

Lighting Condition Adaptation for Perceived Age Estimation

Kazuya Ueki

NEC Soft, Ltd., Japan

Masashi Sugiyama

Tokyo Institute of Technology, Japan

Yasuyuki Ihara

NEC Soft, Ltd., Japan

Abstract

Over the recent years, a great deal of effort has been made to age estimation from face images. It has been reported that age can be accurately estimated under controlled environment such as frontal faces, no expression, and static lighting conditions. However, it is not straightforward to achieve the same accuracy level in real-world environment because of considerable variations in camera settings, facial poses, and illumination conditions. In this paper, we apply a recently-proposed machine learning technique called *covariate shift adaptation* to alleviating lighting condition change between laboratory and practical environment. Through real-world age estimation experiments, we demonstrate the usefulness of our proposed method.

Keywords

face recognition, age estimation, covariate shift adaptation, lighting condition change, Kullback-Leibler importance estimation procedure, importance-weighted regularized least-squares

1 Introduction

In recent years, demographic analysis in public places such as shopping malls and stations is attracting a great deal of attention. Such demographic information is useful for various purposes including designing effective marketing strategies and targeted advertisement based on customers' gender and age. For this reason, a number of approaches have been explored for age estimation from face images [2, 3], and several databases became publicly available recently [1, 6].

The recognition performance of age prediction systems is significantly influenced, e.g., by the type of camera, camera calibration, and lighting variations, and the publicly available databases were mainly collected in semi-controlled environment. For this reason, existing age prediction systems built upon such databases tend to perform poorly in real-world environment.

The situation where training and test data are drawn from different distributions is called *covariate shift* [8, 11, 12]. In this paper, we formulate the problem of age estimation in real-world environment as a supervised learning problem under covariate shift. Within the covariate shift framework, a method called *importance-weighted least-squares* allows us to alleviate the influence of environmental changes, by assigning higher weights to data samples having high test input densities and low training input densities. We demonstrate through real-world experiments that age estimation based on covariate shift adaptation achieves higher accuracy than baseline approaches.

2 Proposed Method

In this section, we formulate the problem of age estimation as a supervised learning problem under covariate shift, and then describe our proposed method.

2.1 Formulation

Throughout this paper, we perform age estimation based not on subjects' real age, but on their *perceived* age. Thus, the 'true' age of the subject y is defined as the average perceived age evaluated by those who observed the subject's face images (the value is rounded-off to the nearest integer).

Let us consider a regression problem of estimating the age y^* of subject \mathbf{x} (face features). We use the following model for regression.

$$f(\mathbf{x}; \boldsymbol{\alpha}) = \sum_{i=1}^{n_{tr}} \alpha_i K(\mathbf{x}, \mathbf{x}_i^{tr}), \quad (1)$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{n_{tr}})^\top$ is a model parameter, $^\top$ denotes the transpose, and $K(\mathbf{x}, \mathbf{x}')$ is a *positive definite kernel* [7].

Suppose we are given labeled training data $\{(\mathbf{x}_i^{tr}, y_i^{tr})\}_{i=1}^{n_{tr}}$. A standard approach to learning the model parameter $\boldsymbol{\alpha}$ would be *regularized least-squares* [4].

$$\min_{\boldsymbol{\alpha}} \left[\frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} (y_i^{tr} - f(\mathbf{x}_i^{tr}; \boldsymbol{\alpha}))^2 + \lambda \|\boldsymbol{\alpha}\|^2 \right], \quad (2)$$

where $\|\cdot\|$ denotes the Euclidean norm, and $\lambda (> 0)$ is the regularization parameter to avoid overfitting.

Below, we explain that merely using regularized least-squares is not appropriate in real-world perceived age prediction, and show how to cope with this problem.

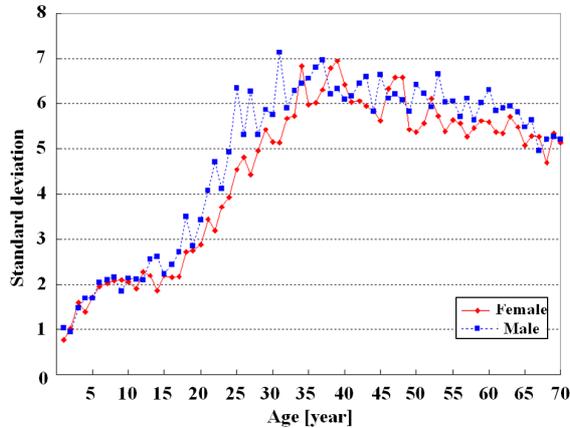


Figure 1: The relation between subjects’ perceived age y^* (horizontal axis) and its standard deviation (vertical axis)

2.2 Incorporating Age Perception Characteristics

Human age perception is known to have heterogeneous characteristics, e.g., it is rare to misjudge the age of a 5-year-old child as 15 years old, but the age of a 35-year-old person is often misjudged as 45 years old. In order to quantify this phenomenon, a large-scale questionnaire survey was carried out in [15]: Each of 72 volunteers were asked to give age labels y to approximately 1000 face images. Figure 1 depicts the relation between subjects’ perceived age y^* and its standard deviation. This shows that the perceived age deviation tends to be small in younger age brackets and large in older age brackets.

In order to match characteristics of our age prediction system to those of human age perception, we weight the goodness-of-fit term in Eq.(2) according to the inverse variance of the perceived age:

$$\min_{\boldsymbol{\alpha}} \left[\frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \frac{(y_i^{tr} - f(\mathbf{x}_i^{tr}; \boldsymbol{\alpha}))^2}{w_{age}(y_i^{tr})^2} + \lambda \|\boldsymbol{\alpha}\|^2 \right], \quad (3)$$

where $w_{age}(y)$ is the standard deviation of the perceived age (see Figure 1 again).

2.3 Coping with Lighting Condition Change

When designing age estimation systems, the environment of recording training face images is often different from the test environment in terms of lighting conditions. Typically, training data are recorded indoors such as a studio with appropriate illumination. On the other hand, in real-world environment, lighting conditions have considerable varieties, e.g., strong sunlight might be cast from a side of faces or there is no enough light. In such situations, age estimation accuracy is significantly degraded.

Let $p_{tr}(\mathbf{x})$ be the training input density and $p_{te}(\mathbf{x})$ be the test input density. When these two densities are different, it would be natural to emphasize the influence of train-

ing samples $(\mathbf{x}_i^{tr}, y_i^{tr})$ which have high similarity to data in the test environment. Such adjustment can be systematically carried out as follows [8, 11, 12]:

$$\min_{\boldsymbol{\alpha}} \left[\frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} w_{imp}(\mathbf{x}_i^{tr}) \frac{(y_i^{tr} - f(\mathbf{x}_i^{tr}; \boldsymbol{\alpha}))^2}{w_{age}(y_i^{tr})^2} + \lambda \|\boldsymbol{\alpha}\|^2 \right], \quad (4)$$

i.e., the goodness-of-fit term in Eq.(3) is weighted according to the *importance function*:

$$w_{imp}(\mathbf{x}) = \frac{p_{te}(\mathbf{x})}{p_{tr}(\mathbf{x})}.$$

The solution of Eq.(4) can be obtained analytically by

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K}^{tr} \mathbf{W}^{tr} \mathbf{K}^{tr} + n_{tr} \lambda \mathbf{I}_{n_{tr}})^{-1} \mathbf{K}^{tr} \mathbf{W}^{tr} \mathbf{y}^{tr}, \quad (5)$$

where \mathbf{K}^{tr} is the kernel matrix whose (i, i') -th element is defined by

$$K_{i,i'}^{tr} = K(\mathbf{x}_i^{tr}, \mathbf{x}_{i'}^{tr}),$$

\mathbf{W}^{tr} is the n_{tr} -dimensional diagonal matrix with (i, i) -th diagonal element defined by

$$W_{i,i}^{tr} = \frac{w_{imp}(\mathbf{x}_i^{tr})}{w_{age}(y_i^{tr})^2},$$

$\mathbf{I}_{n_{tr}}$ is the n_{tr} -dimensional identity matrix, and \mathbf{y}^{tr} is the n_{tr} -dimensional vector with i -th element defined by y_i^{tr} .

When the number of training data n_{tr} is large, we may reduce the number of kernels in Eq.(1) so that the inverse matrix in Eq.(5) can be computed with limited memory; or we may compute the solution numerically by a *stochastic gradient-descent method*.

2.4 Importance-Weighted Cross-Validation (IWCV)

In supervised learning, the choice of models (for example, the basis functions and the regularization parameter) is crucial for obtaining better performance. *Cross-validation* (CV) would be one of the most popular techniques for model selection [9]. CV has been shown to give an *almost* unbiased estimate of the generalization error with finite samples [7], but such almost unbiasedness is no longer fulfilled under covariate shift.

To cope with this problem, a variant of CV called *importance-weighted CV* (IWCV) has been proposed [11]. Let us randomly divide the training set

$$\mathcal{Z} = \{(\mathbf{x}_i^{tr}, y_i^{tr})\}_{i=1}^{n_{tr}}$$

into T disjoint non-empty subsets $\{\mathcal{Z}_t\}_{t=1}^T$ of (approximately) the same size. Let $f_{\mathcal{Z}_t}(\mathbf{x})$ be a function learned from $\mathcal{Z} \setminus \mathcal{Z}_t$ (i.e., without \mathcal{Z}_t). Then the T -fold IWCV (IWCV) estimate of the generalization error is given by

$$\frac{1}{T} \sum_{t=1}^T \frac{1}{|\mathcal{Z}_t|} \sum_{(\mathbf{x}, y) \in \mathcal{Z}_t} \frac{w_{imp}(\mathbf{x})}{w_{age}(y)^2} (f_{\mathcal{Z}_t}(\mathbf{x}) - y)^2,$$

Table 1: Pseudo code of KLIEP. ‘./’ indicates the element-wise division. Inequalities and the ‘max’ operation for vectors are applied in an element-wise manner.

Input: $\{\mathbf{x}_i^{tr}\}_{i=1}^{n_{tr}}, \{\mathbf{x}_j^{te}\}_{j=1}^{n_{te}}$
Output: $\hat{w}(\mathbf{x})$

Choose $\{\mathbf{c}_k\}_{k=1}^b$ as a subset of $\{\mathbf{x}_j^{te}\}_{j=1}^{n_{te}}$;
 $A_{j,k} \leftarrow \exp(-\|\mathbf{x}_j^{te} - \mathbf{c}_k\|^2 / (2\gamma^2))$;
 $b_k \leftarrow \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \exp(-\|\mathbf{x}_i^{tr} - \mathbf{c}_k\|^2 / (2\gamma^2))$;
Initialize $\boldsymbol{\beta} (> \mathbf{0})$ and ε ($0 < \varepsilon \ll 1$);
Repeat until convergence
 $\boldsymbol{\beta} \leftarrow \varepsilon A^\top (\mathbf{1} ./ A \boldsymbol{\beta})$;
 $\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} + (1 - \mathbf{b}^\top \boldsymbol{\beta}) \mathbf{b} / (\mathbf{b}^\top \mathbf{b})$;
 $\boldsymbol{\beta} \leftarrow \max(\mathbf{0}, \boldsymbol{\beta})$;
 $\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} / (\mathbf{b}^\top \boldsymbol{\beta})$;
end

where $|\mathcal{Z}_t|$ denotes the number of samples in the subset \mathcal{Z}_t .

It was proved that IWCV gives an *almost* unbiased estimate of the generalization error even under covariate shift [11].

2.5 Kullback-Leibler Importance Estimation Procedure (KLIEP)

In order to compute the solution (5) or performing IWCV, we need the importance weights $w_{imp}(\mathbf{x}_i^{tr}) = p_{te}(\mathbf{x}_i^{tr}) / p_{tr}(\mathbf{x}_i^{tr})$, which include two probability densities $p_{tr}(\mathbf{x})$ and $p_{te}(\mathbf{x})$. However, since density estimation is a hard problem, a two-stage approach of first estimating $p_{tr}(\mathbf{x})$ and $p_{te}(\mathbf{x})$ and then taking their ratio may not be reliable. Here we describe a method called *Kullback-Leibler Importance Estimation Procedure* (KLIEP) [12], which allows us to directly estimate the importance function $w_{imp}(\mathbf{x})$ without going through density estimation of $p_{tr}(\mathbf{x})$ and $p_{te}(\mathbf{x})$.

Let us model $w_{imp}(\mathbf{x})$ using the following model:

$$\hat{w}_{imp}(\mathbf{x}) = \sum_{k=1}^b \beta_k \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_k\|^2}{2\gamma^2}\right), \quad (6)$$

where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_b)^\top$ is a parameter, and $\{\mathbf{c}_k\}_{k=1}^b$ is a subset of test input samples $\{\mathbf{x}_j^{te}\}_{j=1}^{n_{te}}$. Using the model $\hat{w}_{imp}(\mathbf{x})$, we can estimate the test input density $p_{te}(\mathbf{x})$ by

$$\hat{p}_{te}(\mathbf{x}) = \hat{w}_{imp}(\mathbf{x}) p_{tr}(\mathbf{x}). \quad (7)$$

We determine the parameter $\boldsymbol{\beta}$ in the model (7) so that the Kullback-Leibler divergence



Figure 2: Examples of face images under different lighting conditions (left: standard lighting, middle: dark, right: strong light from a side)

from p_{te} to \hat{p}_{te} is minimized:

$$\begin{aligned} KL(p_{te} \parallel \hat{p}_{te}) &= \int p_{te}(\mathbf{x}) \log \frac{p_{te}(\mathbf{x})}{\hat{p}_{te}(\mathbf{x})} d\mathbf{x} \\ &= \int p_{te}(\mathbf{x}) \log \frac{p_{te}(\mathbf{x})}{p_{tr}(\mathbf{x})} d\mathbf{x} - \int p_{te}(\mathbf{x}) \log \hat{w}_{imp}(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

We ignore the first term (which is a constant) and impose $\hat{w}_{imp}(\mathbf{x})$ to be non-negative and normalized. Then we obtain the following convex optimization problem:

$$\begin{aligned} \max_{\beta} & \left[\sum_{j=1}^{n_{te}} \log \left(\sum_{k=1}^b \beta_k \exp \left(-\frac{\|\mathbf{x}_j^{te} - \mathbf{c}_k\|^2}{2\gamma^2} \right) \right) \right], \\ \text{s.t.} & \begin{cases} \beta_k \geq 0 \text{ for } k = 1, \dots, b, \\ \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \left(\sum_{k=1}^b \beta_k \exp \left(-\frac{\|\mathbf{x}_i^{tr} - \mathbf{c}_k\|^2}{2\gamma^2} \right) \right) = 1. \end{cases} \end{aligned}$$

A pseudo code of KLIEP is described in Table 1. The tuning parameter γ can be optimized based on *likelihood cross-validation* (LCV) [12].

3 Empirical Evaluation

In this section, we experimentally evaluate the performance of the proposed method using in-house face-age datasets.

We use the face images recorded under 17 different lighting conditions: for instance, average illuminance from above is approximately 1000 lux and 500 lux from the front in the standard lighting condition, 250 lux from above and 125 lux from the front in the dark setting, and 190 lux from left and 750 lux from right in another setting (see Figure 2). Note that these 17 lighting conditions are diverse enough to cover real-world lighting conditions. Images were recorded as movies with camera at depression angle 15 degrees. The number of subjects is approximately 500 (250 for each gender). We used a

face detector for localizing the two eye-centers, and then rescaled the image to 64×64 pixels. The number of face images in each environment is about 2500 (5 face images \times 500 subjects).

As pre-processing, a neural network feature extractor [14] was used to extract 100-dimensional features from 64×64 face images. The kernel regression model (1) with the following Gaussian kernel was employed for the extracted 100-dimensional data:

$$K_\sigma(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right).$$

We constructed the male/female age prediction models only using male/female data, assuming that gender classification had been correctly carried out.

We split the 250 subjects into the *training set* (200 subjects) and the *test set* (50 subjects). The training set was used for training the kernel regression model (1), and the test set was used for evaluating its generalization performance. For the test samples $\{(\mathbf{x}_i^{te}, y_i^{te})\}_{i=1}^{n_{te}}$ taken from the test set in the environment with strong light from a side, age-weighted mean square error (WMSE)

$$\text{WMSE} = \frac{1}{n_{te}} \sum_{i=1}^{n_{te}} \frac{(y_i^{te} - f(\mathbf{x}_i^{te}; \hat{\boldsymbol{\alpha}}))^2}{w_{age}(y_i^{te})^2}$$

was calculated as a performance measure. The training test sets were shuffled 5 times in such a way that each subject was selected as a test sample once. The final performance was evaluated based on the average WMSE over the 5 trials.

We compared the performance of the proposed method with the two baseline methods:

Baseline method 1: Training samples were taken only from the standard lighting condition and age-weighted regularized least-squares (3) was used for training.

Baseline method 2: Training samples were taken from all 17 different lighting conditions and age-weighted regularized least-squares (3) was used for training.

The importance weights were not used in these baseline methods. The Gaussian width σ and the regularization parameter λ were determined based on 4-fold CV over WMSE, i.e., the training set was further divided into a training part (150 subjects) and a validation part (50 subjects).

In the proposed method, training samples were taken from all 17 different lighting conditions (which is the same as the baseline method 2). The importance weights were estimated by KLIEP using the training samples and additional *unlabeled* test samples; the hyper-parameter γ in KLIEP was determined based on 2-fold LCV [12]. We then computed the average importance score over different samples for each lighting condition and used the average importance score for training the regression model. The Gaussian width σ and the regularization parameter λ in the regression model were determined based on 4-fold IWCV [11].

Table 2: The test performance measured by WMSE.

| | Male | Female |
|------------------------|-------------|-------------|
| Baseline method 1 | 2.83 | 6.51 |
| Baseline method 2 | 2.64 | 4.40 |
| Proposed method | 2.54 | 3.90 |

Table 2 summarizes the experimental results, showing that, for both male and female data, the baseline method 2 is better than the baseline method 1 and the proposed method is better than the baseline method 2. This illustrates the effectiveness of the proposed method. Note that WMSE for female subjects is substantially larger than that for male subjects. The reason for this would be that female subjects tend to have more divergence such as short/long hair and with/without makeup, which makes prediction harder [16].

4 Summary and Future Works

Lighting condition change is one of the critical causes of performance degradation in age prediction from face images. In this paper, we proposed to employ a machine learning technique called *covariate shift adaptation* for alleviating the influence of lighting condition change. We demonstrated the effectiveness of our proposed method through real-world perceived age prediction experiments.

In the experiments in Section 3, test samples were collected from a particular lighting condition, and samples from the same lighting condition were also included in the training set. Although we believe this setup to be practical, it would be interesting to evaluate the performance of the proposed method when no overlap in the lighting conditions exists between training and test data.

In principle, the covariate shift framework allows us to incorporate not only lighting condition change, but also various types of environment change such as face pose variation and camera setting change. In our future work, we will investigate whether the proposed approach is still useful in such challenging scenarios.

Recently, novel approaches to density ratio estimation for high-dimensional problems have been explored [5, 10, 17, 13]. In our future work, we would like to incorporating these new ideas into our framework of perceived age estimation, and see how the prediction performance can be further improved.

References

- [1] *The FG-NET Aging Database*. <http://www.fgnet.rsunit.com/>.
- [2] Y. Fu, Y. Xu, and T. S. Huang. Estimating human age by manifold analysis of face pictures and regression on aging features. *Proceedings of the IEEE Multimedia and Expo*, pages 1383–1386, 2007.

- [3] G. Guo, G. Mu, Y. Fu, C. Dyer, and T. Huang. A study on automatic age estimation using a large database. *International Conference on Computer Vision in Kyoto (ICCV 2009)*, pages 1986–1991, 2009.
- [4] A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(3):55–67, 1970.
- [5] T. Kanamori, S. Hido, and M. Sugiyama. A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 10:1391–1445, 2009.
- [6] K. J. Ricanek and T. Tesafaye. Morph: A longitudinal image database of normal adult age-progression. *Proceedings of the IEEE 7th International Conference on Automatic Face and Gesture Recognition (FGR 2006)*, pages 341–345, 2006.
- [7] B. Schölkopf and A. J. Smola. *Learning with Kernels*, MIT Press, Cambridge, MA, USA, 2002.
- [8] H. Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2):227–244, 2000.
- [9] M. Stone. Cross-validated choice and assessment of statistical predictions, *Journal of the Royal Statistical Society, Series B*, 36:111–147, 1974.
- [10] M. Sugiyama, M. Kawanabe, P. L. Chui. Dimensionality reduction for density ratio estimation in high-dimensional spaces. *Neural Networks*, 23(1):44–59, 2010.
- [11] M. Sugiyama, M. Krauledat, and K.-R. Müller. Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research*, 8:985–1005, May 2007.
- [12] M. Sugiyama, T. Suzuki, T., S. Nakajima, H. Kashima, P. von Bünau, and M. Kawanabe. Direct importance estimation for covariate shift adaptation. *Annals of the Institute of Statistical Mathematics*, 60(4):699–746, 2008.
- [13] M. Sugiyama, M. Yamada, P. von Bünau, T. Suzuki, T. Kanamori, and M. Kawanabe. Direct density-ratio estimation with dimensionality reduction via least-squares hetero-distributional subspace search, *Neural Networks*, to appear.
- [14] F. H. C. Tivive and A. Bouzerdoum. A gender recognition system using shunting inhibitory convolutional neural networks. *Proceedings of the International Joint Conference on Neural Networks (IJCNN '06)*, pages 5336–5341, 2006.
- [15] K. Ueki, M. Sugiyama, Y. Ihara. A semi-supervised approach to perceived age prediction from face images. *IEICE Transactions on Information and Systems*, to appear.

- [16] K. Ueki, M. Miya, T. Ogawa, T. Kobayashi. Class distance weighted locality preserving projection for automatic age estimation. *Proceedings of IEEE International Conference on Biometrics: Theory, Applications and Systems (BTAS 2008)*, pages 1–5, 2008.
- [17] M. Yamada, M. Sugiyama, G. Wichern, and J. Simm. Direct importance estimation with a mixture of probabilistic principal component analyzers. *IEICE Transactions on Information and Systems*, E93-D(10), 2846–2849, 2010.