

Improving importance estimation in pool-based batch active learning for approximate linear regression

Nozomi Kurihara, Masashi Sugiyama*

Department of Computer Science, Tokyo Institute of Technology, 2-12-1-W8-74 O-okayama, Meguro-ku, Tokyo 152-8552, Japan

ARTICLE INFO

Article history:

Received 10 February 2012
Received in revised form 4 June 2012
Accepted 3 September 2012

Keywords:

Pool-based batch active learning
Approximate linear regression
Covariate shift
Importance-weighted least-squares
P-ALICE
Inclusion probability

ABSTRACT

Pool-based batch active learning is aimed at choosing training inputs from a ‘pool’ of test inputs so that the generalization error is minimized. P-ALICE (Pool-based Active Learning using Importance-weighted least-squares learning based on Conditional Expectation of the generalization error) is a state-of-the-art method that can cope with model misspecification by weighting training samples according to the importance (i.e., the ratio of test and training input densities). However, importance estimation in the original P-ALICE is based on the assumption that the number of training samples to gather is small, which is not always true in practice. In this paper, we propose an alternative scheme for importance estimation based on the inclusion probability, and show its validity through numerical experiments.

© 2012 Elsevier Ltd. All rights reserved.

1. Introduction

The objective of supervised learning is to find an input–output relationship behind training samples (Bishop, 2006; Hastie, Tibshirani, & Friedman, 2001). Once the input–output relationship is successfully learned, outputs for unseen inputs can be predicted, i.e., the learning machine can *generalize*.

When users are allowed to choose the location of training inputs, it is desirable to design the input locations so that the generalization error is minimized. Such a problem is called *active learning* (Settles, 2009) or *experiment design* (Fedorov, 1972; Pukelsheim, 1993), and has been shown to be useful in various application areas such as text classification (Lewis & Gale, 1994; McCallum & Nigam, 1998), age estimation from images (Ueki, Sugiyama, & Ihara, 2010), medical data analysis (Wiens & Guttag, 2010), chemical data analysis (Warmuth et al., 2003), biological data analysis (Liu, 2004), and robot control (Akiyama, Hachiya, & Sugiyama, 2010).

If users are allowed to locate training inputs at any position in the domain, the active learning setup is said to be *population-based* (Kanamori & Shimodaira, 2003; Sugiyama, 2006; Wiens, 2000). On the other hand, if users need to choose training input locations from a pool of finite candidate points, it is said to be *pool-based* (Kanamori, 2007; McCallum & Nigam, 1998; Sugiyama & Nakajima, 2009). Depending on the way training input locations are chosen, active learning is also categorized into *sequential* or *batch* approaches: Training inputs are selected one by one iteratively in the sequential approach (Box & Hunter, 1965;

Sugiyama & Ogawa, 2000), while all training inputs are selected at once in the batch approach (Kiefer, 1959; Sugiyama & Ogawa, 2001). In this paper, we focus on pool-based batch active learning.

Active learning generally induces a *covariate shift* – a situation where training and test input distributions are different (Quiñonero-Candela et al., 2009; Shimodaira, 2000; Sugiyama & Kawanabe, 2012). When a model is correctly specified, covariate shifts do not matter in designing active learning methods. However, for a misspecified model, a covariate shift causes a strong estimation bias and thus classical active learning techniques that require a correct model become unreliable (Fedorov, 1972; Kiefer, 1959).

To cope with the bias induced by the covariate shift, active learning techniques that explicitly take model misspecification into account have been developed (Beygelzimer, Dasgupta, & Langford, 2009; Kanamori, 2007; Kanamori & Shimodaira, 2003; Sugiyama, 2006; Sugiyama & Nakajima, 2009; Wiens, 2000). The key idea of covariate shift adaptation is *importance weighting* – a loss function used for training is weighted according to the importance (i.e., the ratio of test and training input densities). Among the importance-weighted active learning methods, the pool-based batch active learning method for approximate linear regression called P-ALICE (Pool-based Active Learning using Importance-weighted least-squares learning based on Conditional Expectation of the generalization error) was demonstrated to be useful (Sugiyama & Nakajima, 2009).

However, in the original P-ALICE, the number of training samples to gather is assumed to be sufficiently smaller than the size of the sample pool. However, when this assumption is not satisfied, the importance weight used in P-ALICE is not reliable.

* Corresponding author.

E-mail addresses: sugi@cs.titech.ac.jp, sugi@sg.cs.titech.ac.jp (M. Sugiyama).

In this paper, we propose a new method to set the importance weight that does not rely on this assumption. Our new weighting scheme is based on the *inclusion probability* (Horvitz & Thompson, 1952), which allows us to precisely capture the relation between the training and test input distributions. Through experiments, we show that the active learning performance of P-ALICE can be improved by the proposed weighting method when the training sample size is relatively large.

The rest of this paper is structured as follows. In Section 2, we formulate the problem of pool-based active learning and give an overview of P-ALICE. In Section 3, we point out a limitation of importance estimation in P-ALICE, and propose an alternative method. In Section 4, experimental results on toy and benchmark datasets are reported. Finally, concluding remarks are given in Section 5.

2. Problem formulation

In this section, we formulate the problem of pool-based active learning and briefly review the P-ALICE method.

2.1. Pool-based active learning for linear regression

Let us consider a regression problem of learning a real-valued function $f(\mathbf{x})$ defined on $\mathcal{D} \subset \mathbb{R}^d$. For training input-output samples

$$\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}}) \mid y_i^{\text{tr}} = f(\mathbf{x}_i^{\text{tr}}) + \epsilon_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}},$$

where $\{\epsilon_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}}$ are i.i.d. noise with mean zero and unknown variance σ^2 , let us use the following linear regression model:

$$\hat{f}(\mathbf{x}) = \sum_{\ell=1}^t \theta_{\ell} \varphi_{\ell}(\mathbf{x}), \quad (1)$$

where $\{\varphi_{\ell}(\mathbf{x})\}_{\ell=1}^t$ are fixed linearly independent basis functions and $\{\theta_{\ell}\}_{\ell=1}^t$ are parameters to be learned. Let us denote the vector of parameters by $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_t)^{\top}$, where \top denotes the transpose.

The parameter $\boldsymbol{\theta}$ of the regression model is learned by *Weighted Least-Squares (WLS)* with weight function $w(\mathbf{x})$ (>0 for all $\mathbf{x} \in \mathcal{D}$), i.e.,

$$\hat{\boldsymbol{\theta}}_{\text{W}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{i=1}^{n_{\text{tr}}} w(\mathbf{x}_i^{\text{tr}}) (\hat{f}(\mathbf{x}_i^{\text{tr}}) - y_i^{\text{tr}})^2 \right], \quad (2)$$

where the subscript ‘W’ denotes ‘Weighted’. Let \mathbf{X} be the $n_{\text{tr}} \times t$ matrix with the (i, ℓ) -th element

$$X_{i,\ell} = \varphi_{\ell}(\mathbf{x}_i^{\text{tr}}), \quad (3)$$

and let \mathbf{W} be the $n_{\text{tr}} \times n_{\text{tr}}$ diagonal matrix with the i -th diagonal element

$$W_{i,i} = w(\mathbf{x}_i^{\text{tr}}). \quad (4)$$

Then $\hat{\boldsymbol{\theta}}_{\text{W}}$ is given in a closed form as

$$\hat{\boldsymbol{\theta}}_{\text{W}} = \mathbf{L}_{\text{W}} \mathbf{y}^{\text{tr}}, \quad (5)$$

where¹

$$\mathbf{L}_{\text{W}} = (\mathbf{X}^{\top} \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{W}, \quad (7)$$

$$\mathbf{y}^{\text{tr}} = (y_1^{\text{tr}}, y_2^{\text{tr}}, \dots, y_{n_{\text{tr}}}^{\text{tr}})^{\top}. \quad (8)$$

¹ Although we can obtain the analytic solution for Eq. (2), we often face with numerical instability when computing the inverse of the matrix $\mathbf{X}^{\top} \mathbf{W} \mathbf{X}$ in Eq. (7). To avoid this problem, we practically employ a regularization technique (Hoerl & Kennard, 1970; Poggio & Girosi, 1990; Tikhonov & Arsenin, 1977), i.e., Eq. (7) is replaced with

$$\mathbf{L}_{\text{W}} = (\mathbf{X}^{\top} \mathbf{W} \mathbf{X} + \gamma \mathbf{I}_t)^{-1} \mathbf{X}^{\top} \mathbf{W}, \quad (6)$$

where γ is a small positive scalar called the regularization parameter and \mathbf{I}_t is the $t \times t$ identity matrix. In our experiments, we set $\gamma = 10^{-10}$.

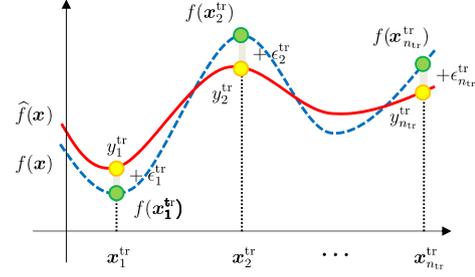


Fig. 1. Regression problem.

Note that the solution $\hat{\boldsymbol{\theta}}_{\text{W}}$ does not depend on the constant scaling of the weight function $w(\mathbf{x})$.

We adopt the squared-loss as the generalization error, i.e., the goodness of a learned function $\hat{f}(\mathbf{x})$ is measured by

$$G = \int (\hat{f}(\mathbf{x}^{\text{te}}) - f(\mathbf{x}^{\text{te}}))^2 p_{\text{te}}(\mathbf{x}^{\text{te}}) d\mathbf{x}^{\text{te}}, \quad (9)$$

where $p_{\text{te}}(\mathbf{x})$ (>0 for all $\mathbf{x} \in \mathcal{D}$) is a probability density function of test input points. The above formulation is summarized in Fig. 1.

Below, we consider a pool-based active learning situation where we are given a ‘pool’ of test input points $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}}$ drawn independently from $p_{\text{te}}(\mathbf{x})$ and choose training input points $\{\mathbf{x}_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}}$ from the pool so that the generalization error (9) is minimized.

2.2. P-ALICE

P-ALICE (Pool-based Active Learning using Importance-weighted least-squares learning based on Conditional Expectation of the generalization error; Sugiyama & Nakajima, 2009) is a pool-based active learning method that chooses training input points one by one (i.e., sampling *without* replacement), with probability proportional to a user-designed *resampling bias function* $b(\mathbf{x})$. Mathematically, the resampling bias function is a strictly positive function defined over the pool of test input samples. P-ALICE finds the resampling bias function that minimizes a generalization error estimator

$$J = \operatorname{tr}(\hat{\mathbf{U}} \mathbf{L}_{\text{W}} \mathbf{L}_{\text{W}}^{\top}),$$

where $\hat{\mathbf{U}}$ is the $t \times t$ matrix with the (ℓ, ℓ') -th element

$$\hat{U}_{\ell,\ell'} = \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \varphi_{\ell}(\mathbf{x}_j^{\text{te}}) \varphi_{\ell'}(\mathbf{x}_j^{\text{te}}), \quad (10)$$

and

$$w(\mathbf{x}_j^{\text{te}}) \propto \frac{1}{b(\mathbf{x}_j^{\text{te}})} \quad (11)$$

is used as a weight in WLS (2).

A more detailed description of P-ALICE is given in [Appendix](#).

3. Improving importance estimation in P-ALICE

In this section, we point out a weakness of P-ALICE and propose an alternative approach.

3.1. Weakness of P-ALICE

In P-ALICE, the importance weight $w(\mathbf{x}_j^{\text{te}})$ is set as Eq. (11), which implies that the training input density $p_{\text{tr}}(\mathbf{x}_j^{\text{te}})$ is proportional to the product of the test input density $p_{\text{te}}(\mathbf{x}_j^{\text{te}})$ and a resampling bias function $b(\mathbf{x})$, i.e.,

$$p_{\text{tr}}(\mathbf{x}_j^{\text{te}}) \propto p_{\text{te}}(\mathbf{x}_j^{\text{te}}) b(\mathbf{x}_j^{\text{te}}).$$

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات