Pattern Information Processing:⁷⁰ Model Selection by Cross-Validation

> Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

Model Parameters

- In the process of parameter learning, we fixed model parameters.
- For example, quadratically constrained least-squares with Gaussian kernel models
 - Gaussian width: c (> 0)
 - Regularization parameter: $\lambda \ (\geq 0)$

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i)$$
$$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2c^2}\right)$$
$$J_{QCLS}(\boldsymbol{\alpha}) = J_{LS}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|^2$$

Different Model Parameters

72

Model parameters strongly affect learned functions.



Determining Model Parameters⁷³

We want to determine the model parameters so that the generalization error (expected test error) is minimized.

$$G = \int_{\mathcal{D}} \left(\hat{f}(\boldsymbol{t}) - f(\boldsymbol{t}) \right)^2 q(\boldsymbol{t}) d\boldsymbol{t}$$
$$\boldsymbol{t} \sim q(\boldsymbol{x})$$

However, f(x) is unknown so the generalization error is not accessible.
q(x) may also be unknown.

Generalization Error Estimation⁷⁴

$$G = \int_{\mathcal{D}} \left(\hat{f}(\boldsymbol{t}) - f(\boldsymbol{t}) \right)^2 q(\boldsymbol{t}) d\boldsymbol{t}$$

Instead, we use a generalization error estimate.



Model Selection

Prepare a set of model candidates.

$$\{\mathcal{M}_i \mid \mathcal{M}_i = (c_i, \lambda_i)\}$$

Estimate generalization error for each model. $\widehat{G}(\mathcal{M}_i)$

Choose the one with the minimum estimated generalization error.

$$\hat{\mathcal{M}} = \operatorname*{argmin}_{\mathcal{M} \in \{\mathcal{M}_i\}_i} \widehat{G}(\mathcal{M})$$

Assumptions

Training input points: x_i ^{i.i.d.} q(x)
Training output values: y_i = f(x_i) + \epsilon_i
Noise \epsilon_i : i.i.d., mean 0, variance \sigma^2

$$\mathbb{E}_{\boldsymbol{\epsilon}}[\epsilon_i] = 0 \qquad \mathbb{E}_{\boldsymbol{\epsilon}}[\epsilon_i \epsilon_j] = \begin{cases} \sigma^2 & (i=j) \\ 0 & (i\neq j) \end{cases}$$

Extra-Sample Method

77

Suppose we have an extra example (x', y')in addition to $\{(x_i, y_i)\}_{i=1}^n$.

$$y' = f(x') + \epsilon' \quad x' \sim q(x) \quad \mathbb{E}_{\epsilon}[\epsilon'] = 0$$
$$\mathbb{E}_{\epsilon}[\epsilon'^2] = \sigma^2$$
$$\mathbb{E}_{\epsilon}[\epsilon'\epsilon_i] = 0, \quad \forall i$$

Test the learned function using the extra example.

$$\widehat{G}_{extra} = \left(\widehat{f}(\boldsymbol{x}') - \boldsymbol{y}'\right)^2$$
$$\widehat{f} \longleftarrow \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$$

Extra-Sample Method (cont.) ⁷⁸

 \widehat{G}_{extra} is unbiased w.r.t. x' and ϵ' (up to σ^2) $\mathbb{E}_{x'}\mathbb{E}_{\epsilon'}[\widehat{G}_{extra}] = G + \sigma^2$

Proof:

$$\mathbb{E}_{\boldsymbol{x}'} \mathbb{E}_{\boldsymbol{\epsilon}'} \left(\hat{f}(\boldsymbol{x}') - f(\boldsymbol{x}') - \boldsymbol{\epsilon}' \right)^2$$

= $\mathbb{E}_{\boldsymbol{x}'} \mathbb{E}_{\boldsymbol{\epsilon}'} \left[(\hat{f}(\boldsymbol{x}') - f(\boldsymbol{x}'))^2 - 2\boldsymbol{\epsilon}' (\hat{f}(\boldsymbol{x}') - f(\boldsymbol{x}')) + \boldsymbol{\epsilon}'^2 \right]$
= $G + \sigma^2$

Gextra may be used for model selection.

 However, in practice, such an extra example is not available (or if we have, it should be included in the original training set!).

Holdout Method

Idea: artificially create an extra sample

- 1. Divide training set $\{(x_i, y_i)\}_{i=1}^n$ into $\{(x_i, y_i)\}_{i \neq j}$ and (x_j, y_j) .
- 2. Train a learning machine using $\{(x_i, y_i)\}_{i \neq j}$ $\hat{f}_j(x) \leftarrow \{(x_i, y_i)\}_{i \neq j}$

3. Test it using the holdout sample (x_j, y_j)

$$\widehat{G}_j = \left(\widehat{f}_j(\boldsymbol{x}_j) - y_j\right)^2$$

Almost Unbiasedness of Holdout⁸⁰ Holdout method is almost unbiased w.r.t. x_{i}, ϵ_{i} :

$$\mathbf{\mathcal{E}}_{\boldsymbol{x}_{j}} \mathbb{E}_{\epsilon_{j}}[\widehat{G}_{j}] = G_{j} + \sigma^{2}$$
$$\approx G + \sigma^{2}$$
$$\mathbf{\mathcal{G}}_{j} = \int_{\mathcal{D}} \left(\widehat{f}_{j}(\boldsymbol{x}) - f(\boldsymbol{x})\right)^{2} q(\boldsymbol{x}) d\boldsymbol{x}$$

 $\widehat{f}_j(\boldsymbol{x}) \approx \widehat{f}(\boldsymbol{x})$ if n is large

However, \widehat{G}_j is heavily affected by the choice of the holdout sample (x_j, y_j) .

Leave-One-Out Cross-Validation⁸¹

Repeat the holdout procedure for all combinations and output the average.

$$\widehat{G}_{LOOCV} = \frac{1}{n} \sum_{j=1}^{n} \widehat{G}_j$$

$$\widehat{G}_j = \left(\widehat{f}_j(\boldsymbol{x}_j) - y_j\right)^2$$

LOOCV is almost unbiased w.r.t. $\{x_i, \epsilon_i\}_{i=1}^n$

 $\mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n} \mathbb{E}_{\{\epsilon_i\}_{i=1}^n} [\widehat{G}_{LOOCV}] \approx \mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n} \mathbb{E}_{\{\epsilon_i\}_{i=1}^n} [G] + \sigma^2$

k-Fold Cross-Validation

Randomly split training set into k disjoint subsets $\{\mathcal{T}_j\}_{j=1}^k$.

$$\widehat{G}_{kCV} = \frac{1}{k} \sum_{j=1}^{k} \widehat{G}_{\mathcal{T}_j}$$
$$\widehat{G}_{\mathcal{T}_j} = \frac{1}{|\mathcal{T}_j|} \sum_{i \in \mathcal{T}_j} \left(\widehat{f}_{\mathcal{T}_j}(\boldsymbol{x}_i) - y_i \right)^2$$
$$\widehat{f}_{\mathcal{T}_j}(\boldsymbol{x}) \longleftarrow \{ (\boldsymbol{x}_i, y_i) \mid i \notin \mathcal{T}_j \}$$

k-fold is easier to compute and more stable.

Advantages of CV

83

 Wide applicability: Almost unbiasedness of LOOCV holds for (virtually) any learning methods
 Practical usefulness: CV is shown to work very well in many practical applications

Disadvantages of CV

 Computationally expensive It requires repeating training of models with different subsets of training samples
 Number of folds

It is often recommended to use k = 5, 10. However, how to optimally choose k is still open.

Closed Form of LOOCV

Linear model
$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{b} \alpha_i \varphi_i(\boldsymbol{x})$$

Quadratically constrained least-squares

 $J_{QCLS}(\boldsymbol{\alpha}) = J_{LS}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|^2$

$$\widehat{G}_{LOOCV} = \frac{1}{n} \|\widetilde{\boldsymbol{H}}^{-1} \boldsymbol{H} \boldsymbol{y}\|^2$$

 $H = I - XL_{QCLS}$ $L_{QCLS} = (X^{\top}X + \lambda I)^{-1}X^{\top}$

H :same diagonal as H but zero for off-diagonal

85

Basic Idea

$$\widehat{G}_{LOOCV} = \frac{1}{n} \|\widetilde{\boldsymbol{H}}^{-1} \boldsymbol{H} \boldsymbol{y}\|^2$$

- Let $\hat{\alpha}^{(j)}$ be the learned parameter without the j-th sample (x_j, y_j) .
- Express $\hat{\alpha}^{(j)}$ in terms of $\hat{\alpha}$ (the learned parameter with all samples).
- Key fact: $(U - uu^{\top})^{-1} = U^{-1} + \frac{U^{-1}uu^{\top}U^{-1}}{1 - u^{\top}U^{-1}u}$
- Proof is homework (although it is rather hard...)!

Homework (cont.)

 (Try to) prove the closed-form expression of leave-one-out crossvalidation score for quadratically constraint least-squares.

$$\widehat{G}_{LOOCV} = \frac{1}{n} \|\widetilde{\boldsymbol{H}}^{-1} \boldsymbol{H} \boldsymbol{y}\|^2$$

Homework (cont.)

88

2. For your own toy 1-dimensional data, perform simulations using

•Gaussian kernel models

 Quadratically-constrained least-squares learning and optimize

Width of Gaussian kernel

Regularization parameter

based on cross-validation. Analyze the results when changing

- Target function
- Number of samples
- Noise level

Suggestions

Please look for software which can solve

- Linearly constrained quadratic programming $\min_{\beta} \left[\frac{1}{2} \langle Q\beta, \beta \rangle + \langle \beta, q \rangle \right]$ subject to $V\beta \leq v$ and $G\beta = g$ • Linearly constrained linear programming
 - Linearly constrained linear programming $\min_{\beta} \langle \beta, q \rangle$

subject to $V\beta \leq v$ and $G\beta = g$

For example, MOSEK, LOQO, or SeDuMi.

The software is not necessarily sophisticated; just elementary one is enough.