IASTED-NCI2004

Feb. 23-25, 2004

Estimating the Error at Given Test Input Points for Linear Regression



Fraunhofer Institut Rechnerarchitektur und Softwaretechnik



Masashi Sugiyama

Fraunhofer FIRST-IDA, Berlin, Germany Tokyo Institute of Technology, Tokyo, Japan

Regression Problem



 $f(\boldsymbol{x})$:Underlying function $\hat{f}(\boldsymbol{x})$:Learned function $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$:Training examples $y_i = f(\boldsymbol{x}_i) + \epsilon_i$ (noise) $\epsilon_i \stackrel{i.i.d.}{\sim}$ mean 0, variance σ^2

From $\{(x_i, y_i)\}_{i=1}^n$, obtain a good approximation $\hat{f}(x)$ to f(x)

Typical Method of Learning

Linear regression model

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{p} lpha_{i} \varphi_{i}(\boldsymbol{x})$$

 α_{i} :Parameters
 $\{\varphi_{i}(\boldsymbol{x})\}_{i=1}^{p}$:Fixed basis functions

Ridge estimation

$$\min_{\{\alpha_i\}} \left[\sum_{i=1}^n \left(\hat{f}(\boldsymbol{x}_i) - y_i \right)^2 + \lambda \sum_{i=1}^p \alpha_i^2 \right]$$

 λ :Ridge parameter (model parameter)

3

Model Selection

—— Underlying function f(x)—— Learned function $\hat{f}(x)$



Choice of the model is crucial for obtaining good learned function $\hat{f}(\mathbf{x})$!

Generalization Error

For model selection, we need a criterion that measures 'closeness' between $\hat{f}(x)$ and f(x):

Generalization error, e.g., $J_{gen}(\lambda) = \int \left(\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2 p(\boldsymbol{x}) d\boldsymbol{x}$

Determine the model
$$\lambda$$

so that an estimator \hat{J}_{gen} of
the unknown generalization
error J_{gen} is minimized.
 $\hat{\lambda} = \underset{\lambda}{\operatorname{argmin}} \hat{J}_{gen}(\lambda)$

 $p(\boldsymbol{x})$:Probability density of test input points



Transductive Inference

Test input points are specified in advance.
We do not have to estimate the entire function *f(x)*, but just estimate the values of the function at the test input points {*t_i*}^{nt}_{i=1}.





- Test error at given test input points is different from the generalization error.
- Model should be chosen so that the test error only at $\{t_i\}_{i=1}^{n_t}$ is minimized.



Small generalization error Large test error



Large generalization error Small test error

Goal of Our Research

We want to estimate the test error at the given test input points!

$$J_{test} = E \sum_{i=1}^{n_t} \left(\hat{f}(\boldsymbol{t}_i) - f(\boldsymbol{t}_i) \right)^2$$

E :Expectation over noise



Setting

Linear regression model

m

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{p} \alpha_i \varphi_i(\boldsymbol{x}) \qquad \alpha_i \text{ :Parameters} \\ \{\varphi_i(\boldsymbol{x})\}_{i=1}^{p} \text{:Fixed basis functions} \\ \text{Linear estimation} \end{cases}$$

$$\hat{lpha} = Xy$$

 $\hat{\boldsymbol{\alpha}} = (\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_p)^\top$ \boldsymbol{X} : A matrix $\boldsymbol{y} = (y_1, y_2, \dots, y_n)^\top$

Realizability

$$f(\boldsymbol{x}) = \sum_{i=1}^{p} \alpha_i^* \varphi_i(\boldsymbol{x})$$

 α_i^* : Unknown true parameters

Bias / Variance Decomposition¹⁰

$$J_{test} = \mathbf{E} \sum_{i=1}^{n_t} \left(\hat{f}(t_i) - f(t_i) \right)^2 \begin{bmatrix} [\mathbf{A}_t]_{i,j} = \varphi_j(t_i) \\ \mathbf{A}_t \hat{\alpha} = (\hat{f}(t_1), \hat{f}(t_2), \dots, \hat{f}(t_{n_t}))^\top \\ = \mathbf{E} \| \mathbf{A}_t \hat{\alpha} - \mathbf{A}_t \alpha^* \|^2 \qquad \alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_p^*)^\top \\ = \| \mathbf{E} \mathbf{A}_t \hat{\alpha} - \mathbf{A}_t \alpha^* \|^2 + \mathbf{E} \| \mathbf{A}_t \hat{\alpha} - \mathbf{E} \mathbf{A}_t \hat{\alpha} \|^2 \\ \mathbf{Bias} \qquad \mathbf{Variance} \\ \hline \mathbf{E} \mathbf{A}_t \hat{\alpha} \qquad \mathbf{Bias} \qquad \mathbf{A}_t \alpha^* \\ \bullet \qquad \mathbf{Variance} \\ \hline \mathbf{Variance} \\ \hline \mathbf{A}_t \hat{\alpha} \qquad \mathbf{Variance} \\ \hline \hline \mathbf$$



Tricks for Estimating Bias

Sugiyama & Ogawa (Neural Comp., 2001) Sugiyama & Müller (JMLR, 2002)

11

$$Bias = ||\mathbf{E} \mathbf{A}_t \hat{\boldsymbol{\alpha}} - \mathbf{A}_t \boldsymbol{\alpha}^*||^2$$

True parameter α^* is unknown.

We utilize an unbiased estimator of the true parameter for estimating the bias.

$$\hat{\boldsymbol{lpha}}_{u} = \boldsymbol{A}^{\dagger} \boldsymbol{y}$$

 $\mathrm{E} \hat{\boldsymbol{lpha}}_{u} = \boldsymbol{lpha}^{*}$

 $\begin{array}{ll} \boldsymbol{A} : \text{Design matrix} & \dagger : \text{Generalized inverse} \\ \boldsymbol{A}_{i,j} = \varphi_j(\boldsymbol{x}_i) \\ \boldsymbol{A} \hat{\boldsymbol{\alpha}} = (\hat{f}(\boldsymbol{x}_1), \hat{f}(\boldsymbol{x}_2), \dots, \hat{f}(\boldsymbol{x}_n))^\top \end{array}$

Unbiased Estimator of Bias¹²



$$Bias = \|\mathbf{E}\mathbf{A}_{t}\hat{\boldsymbol{\alpha}} - \mathbf{A}_{t}\boldsymbol{\alpha}^{*}\|^{2} \stackrel{\boldsymbol{\epsilon} = (\epsilon_{1}, \epsilon_{2}, \dots, \epsilon_{n})^{\top}}{\boldsymbol{z} = (f(\boldsymbol{x}_{1}), f(\boldsymbol{x}_{2}), \dots, f(\boldsymbol{x}_{n}))^{\top}}$$
$$= \|\mathbf{A}_{t}\hat{\boldsymbol{\alpha}} - \mathbf{A}_{t}\hat{\boldsymbol{\alpha}}_{u}\|^{2} - 2\langle \mathbf{A}_{t}\boldsymbol{X}\boldsymbol{z}, \mathbf{A}_{t}\boldsymbol{A}^{\dagger}\boldsymbol{\epsilon}\rangle - \|\mathbf{A}_{t}(\boldsymbol{X} - \boldsymbol{A}^{\dagger})\boldsymbol{\epsilon}\|^{2}$$
$$\stackrel{\mathbf{E}}{\mathbf{E}}$$
$$\widehat{Bias} = \|\mathbf{A}_{t}\hat{\boldsymbol{\alpha}} - \mathbf{A}_{t}\hat{\boldsymbol{\alpha}}_{u}\|^{2} - 0 - \sigma^{2}\mathrm{tr}\left(\mathbf{A}_{t}(\boldsymbol{X} - \boldsymbol{A}^{\dagger})[\mathbf{A}_{t}(\boldsymbol{X} - \boldsymbol{A}^{\dagger})]^{\top}\right)$$
$$\stackrel{\mathbf{E}\widehat{Bias} = Bias}{\mathbf{E}}$$

Unbiased Estimator of Variance¹³

$$Var = E || \mathbf{A}_t \hat{\boldsymbol{\alpha}} - E \mathbf{A}_t \hat{\boldsymbol{\alpha}} ||^2$$
$$= \sigma^2 \operatorname{tr} \left(\mathbf{A}_t \mathbf{X} (\mathbf{A}_t \mathbf{X})^\top \right)$$
$$\sigma^2 : \text{Noise variance}$$

An unbiased estimator of noise variance: $\hat{\sigma}^2 = \frac{\|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{A}^{\dagger}\boldsymbol{y}\|^2}{n-p}$ $\widehat{Var} = \hat{\sigma}^2 \operatorname{tr} \left(\boldsymbol{A}_t \boldsymbol{X} (\boldsymbol{A}_t \boldsymbol{X})^{\top}\right)$

$$\boxed{ E \ \widehat{Var} = Var }$$

Unbiased Estimator of Test Error¹⁴

Adding bias and variance estimators, we have an unbiased estimator of test error.

 $\hat{J}_{test} = \widehat{Bias} + \widehat{Var}$

For simplicity, we ignore constant terms

$$\hat{J}_t = \|\boldsymbol{A}_t \boldsymbol{X} \boldsymbol{y}\|^2 - 2\langle \boldsymbol{A}_t \boldsymbol{X} \boldsymbol{y}, \boldsymbol{A}_t \boldsymbol{A}^{\dagger} \boldsymbol{y} \rangle + 2\hat{\sigma}^2 \operatorname{tr} \left(\boldsymbol{A}_t \boldsymbol{X} (\boldsymbol{A}_t \boldsymbol{A}^{\dagger})^{\top} \right)$$





Unrealizable Cases

So far, we assumed that the model includes the underlying function.

$$f(m{x}) = \sum_{i=1}^p lpha_i^* arphi_i(m{x}) \quad lpha_i^*$$
:Unknown true parameters

We can prove that even when the above assumption is not rigorously fulfilled, \hat{J}_t is still almost unbiased.

 $\mathrm{E}\hat{J}_t \approx J_t$

Simulation: Toy Data Sets

- Basis functions: 10 Gaussian functions centered at equally located points in $[-\pi, \pi]$.
- Target function: sinc-like function (realizable).
- Training examples $\{(x_i, y_i)\}_{i=1}^n$: $x_i \stackrel{i.i.d.}{\sim} U(-\pi, \pi)$ $y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$
- Test input points $\{t_i\}_{i=1}^{50}$: $t_i \overset{i.i.d.}{\sim} U(-\pi,\pi)$
- Ridge estimation is used for learning.







Simulation: DELVE Data Sets ¹⁹

- Training set: 100 randomly selected samples.
- **Test set:** 50 randomly selected samples.
- Basis functions: Gaussian function centered at first 50 training input points.
- Ridge estimation is used for learning.
- Ridge parameter is selected by the proposed method, leave-one-out cross-validation, or an empirical Bayesian method.

Normalized Test Errors

Mean (Standard deviation)

Data set	Proposed method	LOO cross- validation	Empirical Bayes
Boston	1.17 (0.54)	1.26 (0.58)	1.39 (0.59)
Bank-8fm	1.07 (0.29)	1.11 (0.32)	1.09 (0.31)
Bank-8nm	1.09 (0.51)	1.12 (0.56)	1.18 (0.60)
Kin-8fm	1.06 (0.32)	1.17 (0.36)	1.68 (0.48)
Kin-8nm	1.11 (0.27)	1.09 (0.24)	1.15 (0.24)

Red: Best and others with no significant difference by 99% t-test

Proposed method can be successfully applied to transductive model selection!





- Model selection is usually carried out so that estimated generalization error is minimized.
- When test input points are specified in advance (transductive inference), it is natural to choose a model so that the test error only at the test input points is minimized.
- We derived an unbiased estimator of the test error at given test input points.
- Simulation showed the proposed method works well in practical situations.