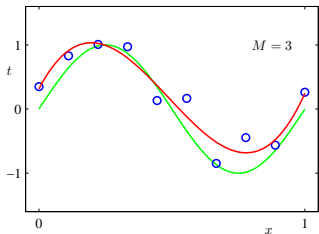
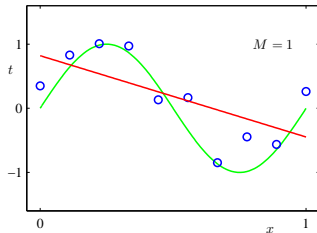
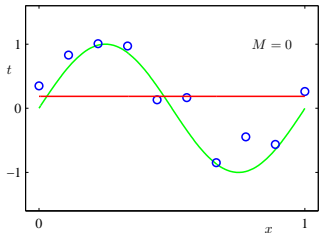
$$f(\mathbf{x}) = \sum_{j=1}^M w_j \phi_j(\mathbf{x}) + w_0 \phi_0(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}),$$

where $\phi_j(\mathbf{x})$ are known as **basis functions** and

$$\begin{aligned}\mathbf{w} &= [w_0, w_1, \dots, w_M]^\top, \\ \phi &= [\phi_0, \phi_1, \dots, \phi_M]^\top.\end{aligned}$$

By using nonlinear basis functions, we allow the function $f(\mathbf{x})$ to be a nonlinear function of the input vector \mathbf{x} (but a linear function of $\phi(\mathbf{x})$).

Polynomial Regression: $y_t = \sum_{j=0}^M w_j \phi_j(x_t) = \sum_{j=0}^M w_j x_t^j$



Basis Functions

- ▶ **Polynomial regression:** $\phi_j(\mathbf{x}) = x^j$.
- ▶ **Gaussian basis functions:** $\phi_j(\mathbf{x}) = \exp\left\{-\frac{\|\mathbf{x}-\boldsymbol{\mu}_j\|^2}{2\sigma^2}\right\}$.
- ▶ **Spline basis functions:** Piecewise polynomials (divide the input space up into regions and fit a different polynomial in each region).
- ▶ Many other possible basis functions: sigmoidal basis functions, hyperbolic tangent basis functions, Fourier basis, wavelet basis, and so on.

Least Squares Method

Given a set of training data $\{(\mathbf{x}_t, y_t)\}_{t=1}^N$, we determine the weight vector \mathbf{w} which minimizes

$$\mathcal{E}_{LS} = \frac{1}{2} \sum_{t=1}^N \{y_t - \mathbf{w}^\top \phi(\mathbf{x}_t)\}^2 = \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{w}\|^2,$$

where $\mathbf{y} = [y_1, \dots, y_N]^\top$ and $\Phi \in \mathbb{R}^{N \times (M+1)}$ is known as the **design matrix** with $\Phi_{tj} = \phi_j(\mathbf{x}_t)$, i.e.,

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_M(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_M(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_M(\mathbf{x}_N) \end{bmatrix}.$$

$\frac{\partial \mathcal{E}_{LS}}{\partial \mathbf{w}} = 0$ leads to the **normal equation** that is of the form

$$\Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{y}.$$

Thus, LS estimate of \mathbf{w} is given by

$$\mathbf{w}_{LS} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y} = \Phi^\dagger \mathbf{y},$$

where Φ^\dagger is known as the **Moore-Penrose pseudo-inverse**.
 $\Phi^T \Phi$ is known as the **Gram matrix**.

Maximum Likelihood

We assume that the target variable y_t is given by a deterministic function $f(\mathbf{x}_t, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x}_t)$ with additive Gaussian noise so that

$$\mathbf{y} = \Phi \mathbf{w} + \epsilon,$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.

The log-likelihood is given by

$$\begin{aligned} \mathcal{L} &= \log p(\mathbf{y} | \Phi, \mathbf{w}) = \sum_{t=1}^N \log p(y_t | \phi(\mathbf{x}_t), \mathbf{w}) \\ &= -\frac{N}{2} \log \sigma^2 - \frac{N}{2} \log 2\pi - \sigma^{-2} \mathcal{E}_{LS}. \end{aligned}$$

Therefore, under **Gaussian noise assumption**, $\mathbf{w}_{ML} = \mathbf{w}_{LS}$.

Ridge Regression

We consider the sum-of-squares error function with a Euclidean norm-based regularizer

$$\mathcal{E} = \underbrace{\frac{1}{2} \|\mathbf{y} - \Phi \mathbf{w}\|^2}_{\text{LSfit}} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{regularizer}}.$$

Solving $\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = 0$ for \mathbf{w} leads to

$$\mathbf{w}_{\text{ridge}} = \left(\lambda I + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y}.$$

Ridge Regression: MAP Perspective

Recall the likelihood:

$$p(\mathbf{y}|\Phi, \mathbf{w}) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I}).$$

Assume a zero-mean Gaussian prior with covariance Σ for parameters \mathbf{w} :

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, \Sigma).$$

Then the posterior over \mathbf{w} ,

$$p(\mathbf{w}|\mathbf{y}, \Phi) = \frac{p(\mathbf{y}|\Phi, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\Phi, \mathbf{w})p(\mathbf{w})d\mathbf{w}},$$

is still Gaussian with mean and mode at

$$\hat{\mathbf{w}} = (\sigma^2\Sigma^{-1} + \Phi^T\Phi)^{-1}\Phi^T\mathbf{y}.$$

When Σ is proportional to identity, i.e., $\Sigma = \lambda^{-1}\sigma^2\mathbf{I}$, this is called [ridge regression](#).

Ridge Regression: Shrinkage

- ▶ Suppose that the SVD of the design matrix $\Phi \in \mathbb{R}^{N \times (M+1)}$ has the form

$$\Phi = \mathbf{U}\mathbf{D}\mathbf{V}^\top, \quad \text{where } \mathbf{D} = \text{diag}(d_1, \dots, d_{M+1}).$$

- ▶ Using the SVD we write the least squares fitted vector as

$$\Phi \mathbf{w}_{LS} = \Phi \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{y} = \mathbf{U}\mathbf{U}^\top \mathbf{y} = \sum_{i=1}^{M+1} \mathbf{u}_i \mathbf{u}_i^\top \mathbf{y},$$

where $\mathbf{U}^\top \mathbf{y}$ are the **coordinates** of \mathbf{y} with respect to the orthonormal basis \mathbf{U} .

- ▶ The ridge solutions are

$$\begin{aligned} \Phi \mathbf{w}_{ridge} &= \Phi \left(\lambda \mathbf{I} + \Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{y} \\ &= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D}\mathbf{U}^\top \mathbf{y} = \sum_{i=1}^{M+1} \mathbf{u}_i \frac{d_i^2}{d_i^2 + \lambda} \mathbf{u}_i^\top \mathbf{y}. \end{aligned}$$

It shrinks the coordinates by the factors $\frac{d_i^2}{d_i^2 + \lambda}$. A greater amount of shrinkage is applied to basis vectors with smaller d_i^2 .

Least Mean Square (LMS)

LMS is a gradient-descent method which minimizes the **instantaneous error** \mathcal{E}_t , where

$$\mathcal{E}_{LS} = \sum_{t=1}^N \mathcal{E}_t = \frac{1}{2} \sum_{t=1}^N \left\{ y_t - \mathbf{w}^\top \phi(\mathbf{x}_t) \right\}^2.$$

The gradient descent method leads to the updating rule for \mathbf{w} that is of the form

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \eta \nabla \mathcal{E}_t \\ &\leftarrow \mathbf{w} + \eta \left\{ y_t - \mathbf{w}^\top \phi(\mathbf{x}_t) \right\} \phi(\mathbf{x}_t), \end{aligned}$$

where $\eta > 0$ is a constant known as **learning rate**.

Recursive (Sequential) LS

We introduce the **forgetting factor** λ to de-emphasize old samples, leading to the following error function

$$\mathcal{E}_{RLS} = \frac{1}{2} \sum_{i=1}^t \lambda^{t-i} (y_i - \phi_i \mathbf{w}_t^\top)^2,$$

where $\phi_t = \phi(\mathbf{x}_t)$.

Solving $\frac{\partial \mathcal{E}_{RLS}}{\partial \mathbf{w}_t} = 0$ for \mathbf{w}_t leads to

$$\left[\sum_{i=1}^t \lambda^{t-i} \phi_i \phi_i^\top \right] \mathbf{w}_t = \left[\sum_{i=1}^t \lambda^{t-i} y_i \phi_i \right].$$

We define

$$\mathbf{P}_t = \left[\sum_{i=1}^t \lambda^{t-i} \phi_i \phi_i^\top \right]^{-1},$$
$$\mathbf{r}_t = \left[\sum_{i=1}^t \lambda^{t-i} y_i \phi_i \right].$$

With these definitions, we have

$$\mathbf{w}_t = \mathbf{P}_t \mathbf{r}_t.$$

The core idea of RLS is to apply the [matrix inversion lemma](#) to develop the [sequential algorithm](#) without matrix inversion.

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{DA}^{-1} \mathbf{B} + \mathbf{C}^{-1})^{-1} \mathbf{DA}^{-1}.$$

The recursion for \mathbf{P}_t is given by

$$\begin{aligned}\mathbf{P}_t &= \left[\sum_{i=1}^t \lambda^{t-i} \phi_i \phi_i^\top \right]^{-1} \\ &= \left[\lambda \sum_{i=1}^{t-1} \lambda^{t-1-i} \phi_i \phi_i^\top + \phi_t \phi_t^\top \right]^{-1} \\ &= \frac{1}{\lambda} \left[\mathbf{P}_{t-1} - \frac{\mathbf{P}_{t-1} \phi_t \phi_t^\top \mathbf{P}_{t-1}}{\lambda + \phi_t^\top \mathbf{P}_{t-1} \phi_t} \right]. \quad (\text{matrix inversion lemma})\end{aligned}$$

Thus, the updating rule for \mathbf{w} is given by

$$\begin{aligned}\mathbf{w}_t &= \mathbf{P}_t \mathbf{r}_t \\ &= \frac{1}{\lambda} \left[\mathbf{P}_{t-1} - \frac{\mathbf{P}_{t-1} \phi_t \phi_t^\top \mathbf{P}_{t-1}}{\lambda + \phi_t^\top \mathbf{P}_{t-1} \phi_t} \right] [\lambda \mathbf{r}_{t-1} + y_t \phi_t] \\ &= \mathbf{w}_{t-1} + \underbrace{\frac{\mathbf{P}_{t-1} \phi_t}{\lambda + \phi_t^\top \mathbf{P}_{t-1} \phi_t}}_{\text{gain}} \underbrace{\left[y_t - \phi_t^\top \mathbf{w}_{t-1} \right]}_{\text{error}}.\end{aligned}$$

Expected Loss

The squared loss L is given by

$$L(y, f(\mathbf{x})) = (f(\mathbf{x}) - y)^2.$$

The expected loss is computed by

$$\begin{aligned}\mathbb{E}\{L(y, f(\mathbf{x}))\} &= \int \int (f(\mathbf{x}) - y)^2 p(\mathbf{x}, y) dx dy \\ &= \mathbb{E}\left\{(f(\mathbf{x}) - \mathbb{E}\{y|\mathbf{x}\} + \mathbb{E}\{y|\mathbf{x}\} - y)^2\right\} \\ &= \mathbb{E}\left\{(f(\mathbf{x}) - \mathbb{E}\{y|\mathbf{x}\})^2 + (\mathbb{E}\{y|\mathbf{x}\} - y)^2\right. \\ &\quad \left.+ \underbrace{2\mathbb{E}\{(f(\mathbf{x}) - \mathbb{E}\{y|\mathbf{x}\})(\mathbb{E}\{y|\mathbf{x}\} - y)\}}_0\right\}.\end{aligned}$$

The expected loss can be written as

$$\mathbb{E}\{L(y, f(\mathbf{x}))\} = \underbrace{\int (f(\mathbf{x}) - \mathbb{E}\{y|\mathbf{x}\})^2 p(\mathbf{x})d\mathbf{x}}_{\text{first term}} + \underbrace{\int (\mathbb{E}\{y|\mathbf{x}\} - y)^2 p(\mathbf{x})d\mathbf{x}}_{\text{second term}}.$$

- ▶ The function $f(\mathbf{x})$ appears only in the first term which will be minimized when $f(\mathbf{x}) = \mathbb{E}\{y|\mathbf{x}\}$.
- ▶ The second term is the variance of the distribution of y , averaged over \mathbf{x} , representing the intrinsic variability of the target data and can be regarded as noise. Because it is independent of $f(\mathbf{x})$, it represents the irreducible minimum value of the loss function.

Bias-Variance Decomposition

Consider the integrand of the first term, which for a particular data set \mathcal{D} takes the form $(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}\{y|\mathbf{x}\})^2$.

Then we have

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left\{ (f(\mathbf{x}; \mathcal{D}) - \mathbb{E}\{y|\mathbf{x}\})^2 \right\} &= \underbrace{(\mathbb{E}_{\mathcal{D}}\{f(\mathbf{x}; \mathcal{D})\} - \mathbb{E}\{y|\mathbf{x}\})^2}_{\text{(bias)}^2} \\ &+ \underbrace{\mathbb{E}_{\mathcal{D}} \left\{ (f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}\{f(\mathbf{x}; \mathcal{D})\})^2 \right\}}_{\text{variance}}. \end{aligned}$$

expected loss = (bias)² + variance + noise.

Bias-Variance Dilemma

There is a trade-off between bias and variance:

- ▶ flexible models: low bias but high variance
- ▶ rigid models: high bias but low variance

The model with the optimal predictive capability is the one that leads to the best balance between bias and variance.

Example

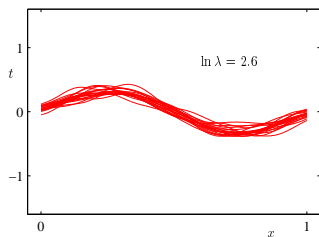
We consider the sinusoidal data set. We generate 100 data sets, each containing $N = 25$ data points, independently from the sinusoidal curve $h(x) = \mathbb{E}\{y|x\} = \sin(2\pi x)$. The data sets are indexed by $l = 1, \dots, L$ where $L = 100$. For each data set $\mathcal{D}^{(l)}$, we fit a model with 24 Gaussian basis functions through a method of ridge regression to give a prediction function $f^{(l)}(x)$.

$$\bar{f}(x) = \frac{1}{L} \sum_{l=1}^L f^{(l)}(x),$$

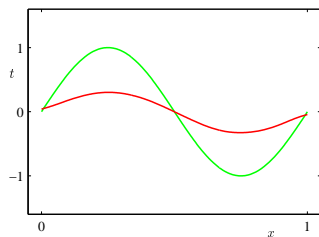
$$(\text{bias})^2 = \frac{1}{N} \sum_{t=1}^N [\bar{f}(x_t) - h(x_t)]^2,$$

$$\text{variance} = \frac{1}{N} \sum_{t=1}^N \frac{1}{L} \sum_{l=1}^L [f^{(l)}(x_t) - \bar{f}(x_t)]^2.$$

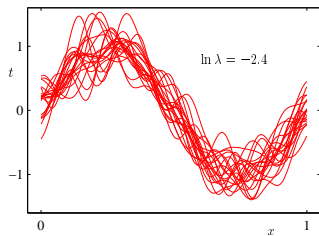
Rigid model (1st row) and flexible model (2nd row)



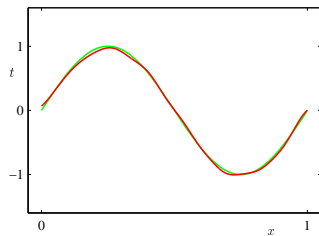
low variance



high bias



high variance



low bias

Bayesian Linear Regression

The **likelihood function** is given by

$$p(\mathbf{y} | \Phi, \mathbf{w}) = \prod_{t=1}^N p(y_t | \phi_t, \mathbf{w}),$$

where $p(y_t | \phi_t, \mathbf{w}) = \mathcal{N}(y_t | \mathbf{w}^\top \phi_t, \beta^{-1})$.

Assuming **Gaussian prior** for \mathbf{w} , i.e., $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mu_0, \Sigma_0)$, the **posterior distribution over \mathbf{w}** is again **Gaussian** that is of the form

$$p(\mathbf{w} | \mathbf{y}, \Phi) = \mathcal{N}(\mathbf{w} | \mu_N, \Sigma_N),$$

where

$$\begin{aligned}\mu_N &= \Sigma_N(\Sigma_0^{-1}\mu_0 + \beta\Phi^\top\mathbf{y}), \\ \Sigma_N^{-1} &= \Sigma_0^{-1} + \beta\Phi^\top\Phi.\end{aligned}$$

For the sake of simplicity, we consider a particular form of Gaussian prior (**zero mean isotropic Gaussian prior**),

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} \mathbf{I}).$$

Then the corresponding posterior distribution over \mathbf{w} is given by

$$p(\mathbf{w} | \mathbf{y}) = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N),$$

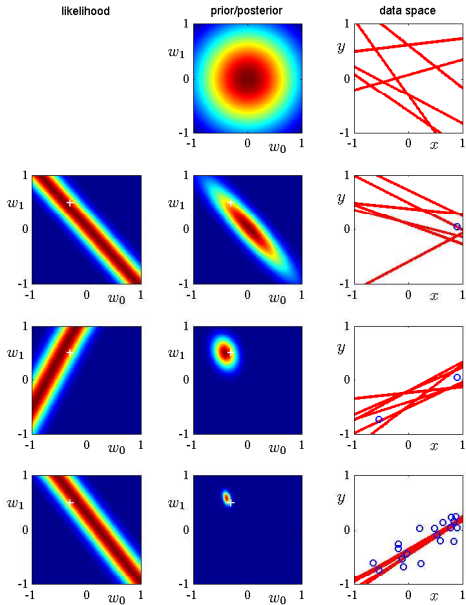
where

$$\begin{aligned} \boldsymbol{\mu}_N &= \beta \boldsymbol{\Sigma}_N \boldsymbol{\Phi}^\top \mathbf{y}, \\ \boldsymbol{\Sigma}_N^{-1} &= \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^\top \boldsymbol{\Phi}. \end{aligned}$$

An Example of Sequential Bayesian Learning

Consider a linear model $y = w_0 + w_1x$ where only two parameters are involved. Next slide illustrates the sequential nature of Bayesian learning, showing that the posterior distribution over parameters become sharper as more data points are observed.

- ▶ **Left-hand column:** likelihood, $p(y|x, \mathbf{w})$.
- ▶ **Middle column:** prior/posterior, $p(\mathbf{w}|D_t)$.
- ▶ **Right-hand column:** samples of function $y = w_0 + w_1x$ where \mathbf{w} are drawn from the posterior.



Predictive Distribution

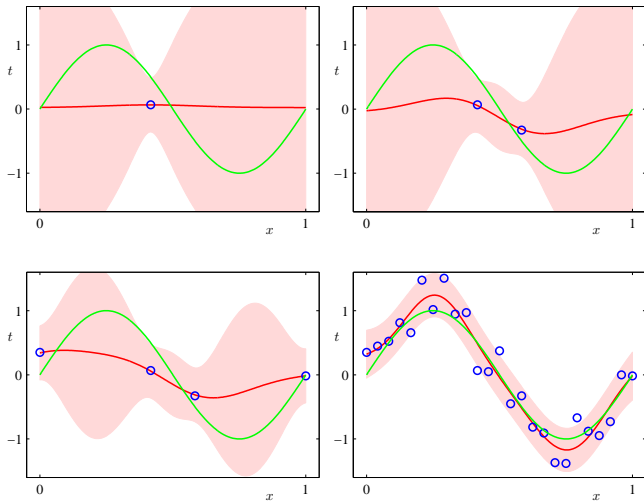
Make predictions of y_* for new values \mathbf{x}_* (i.e., ϕ_*). Let $\mathcal{D} = \{\Phi, \mathbf{y}\}$. Then the predictive distribution is given by

$$\begin{aligned} p(y_* | \phi_*, \mathcal{D}, \alpha, \beta) &= \int p(y_* | \phi_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathcal{D}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(y_* | \boldsymbol{\mu}_N^\top \phi_*, \sigma_N^2(\mathbf{x}_*)), \end{aligned}$$

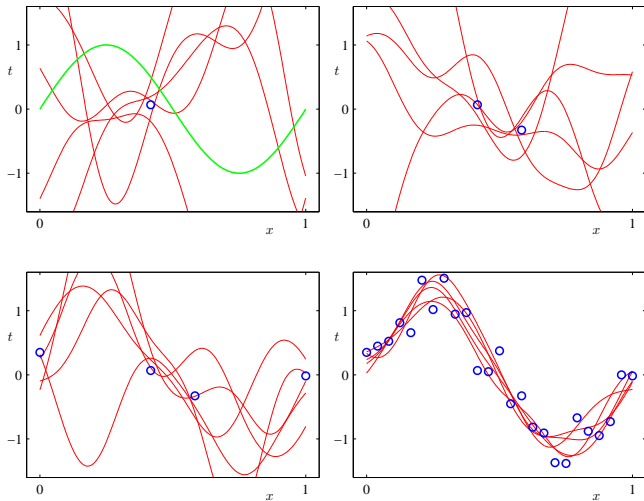
where

$$\begin{aligned} \boldsymbol{\mu}_N &= \beta \boldsymbol{\Sigma}_N \Phi^\top \mathbf{y}, \quad \boldsymbol{\Sigma}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^\top \Phi, \\ \sigma_N^2(\mathbf{x}_*) &= \underbrace{\frac{1}{\beta}}_{\text{noise on the data}} + \underbrace{\phi_*^\top \boldsymbol{\Sigma}_N \phi_*}_{\text{uncertainty associated with parameters } \mathbf{w}}. \end{aligned}$$

Examples of Predictive Distributions



Plots of $f(x; \mathbf{w})$ using samples from $p(\mathbf{w} | \mathbf{y})$



Bayesian Model Comparison

- ▶ Avoid the over-fitting associated with maximum likelihood by marginalizing over the model parameters instead of making point estimates of their values.
- ▶ Models can be compared directly on the training data, without the need for a validation set.
- ▶ Avoids the multiple training runs for each model associated with cross-validation.

Suppose that we wish to compare a set of L models, $\{\mathcal{M}_i\}$ for $i = 1, \dots, L$. (a model refers to a probability distribution over the observed data \mathcal{D})

Given a training set \mathcal{D} , we then wish to evaluate the posterior distribution

$$p(\mathcal{M}_i|\mathcal{D}) \propto \underbrace{p(\mathcal{M}_i)}_{\text{prior}} \underbrace{p(\mathcal{D}|\mathcal{M}_i)}_{\text{evidence}}.$$

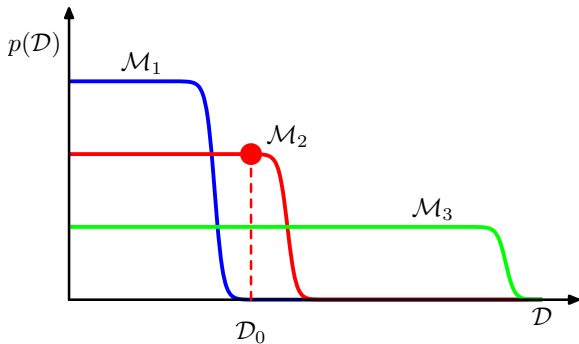
- ▶ $p(\mathcal{M}_i)$ is a prior probability distribution to express our uncertainty. The data is generated from one of these models but we are uncertain which one.
- ▶ $p(\mathcal{D}|\mathcal{M}_i)$ is the **model evidence** which expresses the preference shown by the data for different models. This is also known as **marginal likelihood**, since it can be viewed as a likelihood function over the space of models, in which the parameters have been marginalized out.

- ▶ **Bayes factor** is the ratio of model evidences for two models:

$$\frac{p(\mathcal{D}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_j)}.$$

- ▶ **Model selection**: Choose the single most probable model. For a model governed by a set of parameters \mathbf{w} , the model evidence is given by

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)d\mathbf{w}.$$



- ▶ \mathcal{M}_1 is the simplest and \mathcal{M}_3 is the most complex.
- ▶ For the particular observed data set \mathcal{D}_0 , the model \mathcal{M}_2 with intermediate complexity has the largest evidence.

Marginal Likelihood: Hyperparameter Estimation

We estimate hyperparameters α and β through maximizing the marginal likelihood.

The marginal likelihood is given by

$$p(\mathbf{y} | \Phi, \alpha, \beta) = \int p(\mathbf{y} | \Phi, \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w}.$$

Marginal likelihood maximization is illustrated in detail in Sec. 3.5.1 and 3.5.2.