

Linear Regression and Classification

Wenting Tu

SHUFE, SIME

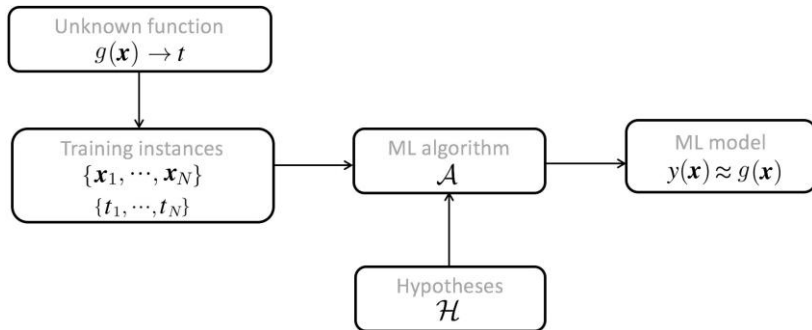
Machine Learning and Deep Learning

Outline

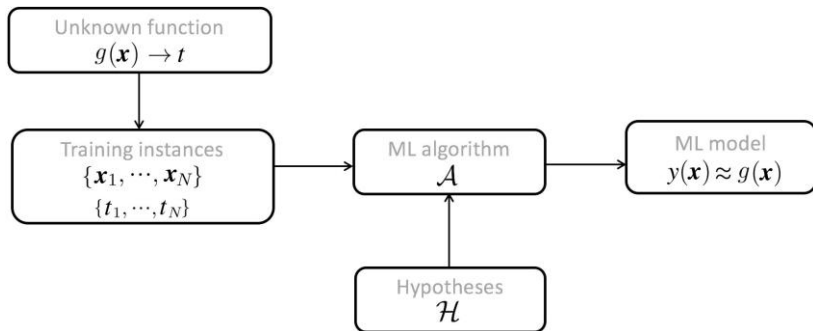
Linear Regression

Linear Classification

Definition



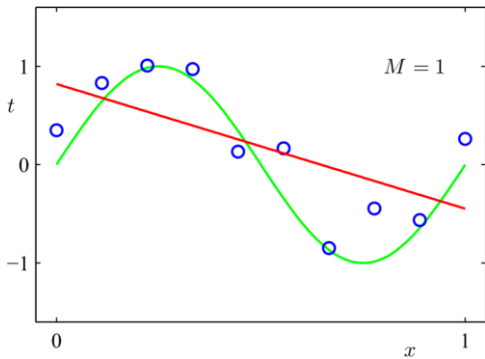
Definition



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

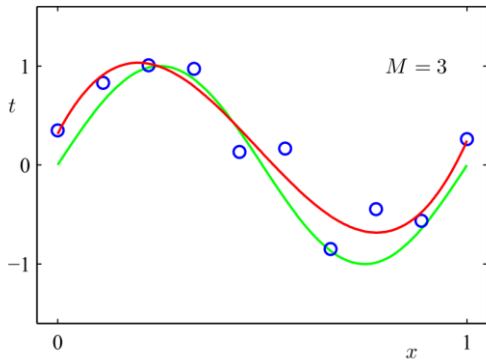
$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_D x_D$$

Illustration



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

Illustration



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

Linear Basis Function Models

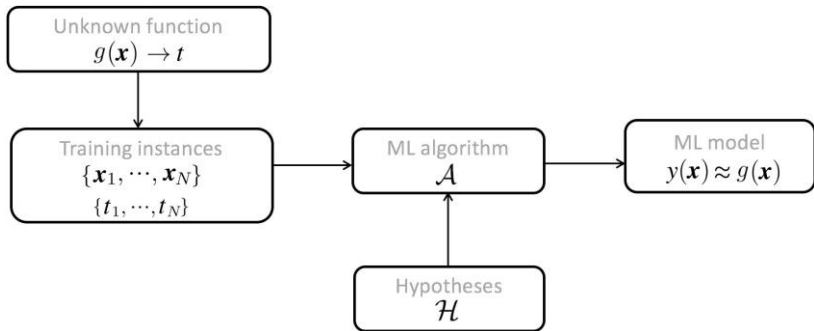
$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad \phi_0(\mathbf{x}) = 1$$

$$\phi_j(x) = x^j$$

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

$$\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right) \quad \sigma_a = \frac{1}{1 + \exp(-a)}$$



Least-squares

- Loss function

$$L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2 \quad (\text{squared residuals})$$

- Empirical risk

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

Least-squares

- Solution

$$\mathbf{w}^{\star} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

Least-squares and Maximum Likelihood

- Review for maximum-likelihood estimation (MLE)

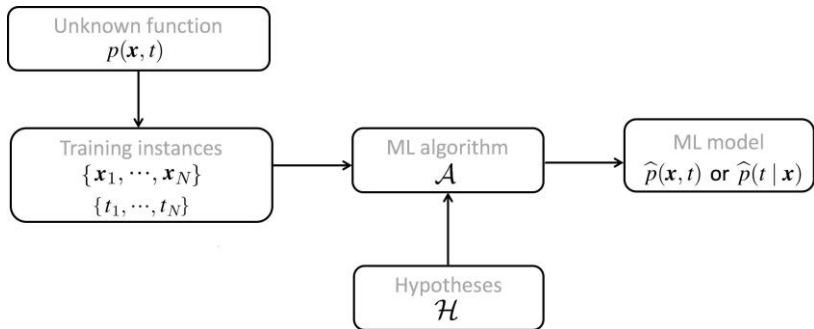
$$q(\mathbf{x}; \boldsymbol{\theta}) \longrightarrow p(\mathbf{x})$$

$$\mathcal{D} = \{\mathbf{x}\}_{i=1}^n$$

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathcal{D} \mid \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^n q(\mathbf{x}; \boldsymbol{\theta})$$

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \log L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \left[\sum_{i=1}^n \log q(\mathbf{x}_i; \boldsymbol{\theta}) \right]$$

$$\hat{p}(\mathbf{x}) = q(\mathbf{x}; \hat{\boldsymbol{\theta}}_{\text{ML}})$$



Least-squares and Maximum Likelihood

- Relation between LS and MLE

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

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

$$\mathbb{E}[t \mid \mathbf{x}] = \int t p(t \mid \mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$$

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$\begin{aligned} \ln p(\mathbf{t} \mid \mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \end{aligned}$$

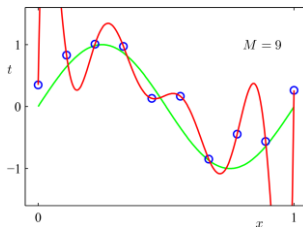
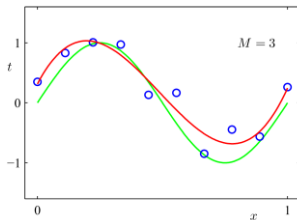
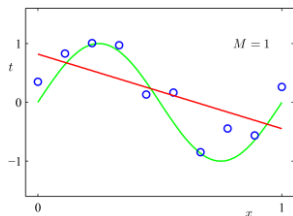
$$\mathbf{w}_{MLE}^* = \arg \max_{\mathbf{w}} \ln p(\mathbf{t} \mid \mathbf{w}, \beta)$$

$$\mathbf{w}_{MLE}^* = \mathbf{w}_{LS}^*$$

Underfitting and Overfitting

- Illustration

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



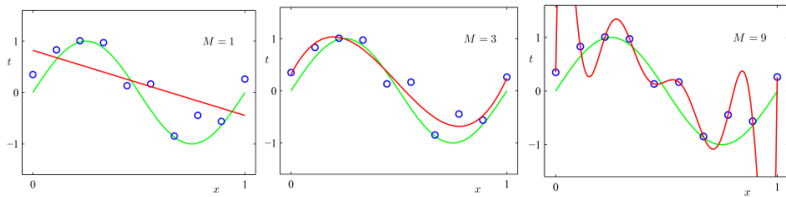
Bias-Variance Decomposition

expected loss = (bias)² + variance + noise

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

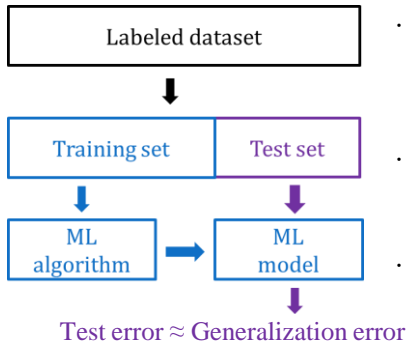
$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} \left[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 \right] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$



Model Evaluation

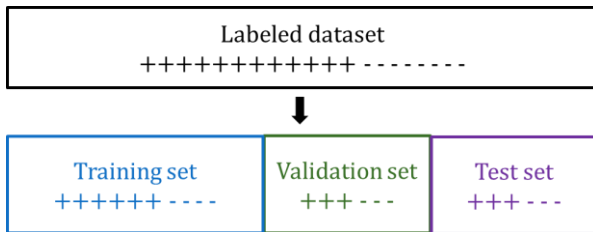
- Generalization ability



- When learning a model, you should pretend that you don't have the test data yet. If the test-set labels influence the learned model in any way, accuracy estimates will be biased.
- Your test set should be large enough to detect meaningful changes in the accuracy of your algorithm, but not necessarily much larger.
- When randomly selecting training or validation sets, we may want to ensure that class proportions are maintained in each selected set.

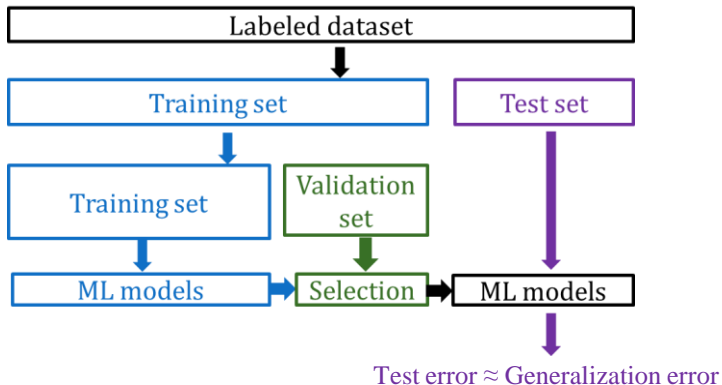
Model Evaluation

- Validation set



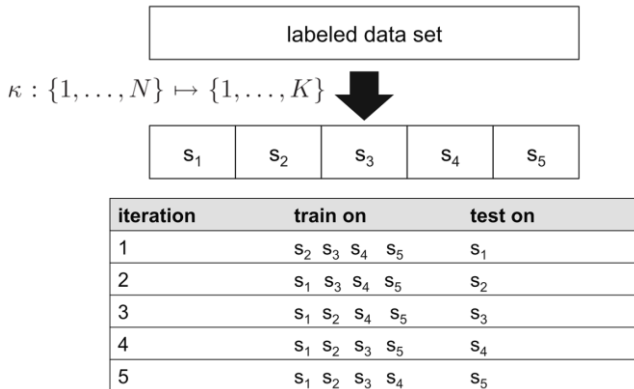
Model Evaluation

- Tuning hyperparameters



Model Evaluation

- Cross validation



The K results can then be averaged to produce a single estimation.

CV makes efficient use of the available data for testing

Model Evaluation

- Performance Evaluation for Regression

$$MAE = \frac{1}{n} \sum_{i=1}^n |t_i - f(\mathbf{x}_i)|$$

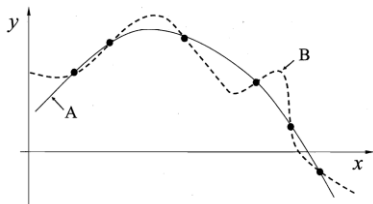
$$MSE = \frac{1}{n} \sum_{i=1}^n (t_i - f(\mathbf{x}_i))^2$$

$$RMSE = \sqrt{MSE}$$

Model Selection

- Occam's razor

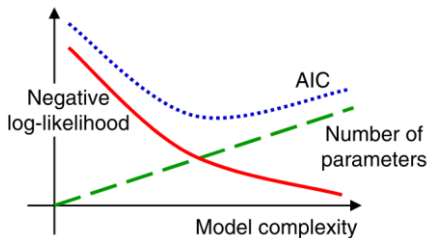
Suppose there exist two explanations for an occurrence. In this case the one that requires the smallest number of assumptions is usually correct.



Model Selection

- Akaike information criterion (AIC)

$$\ln p(\mathcal{D}|\mathbf{w}_{\text{ML}}) - M$$



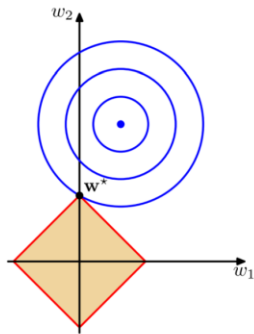
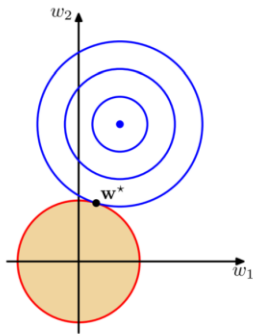
Regularization

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

LASSO

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

$$= \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \Phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



Ridge Regression

$$\begin{aligned} & E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \\ &= \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ \mathbf{w}_{\text{ridge}}^* &= (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned}$$

Ridge Regression and Maximum a Posteriori Estimation

- Review for maximum a posteriori (MAP) estimation

Likelihood $p(\mathcal{D}|\boldsymbol{\theta})$

Prior $p(\boldsymbol{\theta})$

Posterior $p(\boldsymbol{\theta}|\mathcal{D})$

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}|\mathcal{D})$$

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left(\sum_{i=1}^n \log q(\mathbf{x}_i|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \right)$$

Ridge Regression and Maximum a Posteriori Estimation

- Relation between ridge regression and MAP

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0)$$

$$p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t}), \quad \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

Ridge Regression and Maximum a Posteriori Estimation

- Relation between ridge regression and MAP

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

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

$$p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}, \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

$$\begin{aligned} \mathbf{w}_{MAP}^* &= \arg \max_{\mathbf{w}} \ln p(\mathbf{w} \mid \mathbf{t}) \\ &= \arg \max_{\mathbf{w}} \left(-\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const} \right) \end{aligned}$$

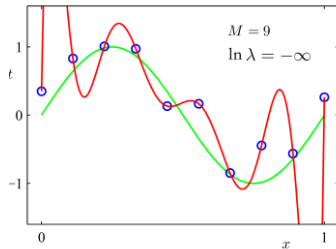
$$\mathbf{w}_{MAP}^* = \mathbf{w}_{ridge}^*$$

Ridge Regression

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \sum_{j=0}^M w_j x^j \right\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

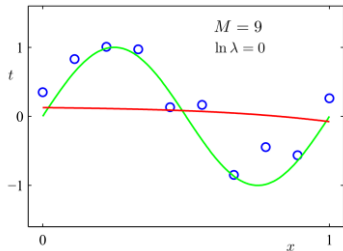
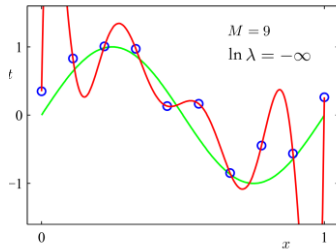
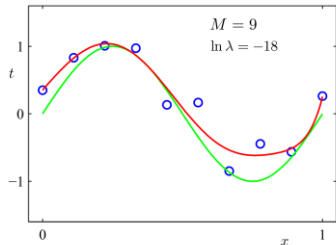
$$M = 9$$

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01



Ridge Regression

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \sum_{j=0}^M w_j x^j \right\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$



Bayesian Models

- Bayesian estimation

$$\int \boldsymbol{\theta} p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \quad (\text{posterior expectation})$$

- Bayesian inference

$$\begin{aligned}\hat{p}_{\text{Bayes}}(\mathbf{x}) &= \int q(\mathbf{x}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \\ &= \int q(\mathbf{x}|\boldsymbol{\theta}) \frac{p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})} = \int q(\mathbf{x}|\boldsymbol{\theta}) \frac{\prod_{i=1}^n q(\mathbf{x}_i|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int \prod_{i=1}^n q(\mathbf{x}_i|\boldsymbol{\theta}') p(\boldsymbol{\theta}') d\boldsymbol{\theta}'} d\boldsymbol{\theta}\end{aligned}$$

Bayesian Linear Regression

$$p(t \mid \mathbf{t}, \alpha, \beta) = \int p(t \mid \mathbf{w}, \beta) p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) d\mathbf{w} \quad (\text{Predictive distribution})$$

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

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