

# Multivariate Linear Regression

## Nail Bed Images

Nov 14, 2022

### 1. Introduction

Anemia is as a life-threatening disease that affects 2 billion, or approximately 1 in 3 people world wide (Mannino 2018). Patients suffering from chronic anemia require frequent monitoring of indicators such as the Hgb concentration in blood to track the progression of their disease. Despite high prevalence, however, the current diagnostic process requires blood tests which causes discomfort and trauma in patients, in addition to incurring high monetary cost, especially in regions of lower socio-economic development (Safiri 2021). Therefore, we aim to create machine learning models that enable non-invasive inexpensive diagnosis using patients' nail bed images to predict their Hgb concentration.

Previous work focuses on predicting Hgb concentration on three regressions of interest: the fingernail beds, the conjunctiva and the palmar creases. We aim to extend the previous research from several perspectives: The previous study found no significant correlation between blue pixel intensity and gold standard measured Hgb, which we aim to investigate further using more complex machine learning algorithms such as random forests. The study finds that machine learning techniques such as convoluted neural networks do not improve Hgb level measurement accuracy given the current sample size of the study population; we would like to refine the algorithm and test on a wider range of data as well as techniques such as neural networks and Bayesian regressions. We aim to develop quality control algorithms accounting for other common irregularities in images of fingernails that could lead to inaccurate Hgb level estimation including presence of abnormal fingernail bed pigmentation, abnormal imaging brightness, and lack of image focus, using convolution of fingernail bed images with edge detection kernels to detect edges within the fingernail beds corresponding to abrupt color changes caused by abnormal fingernail bed pigmentation which results in improved prediction accuracy.

### 2. Methodology

#### 2.1 Multilinear Regression

##### Data

Our data are nail bed images collected from patients enrolled in Dr. Nirmish Shah's clinic. Each patient has four images corresponding to different fingers. The images are processed such that the nail bed is captured in a bounding box while the background of the image is discarded. Colour information of each pixel is extracted from the bounded nail bed images as features.

We know that each pixel can be represented by RGB values, but the RGB colour space contains both colour information and the light information, which is different for each image since photos are taken at different times and settings. To eliminate the inconsistency caused by variation in lightning and background, we used two other colour spaces (HSV, LAB) to separate the colour information from the lightning information. We computed the mean of each value/channel across the bounded nail bed image for each of these three colour spaces (HSV, LAB, RGB) and used them as our model input.

Our response variable is blood hemoglobin concentration associated with each nail bed.

Our predictor variables are:

- mean\_H: mean value of hue (the color component / base pigment) of HSV color space.

- mean\_S: mean value of saturation (amount of color / depth of the pigment / dominance of hue) of HSV color space.
- mean\_V: mean value of value (brightness of the color) of HSV color space.
- mean\_L: mean value of lightness from black to white on a scale of 0 - 100 of LAB color space.
- mean\_A: mean value of representation of greenness to redness on a scale of -128 to +127 of LAB color space.
- mean\_B: mean value of representation of blueness to yellowness on a scale of -128 to +127 of LAB color space.
- mean\_Prop\_R: mean value of redness of RGB color space.
- mean\_Prop\_G: mean value of greenness of RGB color space.
- mean\_Prop\_B: mean value of blueness of RGB color space.

Since the dataset is private information of the patients enrolled in Dr. Shah's clinic, we would not include the details of the data in this analysis.

## Glimpse of Data

```
## Rows: 72
## Columns: 15
## $ Image_URL      <chr> "https://storage.labelbox.com/cklyhytg2ulao0774uvyw0poy%~
## $ xmin           <dbl> 248, 340, 469, 644, 347, 511, 640, 744, 201, 322, 458, 6~
## $ xmax           <dbl> 292, 396, 527, 696, 385, 553, 682, 774, 248, 376, 509, 6~
## $ ymin           <dbl> 537, 465, 414, 426, 398, 413, 472, 532, 767, 796, 700, 4~
## $ ymax           <dbl> 582, 523, 474, 482, 461, 472, 522, 570, 803, 838, 726, 5~
## $ Mean_H         <dbl> 8.1922, 13.0864, 12.3233, 13.1862, 11.3486, 11.9142, 10.~
## $ Mean_S         <dbl> 0.2842, 0.2684, 0.2863, 0.2962, 0.2670, 0.2686, 0.2695, ~
## $ Mean_V         <dbl> 0.7824, 0.8003, 0.7863, 0.7895, 0.7871, 0.7820, 0.7820, ~
## $ Mean_L         <dbl> 68.7689, 68.3792, 68.0705, 68.5229, 68.4438, 68.2958, 67~
## $ Mean_A         <dbl> 15.4397, 13.9356, 15.3458, 14.7066, 14.6354, 14.0090, 14~
## $ Mean_B         <dbl> 10.6895, 12.0016, 11.8946, 14.3816, 11.9706, 11.9289, 11~
## $ Mean_Prop_R    <dbl> 0.3999, 0.3980, 0.3943, 0.3994, 0.3972, 0.3997, 0.3977, ~
## $ Mean_Prop_G    <dbl> 0.3077, 0.3125, 0.3147, 0.3140, 0.3113, 0.3115, 0.3107, ~
## $ Mean_Prop_B    <dbl> 0.2925, 0.2895, 0.2910, 0.2866, 0.2915, 0.2888, 0.2916, ~
## $ concentration <dbl> 9.5, 9.5, 9.5, 9.5, 9.5, 9.5, 9.5, 9.5, 9.2, 9.2, 9.2, 9~
```

## Cross Validation

In order to test how sensitive our estimated coefficients are to the training samples, we conducted a 10-fold cross validation.

## Summarize assessment CV metrics

| id     | .metric | .estimator | .estimate | .config              |
|--------|---------|------------|-----------|----------------------|
| Fold01 | rmse    | standard   | 1.444     | Preprocessor1_Model1 |
| Fold01 | rsq     | standard   | 0.004     | Preprocessor1_Model1 |
| Fold02 | rmse    | standard   | 1.118     | Preprocessor1_Model1 |
| Fold02 | rsq     | standard   | 0.485     | Preprocessor1_Model1 |
| Fold03 | rmse    | standard   | 1.475     | Preprocessor1_Model1 |
| Fold03 | rsq     | standard   | 0.084     | Preprocessor1_Model1 |
| Fold04 | rmse    | standard   | 1.230     | Preprocessor1_Model1 |
| Fold04 | rsq     | standard   | 0.129     | Preprocessor1_Model1 |

| id     | .metric | .estimator | .estimate | .config              |
|--------|---------|------------|-----------|----------------------|
| Fold05 | rmse    | standard   | 1.105     | Preprocessor1_Model1 |
| Fold05 | rsq     | standard   | 0.721     | Preprocessor1_Model1 |
| Fold06 | rmse    | standard   | 1.587     | Preprocessor1_Model1 |
| Fold06 | rsq     | standard   | 0.033     | Preprocessor1_Model1 |
| Fold07 | rmse    | standard   | 1.316     | Preprocessor1_Model1 |
| Fold07 | rsq     | standard   | 0.258     | Preprocessor1_Model1 |
| Fold08 | rmse    | standard   | 0.730     | Preprocessor1_Model1 |
| Fold08 | rsq     | standard   | 0.444     | Preprocessor1_Model1 |
| Fold09 | rmse    | standard   | 2.306     | Preprocessor1_Model1 |
| Fold09 | rsq     | standard   | 0.094     | Preprocessor1_Model1 |
| Fold10 | rmse    | standard   | 1.125     | Preprocessor1_Model1 |
| Fold10 | rsq     | standard   | 0.525     | Preprocessor1_Model1 |

## Summarize model fit CV metrics

| mean_adj_rsq | mean_aic | mean_bic |
|--------------|----------|----------|
| 0.364        | 162.362  | 183.081  |

| term        | estimate  | std.error | statistic | p.value |
|-------------|-----------|-----------|-----------|---------|
| (Intercept) | 1356.390  | 3815.071  | 0.356     | 0.724   |
| Mean_H      | 0.004     | 0.005     | 0.741     | 0.463   |
| Mean_S      | 31.719    | 18.789    | 1.688     | 0.098   |
| Mean_V      | 16.880    | 25.667    | 0.658     | 0.514   |
| Mean_L      | -0.126    | 0.245     | -0.512    | 0.611   |
| Mean_A      | -0.435    | 0.272     | -1.596    | 0.118   |
| Mean_B      | -0.257    | 0.328     | -0.783    | 0.438   |
| Mean_Prop_R | -1306.425 | 3818.067  | -0.342    | 0.734   |
| Mean_Prop_G | -1435.883 | 3818.314  | -0.376    | 0.709   |
| Mean_Prop_B | -1319.198 | 3813.765  | -0.346    | 0.731   |

| RMSE_train | RMSE_test |
|------------|-----------|
| 1.041681   | 1.445858  |

## Discussion of Model Output

We have derived the following linear model:

$$\hat{H}gb = 1356.390 + 0.004 \times \text{Mean}_H + 31.719 \times \text{Mean}_S \quad (1)$$

$$+ 16.880 \times \text{Mean}_V - 0.126 \times \text{Mean}_L - 0.435 \times \text{Mean}_A \quad (2)$$

$$- 0.257 \times \text{Mean}_B - 1306.425 \times \text{Mean\_Prop}_R \quad (3)$$

$$- 1435.883 \times \text{Mean\_Prop}_G - 1319.198 \times \text{Mean\_Prop}_B \quad (4)$$

We notice that while variables **Mean\_L**, **Mean\_A**, **Mean\_B**, **Mean\_Prop\_R**, **Mean\_Prop\_G**, and **Mean\_Prop\_B** are negatively associated with the response variable **Hgb** concentration, **Mean\_H**, **Mean\_S** and **Mean\_V** are positively

associated with the response variable. This means as mean value of hue, or mean value of saturation, or mean value of brightness of HSV color space increases, the hemoglobin concentration level also tends to increase. On the other hand, as the mean value of brightness, redness, greenness or blueness of RGB color space increases, the hemoglobin concentration level tends to decrease.

We found that out of all the predictors, **Mean\_S** tends to associates most significantly with the response variable with the smallest p value. For each 1 unit increase in the mean value of saturation, Hgb concentration is expected to increase by 31.719 on average, keeping all else constant. The second most significant predictor is **Mean\_A**, by which we found that for each 1 unit increase in the mean value of representation of greenness to redness, Hgb concentration is expected to decrease by 0.435 on average, keeping all else constant.

## 2.2 Principle Components Analysis (PCA)

### Data

In the first part of our research, we used the TBND\_V2 (Transient Biometrics Nails Dataset) on Kaggle

1

, which contains unlabeled nail bed images (Barbosa, Theoharis, and Abdallah 2019).

### Experiment on TBND\_V2 Dataset

We initially experimented with the TBND\_V2 dataset found on Kaggle, since there is not enough data from the patients. Given that the data is unlabelled, we tried unsupervised clustering methods on the data, hoping to gain some insights on classification of nail bed images.

To get features, we used VGG16, a convolutional neural network (CNN) in our model. It extracts features from the input images, turning each image into feature vectors (4096 by 1). We removed the final (prediction) layer from the neural network manually, and the new output layer is a fully-connected layer with 4,096 individual nodes. We do this by specifying the “outputs” argument when initialising the model. We therefore get input of our model by using the neural net VGG16 as a feature extractor for the image data.

We then performed a principal component analysis (PCA) on the feature vectors to reduce the dimension of the feature space. For each of the 93 image samples, we now have a corresponding 1 by 4096 feature vector. This means that our model needs to process a 93 by 4096 matrix. To reduce the computational and complexity cost of processing high-dimensional data, we performed principal component analysis (PCA) on the matrix for dimension reduction. We set the parameter to 50 to obtain the top 50 principal components of the feature vector. The principal components are by default sorted in descending order. This means that the first principal component will be able to explain the most variability in the feature vector. It’s a linear combination of the feature variables, and its direction captures the most variability. Thus, PCA helps us to reduce the dimension of the features from 4096 to 50 while preserving as much information in the original data as possible.

After getting features, we used Kmean clustering, an unsupervised algorithm that is commonly used in exploratory data analysis, to perform the clustering. The Kmean clustering works as follows: Initialise the centre of each k cluster by shuffling the dataset and randomly selecting K data points without replacement. Iterate the following steps until the assignment of data points to clusters is no longer changing: Compute the sum of the Euclidean distance squared between data points and all centres. Assign each data point to the closest cluster; each cluster is represented by its unique centre point. Compute the cluster’s centroid by taking the average of all data points assigned to that cluster.

In our model, we set the hyperparameter k to be 5 and clustered all samples into 5 categories. Based on the features we extracted, each cluster will contain images that are visually similar.

### PCA Results

Based on the documentation, the `explained_variance_ratio_` function returns the percentage of variability explained by each of the selected components. Running this function gives us the amount of variability that

is explained by all the PCs (0.1015 is explained by the first PC, 0.0894 by the second, and 0.0781 by the third etc.)

```
In [22]: pca.explained_variance_ratio_

Out[22]: array([0.10151978, 0.08946814, 0.07805712, 0.06449335, 0.05608589,
0.04782138, 0.04121738, 0.0316945 , 0.03162053, 0.02549183,
0.0250266 , 0.02351875, 0.0219821 , 0.01871412, 0.0177806 ,
0.01744947, 0.01555375, 0.01349313, 0.0124064 , 0.01192693,
0.01143915, 0.01080564, 0.01020705, 0.00966431, 0.00932489,
0.00897332, 0.0084416 , 0.00756158, 0.00751938, 0.00715409,
0.00660651, 0.00646016, 0.0060511 , 0.0058278 , 0.00565023,
0.00556989, 0.00544987, 0.00497867, 0.0046719 , 0.0045056 ,
0.00435844, 0.00414671, 0.0039429 , 0.00383598, 0.00379014,
0.00361596, 0.00351077, 0.0034277 , 0.00326044, 0.00311419,
0.00309905, 0.00298565, 0.0029144 , 0.00277587, 0.00269422,
0.00252586, 0.00243966, 0.00242474, 0.00231381, 0.00221256,
0.00215571, 0.00209832, 0.00196334, 0.00195042, 0.00183558,
0.00178786, 0.00176595, 0.00171314, 0.00166282, 0.00162487,
0.00158268, 0.00155565, 0.00142908, 0.00142209, 0.00140556,
0.00135851, 0.00133556, 0.0012662 , 0.00124887, 0.00123438,
0.00121307, 0.00119807, 0.00115303, 0.00109622, 0.00106344,
0.00101443, 0.0009791 , 0.00093018, 0.00088837, 0.00087974],
dtype=float32)
```

Figure 1: PCA variance ratios

Then we generated a bar chart to represent the variability explained by different principal components, as well as the cumulative step plot to represent the variability explained by the first most important components.

The neural net in the model functions as a feature extractor for the image data, which is the input of our model. To be more specific, we used VGG16, a convolutional neural network (CNN) in our model. It extracts features from the input images, turning each image into feature vectors (4096 by 1). We removed the final (prediction) layer from the neural network manually, and the new output layer is a fully-connected layer with 4,096 individual nodes. We do this by specifying the “outputs” argument when initialising the model.

The PCA reduces dimensions of our feature vectors. For each of the 93 image samples, we now have a corresponding 1 by 4096 feature vector. This means that our model needs to process a 93 by 4096 matrix. To reduce the computational and complexity cost of processing high-dimensional data, we performed principal component analysis (PCA) on the matrix for dimension reduction. We set the parameter to 50 to obtain the top 50 principal components of the feature vector. The principal components are by default sorted in descending order. This means that the first principal component will be able to explain the most variability in the feature vector. It’s a linear combination of the feature variables, and its direction captures the most variability. Thus, PCA helps us to reduce the dimension of the features from 4096 to 50 while preserving as much information in the original data as possible.

In our model, we set the hyperparameter k to be 5 and clustered all samples into 5 categories. Based on the features we extracted, each cluster will contain images that are visually similar.

### 3. Conclusion

In this study, we explored primarily two methods in predicting hemoglobin concentration level based on parameters derived from patients’ nail beds images: (1) multilinear regression, (2) principle components analysis. While the first model provides reasonable train and and test error, the p value of the predictors seems to be large. Our second model yields a reasonable explained variance ratio with respect to the scale of principle components. Limitations of our work includes that the data from the clinic lacks the label, therefore

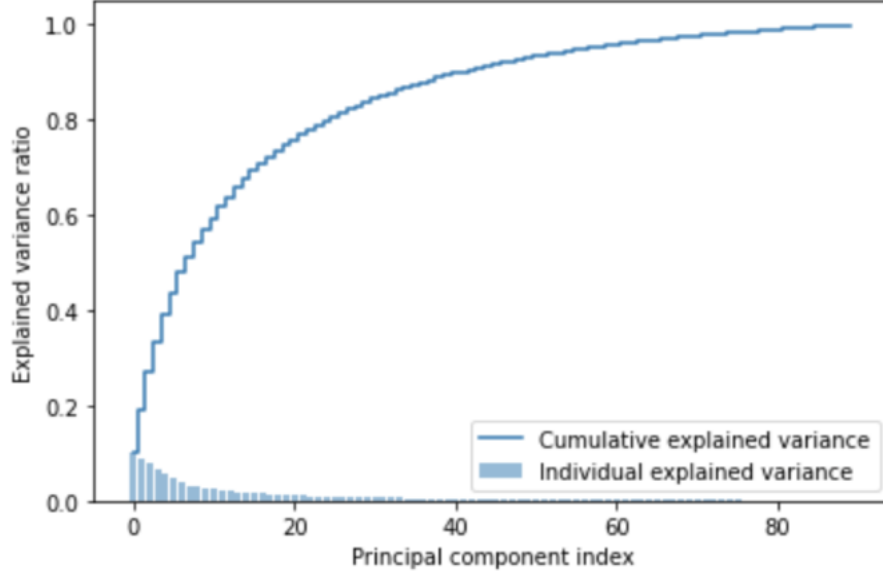


Figure 2: PCA cumulative step plot

we are unable to use classification methods such as random forest or SVM. Our future work includes fitting Bayesian regressions on the data, exploring more machine learning methods such as random forest or support vector machines, and using shrinkage methods such as ridge regression and LASSO.

## References

- Barbosa, Igor Barros, Theoharis Theoharis, and Ali E. Abdallah. 2019. “TBND\_V2.” Kaggle. <https://doi.org/10.34740/KAGGLE/DS/309682>.
- Mannino, Robert G. 2018. “A NONINVASIVE, IMAGE-BASED SMARTPHONE APP FOR DIAGNOSING ANEMIA,” 12.
- Safiri, Saeid. 2021. “Burden of Anemia and Its Underlying Causes in 204 Countries and Territories, 1990–2019: Results from the Global Burden of Disease Study 2019,” 1. <https://jhoonline.biomedcentral.com/articles/10.1186/s13045-021-01202-2>.