

Improving Tone Prediction Accuracy in Calibration for Color Electrophotography Part II – Principal Component Regression

Yan-Fu Kuo^a, Chao-Lung Yang^b, George T.-C. Chiu^a, Yuehwern Yih^b, and Jan P. Allebach^c; *a: School of Industrial Engineering, b: School of Mechanical Engineering, c: School of Electrical and Computer Engineering, Purdue University, W Lafayette, IN 47907*

Abstract

This paper presents results of improved tone prediction accuracy in calibration through a principal component analysis based regression approach for color electrophotography (EP). During calibration, multiple color patches of the same primary color at different halftone levels are printed on a belt with black exterior and are measured using on-board sensors. Regression models are often developed to predict the primary color tone values on output media from these on-board sensor measurements. The prediction accuracy of the regression models directly impact the quality and consistency of the calibration. Analyses have revealed a high degree of correlation among the color patch measurements, which results in using multicollinear measurements as explanatory variables during regression analysis to identify model coefficients. It is well known that using collinear explanatory variables during regression analysis will result in sub-optimal model coefficients that will degrade the prediction accuracy. In this study, a principle component regression (PCR) approach is applied to tackle the potential issue with collinear measurement data in model coefficient estimation. The experimental results show the resulting PCR models provides 25% improvement on average in root-mean-squared predication accuracy over separate ordinary least square regression models.

Introduction

A color electrophotographic (EP) printing system typically uses four primary colors – cyan, magenta, yellow, and black. Calibrations are performed to maintain consistent color reproduction under different operating conditions. During a calibration, multiple patches of different halftone levels of the same primary color are printed on an intermediate media and are measured with on-board sensors. Calibration models are used to predict the primary color tone values on output media from these on-board sensor measurements. In this study, our aim is to improve the prediction accuracy of the calibration models through a principal component regression (PCR) approach for color EP systems.

The calibration models are developed with life test data. In typical life test, additional color patches are printed on output media immediately following calibration. Their tone values, output tone values, are measured off-line with measurement devices, such as photometers. Calibration models are developed to map the on-board sensor measurements to the output tone values [1]. During calibration, output tone values are predicted based on the on-board sensor measurements using the calibration model. Appropriate tone correction is performed to regulate the color reproduction.

Halftoned color images composed of arrays of closely spaced micro-dots. Changes in operating conditions or different EP parameter settings will impact the sizes of the micro-dots. Assuming operating conditions and EP parameter settings have consistent impacts on the sizes of the micro-dots for a given print, the effects of micro-dot size fluctuation can be detected by the on-board sensor from different halftone patches of the same color. This results in increased correlation among the on-board sensor measurements, i.e., multicollinearities. It is well known that using the collinear measurements as explanatory variables to identify model coefficients directly through ordinary least-square regression (OLSR) will result in sub-optimal model coefficients that will degrade prediction accuracy [2]. Hence, existing calibration models are developed using a single-response regression approach, i.e., the output tone value at a particular halftone level is regressed only with the on-board sensor measurement at the same halftone level.

Recent researches in regression analysis have shown improved prediction accuracy of regression models using multiple explanatory variables as compared to single-response regression models [3, 4]. To address the multicollinearities associated with multiple on-board sensor measurements, in this study, a principal component regression (PCR) approach is proposed [5]. PCR avoids the numerical issues associated with the OLSR by transforming multicollinear sensor measurements into a set of orthogonal principal components (PC) basis. In addition, it achieves biased regression by determining an optimal subset of PC's to be retained while discarding PC's that are less statistically significant. To illustrate the utility of the proposed approach, a first-order linear calibration model for an off-the-shelf in-line color EP printer is developed using existing life test data. Cross-validation results demonstrate 25% improvement in prediction accuracy compared with the existing calibration model.

Method

Calibration Model

Since each primary color is reproduced independently for a single-pass color EP process, a calibration model is developed for each primary color. A calibration model can be written as $y = f(z, d)$, where y is tone values on paper, z is on-board sensor measurements, and d is uncontrollable but measurable factors/disturbances, such as temperature and humidity. The tone values, y , are the measured intensities of the reproduced color patches printed at designated halftone levels. In this study, a tone value is defined as the Euclidian distance (ΔE) in CIE $L^*a^*b^*$ space between the color point of a primary color printed at a particular halftone level and the substrate appearance color. A static linear calibration model is used in this study.

Problem Formulation

Life test data are used to identify the calibration models. For one observation, a set of on-board sensor measurements, measurable disturbances, and the corresponding tone values measured on paper are collected. Denote $z_{ij} \in \mathbb{R}$ as the j^{th} on-board sensor measurement in the i^{th} observation, $y_{ij} \in \mathbb{R}$ as the j^{th} tone value measurement in the i^{th} observation, and $d_{ij} \in \mathbb{R}$ as the j^{th} measurable disturbances in the i^{th} observation. Denote \mathbf{G} as the calibration model. The calibration model, \mathbf{G} , is a linear transformation relating the tone value measurements, \mathbf{y} , to the sensor measurements, \mathbf{z} , and the disturbances, \mathbf{d} .

Consider p on-board sensor measurements, q measurable disturbances, and w tone value measurements are made in one observation, and n observations are made in the life test. Denote $\mathbf{Z} = [z_{ij}] \in \mathbb{R}^{n \times p}$ as the sensor measurement matrix and $\mathbf{D} = [d_{ij}] \in \mathbb{R}^{n \times q}$ as the measurable disturbance matrix. Let $\mathbf{X} \in \mathbb{R}^{n \times r}$ denote the explanatory variable matrix, which is the augmented matrix consisting of matrices \mathbf{Z} and \mathbf{D} , i.e., $\mathbf{X} = [\mathbf{Z} | \mathbf{D}]$ and $r = p + q$. Denote $\mathbf{Y} = [y_{ij}] \in \mathbb{R}^{n \times w}$ as the response variable matrix containing the tone value measurements. Denote $\mathbf{G} \in \mathbb{R}^{r \times w}$ as the calibration model that can be written as $\mathbf{Y} = \mathbf{X}\mathbf{G}$. Note that the matrices are assumed to be centered column-wise. Hence no constant or intercept term is required in the regression model development.

Ordinary Least Square Regression

Consider a standard multivariate regression model:

$$\mathbf{Y} = \mathbf{X}\mathbf{G} + \mathbf{E},$$

where the error matrix, \mathbf{E} , satisfies the usual assumption of being independent and identically-distributed. The number of observations is usually larger than the number of calibration color patches printed in a calibration, i.e., $n \gg r$. The OLSR solution to the over-determined problem stated above is to minimize the loss function, i.e.,

$$\mathbf{G} = \arg \min \|\mathbf{Y} - \mathbf{X}\mathbf{G}\|^2 = \arg \min \left\| \mathbf{Y} - [\mathbf{Z} | \mathbf{D}] \begin{bmatrix} \mathbf{G}_z \\ \mathbf{G}_d \end{bmatrix} \right\|^2, \quad (1)$$

where the calibration model, \mathbf{G} , can be split into two matrices \mathbf{G}_z and \mathbf{G}_d with proper dimensions corresponding to the sensor measurement matrix, \mathbf{Z} , and the measurable disturbance matrix, \mathbf{D} , respectively. The OLSR solution to Eq. (1) is given by

$$\mathbf{G} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (2)$$

The explanatory variable matrix, \mathbf{X} , is not of full rank since the column vectors are collinear. The calculation of the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ is computational challenging especially with matrices that with lower conditioning number. This yields larger variance in model coefficient estimation.

Variance Inflation Factor

Variance inflation factor [6] (VIF) is commonly used to detect multicollinearity among explanatory variables. It is defined as

$$(VIF)_j = \frac{1}{1 - R_j^2}, \quad (3)$$

where R_j^2 stands for the unadjusted coefficient of determination of the j^{th} explanatory variable when it is predicted by the other explanatory variables included in the model. Suppose the j^{th} explanatory variable is linearly correlated to any the other explanatory in the model, R_j^2 is large and, consequently, the VIF value is large. Values of VIF that exceed 10 are often regarded as indicating strong multicollinearity [7] among the explanatory variables and OLSR may not be a good regression approach.

Principal Component Regression

The key idea of PCR is to linearly transform the multicollinear sensor measurement matrix, \mathbf{Z} , to a principle component (PC) matrix that consists of a set of orthogonal vectors. The model coefficient estimation can be directly carried out following Eq. (2). A singular value decomposition (SVD) on the sensor measurement matrix, \mathbf{Z} , is performed as the first step to calculate the PC matrix, i.e.,

$$\mathbf{Z} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^T,$$

where $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p) \in \mathbb{R}^{n \times p}$ is a diagonal matrix of singular values σ_i associated with the principal component PC_i , and $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{p \times p}$ are left and right unitary matrices of corresponding singular vectors, \mathbf{u}_i and \mathbf{v}_i , respectively. The PC matrix, $\mathbf{\Psi} \in \mathbb{R}^{n \times p}$, can be obtained by multiplying the sensor measurement matrix, \mathbf{Z} , with the right unitary matrix, \mathbf{V} , i.e., $\mathbf{\Psi} = \mathbf{Z}\mathbf{V}$. The principal component PC_i is a linear combination of the raw sensor measurements with the coefficients in the associated row vector, \mathbf{v}_i .

Next, the PC matrix is augmented with the disturbance matrix as the explanatory variable matrix, i.e., $\mathbf{X} = [\mathbf{\Psi} | \mathbf{D}]$, in the subsequent multivariate regression. The loss function is

$$\mathbf{\Gamma} = \arg \min \left\| \mathbf{Y} - [\mathbf{\Psi} | \mathbf{D}] \begin{bmatrix} \mathbf{\Gamma}_\Psi \\ \mathbf{\Gamma}_D \end{bmatrix} \right\|^2, \quad (4)$$

where $\mathbf{\Gamma} \in \mathbb{R}^{r \times w}$ is the coefficient matrix to be determined. The coefficient matrix, $\mathbf{\Gamma}$, can be split into two matrices $\mathbf{\Gamma}_z$ and $\mathbf{\Gamma}_d$ with proper dimensions corresponding to the PC matrix, $\mathbf{\Psi}$, and the disturbance measurement matrix, \mathbf{D} , respectively. Since the PC matrix, $\mathbf{\Psi}$, is of full rank, the solution of the coefficient matrix, $\mathbf{\Gamma}$, in Eq. (4) can be carried out directly following Eq. (2). Matching the response variable matrix, \mathbf{Y} , in Eq. (1) and (4), one can obtain

$$\mathbf{Y} = \mathbf{\Psi} \mathbf{V}^T \mathbf{V}^T \mathbf{\Gamma}_\Psi + \mathbf{D} \mathbf{\Gamma}_D = \mathbf{Z} \mathbf{G}_z + \mathbf{D} \mathbf{G}_d. \quad (5)$$

Therefore the calibration model is

$$\mathbf{G} = \begin{bmatrix} \mathbf{G}_z \\ \mathbf{G}_d \end{bmatrix} = \begin{bmatrix} \mathbf{V} \mathbf{\Gamma}_\Psi \\ \mathbf{\Gamma}_D \end{bmatrix}.$$

Biased Principal Component Regression

The noise in the sensor measurements can result in bias in regression analysis and, consequently, can increase the uncertainty in model coefficient estimation. Biased PCR excludes the PC's regarded as noise from being used in the regression. Assume the noise in the sensor measurement matrix, \mathbf{Z} , is additive. It can be decomposed into two matrices – an exact signal matrix \mathbf{S} , and a noise perturbation matrix \mathbf{N} – so that

$$\mathbf{Z} = \mathbf{S} + \mathbf{N} = \mathbf{U}_S \mathbf{\Sigma}_S \mathbf{V}_S^T + \mathbf{U}_N \mathbf{\Sigma}_N \mathbf{V}_N^T,$$

where $\mathbf{\Sigma}_S$, \mathbf{U}_S , and \mathbf{V}_S , and $\mathbf{\Sigma}_N$, \mathbf{U}_N , and \mathbf{V}_N are the singular value matrix, left matrix, and right matrix associated with the SVD of the signal matrix, \mathbf{S} , and the noise perturbation matrix, \mathbf{N} , respectively. If PC_i is regarded as noise, the corresponding singular value, σ_i , left singular vector, \mathbf{u}_i , and right singular vectors, \mathbf{v}_i , are put to the noise perturbation matrix, \mathbf{N} . Then the biased PC matrix, $\mathbf{\Psi}_S = \mathbf{Z}\mathbf{V}_S$, is used in the subsequent regression.

Principal Component Selection

A forward selection algorithm is used to determine the PC's to be included in the regression. This study utilizes forward selection proposed by Xie and Kalivas [8]. The forward selection tries out the PC's one by one and includes one PC in the model if it is statistically significant to the response variables. Bayesian information criterion (BIC) [6]

$$n \cdot \ln\left(\frac{RSS}{n}\right) + k \cdot \ln(n), \quad (6)$$

is used as the selection criterion, where RSS is residual sum of squares from the estimated model and k is number of the PC's to be included in the forward selection. BIC is known to be more parsimonious compared to other information criterion. Hence the chance of over-fitting can be reduced with using BIC as the selection criterion. The PC selection procedure can be summarized in the following four steps:

- Step 1: Compute all the PC's through SVD.
- Step 2: Determine the first PC producing the minimum selection criterion following Eq. (6). Call this the first PC subset.
- Step 3: Identify the second PC subset as the subset of PC's providing the minimum selection criterion from all possible combinations containing the first PC subset and one more PC that has not been included in the first PC subset. Compute the selection criterion of the second PC subset following Eq. (6).
- Step 4: The process stops when the selection criterion of the second subset is larger than that of the first subset or when all PC's are included in the regression. Otherwise, let the second subset be the first subset and continue from step 3.

Preferably, the PC selection can be performed for an individual response variable to preserve the freedom of PC retention. Then the individual response variable is regressed with the selected PC's. The calibration model is the integration of the coefficients from the regression analyses. The signal matrix consisting of selected PC's for the m^{th} response variable can be expressed as

$$\mathbf{S}^{(m)} = \mathbf{U}_S^{(m)} \mathbf{\Sigma}_S^{(m)} (\mathbf{V}_S^{(m)})^T. \quad (7)$$

The biased PC matrix of the m^{th} response variable, $\mathbf{\Psi}_S^{(m)}$, can be obtained by multiplying the sensor measurements matrix, \mathbf{Z} , with the right unitary matrix, $\mathbf{V}_S^{(m)}$ from Eq. (7), i.e., $\mathbf{\Psi}_S^{(m)} = \mathbf{Z}\mathbf{V}_S^{(m)}$. Let $\mathbf{y}^{(m)} \in \mathbb{R}^n$ denotes the m^{th} column vector in the response variable matrix, \mathbf{Y} . The loss function to be minimized for the m^{th} response variable is

$$\gamma^{(m)} = \arg \min \left\| \mathbf{y}^{(m)} - [\mathbf{\Psi}_S^{(m)} | \mathbf{D}] \begin{bmatrix} \gamma_S^{(m)} \\ \gamma_D^{(m)} \end{bmatrix} \right\|^2, \quad (8)$$

where $\gamma^{(m)} \in \mathbb{R}^r$ is the coefficient vector corresponding to $\mathbf{y}^{(m)}$ to be determined in the regression. The coefficient vector, $\gamma^{(m)}$, can be split into two vectors, $\gamma_S^{(m)}$ and $\gamma_D^{(m)}$, with proper dimensions corresponding to the biased PC matrix, $\mathbf{\Psi}_S^{(m)}$, and the disturbance measurement matrix, \mathbf{D} , respectively. The calibration model can be obtained by combining the product vectors of multiplication of the coefficient vectors, $\gamma_S^{(m)}$, from Eq. (8) and the associated right unitary matrix, $\mathbf{V}_S^{(m)}$, from Eq. (7) with the coefficient vectors, $\gamma_D^{(m)}$, i.e.,

$$\mathbf{G} = \begin{bmatrix} [\mathbf{V}_S^{(1)} \gamma_S^{(1)} | \dots | \mathbf{V}_S^{(w)} \gamma_S^{(w)}] \\ [\gamma_D^{(1)} | \dots | \gamma_D^{(w)}] \end{bmatrix}.$$

Experiment

The proposed PCR procedure is performed on an off-the-shelf one-pass color EP printer. The printer prints and measures nine calibration patches at different halftone levels for each primary color during a calibration, i.e., $p=9$. These halftone levels are labeled as HL_j , $j=1 \dots 9$, from light to dark. Calibration patches identical to those printed in calibration are printed on output media for each primary color immediately following a calibration. Their tone value measurements are made with spectrophotometers (X-Rite® DTP-70). A commercial white paper (Xerox®) is used as the output media.

The experiment is performed on twenty printers with several consumable sets across a wide range of environmental conditions. The temperature ranges from 15 to 30°C, and the relative humidity ranges from 10 to 80%. Totally more than four hundred observations, i.e., $n > 400$, are made. The models are identified following the proposed PCR procedure using Matlab®. Humidity ratio is chosen to be the measurable disturbance.

VIF of the Sensor Measurements

Table 1 lists the variance inflation factor (VIF) of the experimental sensor measurements based on Eq. (3). The VIF values indicate a high degree of multicollinearity among the sensor measurements, especially for those sensor measurements of the color patches printed with halftone levels in the mid-tone and shadow range. Particularly, the VIF value can be as large as fifty for yellow. The large VIF values demonstrate the necessity to conduct the model development through the PCR.

TABLE I: Variance inflation factor values of the sensor measurements at each halftone level (HL)

	Cyan	Magenta	Yellow	Black
HL_1	1.3	1.2	1.3	1.3
HL_2	3.9	2.6	4.6	2.2
HL_3	7.4	5.3	7.6	4.7
HL_4	8.8	7.2	22.5	8.6
HL_5	11.6	13.8	26.3	10.7
HL_6	17.9	14.6	27.0	14.9
HL_7	19.9	16.4	53.0	23.7
HL_8	18.9	16.2	31.9	15.8
HL_9	11.2	7.6	14.1	12.9

Model Comparison

PCR models are developed through the proposed methods and are compared with separate OLSR models in which a tone value is regressed with single sensor measurement and the humidity ratio at each halftone level. A statistical F-test is conducted to compare the two types of models. The results show the PCR models are significantly better than the OLSR models at all halftone levels for all colors at 99% confidence level.

A 10-fold cross-validation (CV) without replacement is performed with the experimental data (see Fig. 1). Overall the PCR models provide 25% improvement on average over the OLSR models.

Conclusion

A PCR method to improve tone prediction accuracy of calibration models for color EP systems is proposed in this work. A high degree of multicollinearity among calibration color patch measurements is verified through experiments and statistical analyses. This motivates using PCR for calibration model identification. The proposed method includes a forward selection algorithm to determine the optimal subset of PC's to be retained in biased PCR. The effectiveness of the proposed PCR method is verified with experimental data collected under different environmental conditions and consumable usage levels. Statistical tests demonstrate the proposed PCR models outperform separate OLSR models. The PCR models provide 25% improvement on average in root-mean-squared prediction accuracy over OLSR models based on cross-validation.

Acknowledgements

We gratefully acknowledge the support from the Hewlett-Packard Company. We would like to specially thank Dennis Abramssohn and Jeff Trask for their valuable guidance in this research.

References

- [1] D. A. Johnson, "Calibration of printing devices," U. S. Patent 6,982,812 (2006).
- [2] D. C. Montgomery and D. J. Friedman, "Prediction using regression models with multicollinear predictor variables," *IIE Transactions*, 25: 73 (1993).
- [3] L. Breiman and J. H. Friedman, "Predicting multivariate responses in multiple linear regression," *J. Roy. Statist. Soc. B*, 59: 3 (1997).
- [4] R. D. De Veaux and L. H. Ungar, "Multicollinearity: a tale of two nonparametric regressions," in *Selecting Models from Data: AI and Statistics*, (Springer-Verlag, New York, NY, 1994), pg. 293-302.
- [5] W. F. Massy, "Principal components regression in exploratory statistical research," *J. of the American Statistical Association*, 60: 234 (1965).
- [6] M. H. Kutner, C. J. Nachtsheim, J. Neter, and W. Li, *Applied linear regression models*, 4th Ed., (McGraw-Hill Irwin, New York, 2004).
- [7] P. J. Curran, S. G. West, and J. F. Finch, "Model-dependent variance inflation factor cutoff values," *Quality Engineering*, 14: 391 (1996).
- [8] Y. Xie and J. Kalivas, "Evaluation of principal component selection methods to form a global prediction model by principal component regression," *Analytica Chimica Acta*, 348: 19 (1997).

Author Biography

Yan-Fu Kuo is currently a Ph.D. student at School of Mechanical Engineering, Purdue University. He worked as an intern engineer at Hewlett-Packard Company in 2007 and 2008, developing calibration tools to improve color consistency for electrophotography printers.

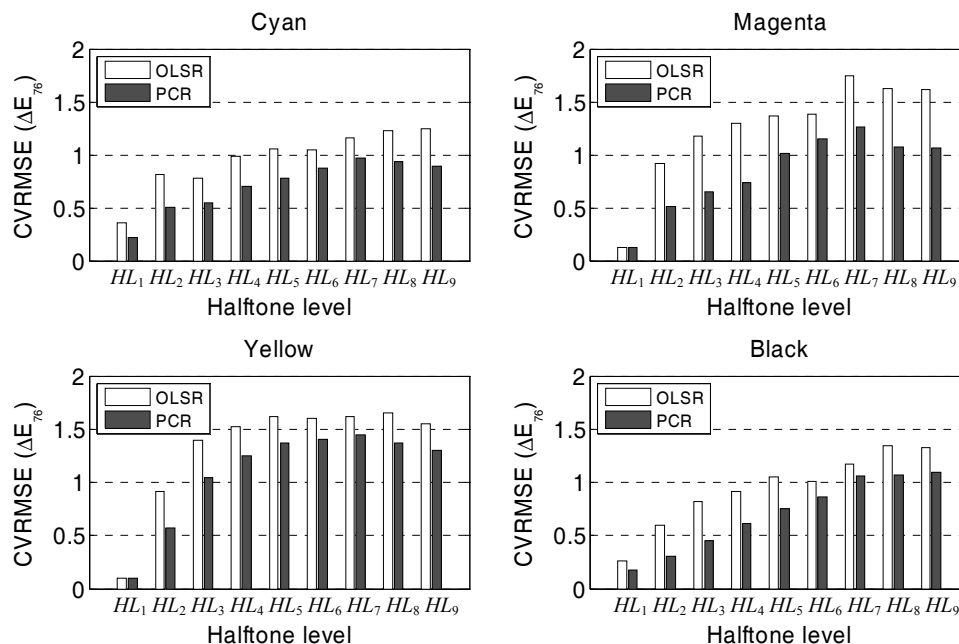


Figure 1. Cross-validation root-mean-squared errors (CVMSE) of the ordinary least square regression (OLSR) models and the principal component regression (PCR) models at each halftone level (HL).