

Dissimilarity Based Regularized Deep Learning Model for Information Charts

Prerna Mishra
Department of Computer Science
DSPM-IIITNR
Raipur, CG, INDIA
Email: prerna@iiitnr.edu.in

Santosh Kumar
Department of Computer Science
DSPM-IIITNR
Raipur, CG, INDIA
Email: santosh@iiitnr.edu.in

Mithilesh Kumar Chaube
Department of Mathematics
DSPM-IIITNR
Raipur, CG, INDIA
Email: mithiles@iiitnr.edu.in

Abstract—The charts are very much convenient way to represent the complex data into simple pictorial based representation. Every chart type has variations in its characteristics, structure, and appearances making every type and subtype of chart different from each other on its physical outlook. Classification of such similar outlook charts still remains an untouched area. This paper presents a model that computes chart dissimilarity index, which is amalgamated with regularization on input layers of the learning model. Thus, all structural variations of charts are integrated into the model which produces 96.66% accuracy rate outperforming existing state-of-the-art models.

We proposed a novel approach to learn structure invariant dissimilar features using regularized learning techniques, by incorporating the dissimilarity index with the learning model to ease in learning dissimilar and hidden features of chart images.

Index Terms—Deep Learning, Chart Image Classification, Dissimilarity Index

I. INTRODUCTION

With the advancements in technology, data visualization has been continuously playing a crucial role in data analysis which is performed by machine-based experts and analysts. Charts represent, from delicate to complex information or data in various styles and orientations. Chart interpretation is as essential as textual interpretation to understand the information or knowledge portrayed within the document automatically [1][2]. Chart being the most influential presenter in infographics streams, it is the youngest field on research ground in terms of chart recognition, interpretation, visualization, and summarization. Different chart types require different visualization and interpretation mechanisms can be easily done by humans [3]. However, the same is complicated for the machine because it needs a significant number of constraints and parameters to be tuned [4].

On interpreting chart images automatically, significant issues lie in the variance within the structure, style, and visual appearance, making it very tedious for a machine to classify these charts using the same methodology. As per the available literature, every chart type requires a different methodology for classification and interpretation due to its diversified structure and orientation [5]. Every chart has its own structure and subtypes; thus, the same methodology fails to classify such

chart types.

A few works evaluated the classification of various pre-trained models in comparison with conventional classifiers over different types of charts [6], [7], [8], [9], [10] [5], [14], [15], [16]. Results showed that CNN models outperformed other classifiers giving highest accuracy. For raster formatted graphics authors employed a multi-class symbol shape classifier for graphical component classification generating multi-modal representations [16]. While an enhanced version of the pre-trained LeNet model was presented for classifying few types of chart classes [5]. In similar directions, a novel approach ReVision was suggested for beautifying each bitmap chart images automatically [5]. This system redesigns the chart image for better visual effects. The approach works in three stages, i.e., data extraction, chart classification by SVM using low-level image features and lastly, the redesign. System was tested on web collected chart images, giving satisfactory classification accuracy. A step ahead, a novel framework Deep Chart, was presented for classifying web collected chart image by using both convolution network and deep belief networks[14]. However, in literature few works on chart image classification [6], [7], [8], [9], [10], [14] depends only upon the primary and different chart types but not on its similar structural variants and subtypes; thus, significant features were avoided.

Several parameters, such as the orientation of charts, change of data, a variety of chart styles, and varied data types have been identified in the literature. Among all of these, chart style and orientation are some of the most challenging parameters in classification. The paper proposed a model to overcome such challenges faced during the classification of similar and different chart classes by incorporating the regularization on input layers of the CNN model. Regularization has effectively aided the process enabling the model to understand even the minute difference [9] within the chart and its subtypes, i.e., within classes and between classes difference in a better way.

A. Motivation and Research Contribution

Through the study of the literature, it was found that there are five deficiencies in the existing approaches. First, diverse chart types lead to difficulty in employing the same model

for all chart types, even though two chart images belong to the same chart class, their chart styles and structures can be significantly different [10] for example, stacked bar chart and vertical bar chart. Second, methods suggested until now are based on assumptions and consideration of constrained chart images using predefined architecture. Third, the methods suggested are focused only on local or low-level features, which are not useful in global chart interpretation. Fourth, there is no readily available database of charts, which consists of various charts and their subtypes. Fifth, similar interpretation might be generated for the charts even though they belong to different categories. Significant research contributions are given as follows:

- 1) A learning model utilizing chart dissimilarity based regularization to modify the cost function of CNN architecture improving the recognition accuracy of charts images having diversified structures and styles.
- 2) An image database of charts over 5280 images pertaining to 10 different chart type is handcrafted. Moreover, the performance of the current state of the art methods are evaluated across various types of chart images on a standard benchmark.

The proposed methodology is explained in Section II. Section III highlights the specifications of the dataset and discussion of results, followed by the conclusion.

II. PROPOSED MODEL

This section gives a brief overview of the architecture and methods used for chart image classification. In the model, a cost function is modified by adding a regularizer based on the dissimilarity index of the images on the input layers of the proposed CNN model, as seen in figure 1. The proposed model aims to classify multiple chart classes of varied and similar types and styles. To achieve this, a CNN model is regularised on input layers based on dissimilarity index of chart images, as shown in figure 2.

Let a be input nodes, h be the hidden nodes and c be the output nodes. Let W^l be the weight matrix between input layer l and layer $l+1$, i.e., it contains weights in l input layers and b_l be the bias of each layer l . Let Y_l denote the vector of input and Z_l be the output vector for each layer l . For training sets (u,v) such that $u_{a \times 1}$ and $v_{c \times 1}$ are of size a and c , respectively. The energy function of the CNN model can be given as follows:

$$\begin{aligned} Z &= \sum_1^n W^l \times u^l + b^l \\ Z &= f(W^l \times u^l + b^l) \end{aligned} \quad (1)$$

Where W^l is the weight for each input, u is input activations, $f(\cdot)$ is a non-linear function, and b is a biased term. A model tries to learn the function f such that the difference between predicted and actual output is minimal. The probability distribution p_i of a CNN can be stated as,

$$p_i = \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}} \quad (2)$$

Where a is an N vector and p_i is always positive in the range $(0,1)$. Using the training data, CNN is trained to maximize the likelihood, i.e., the cost function indicates the difference of actual output distribution and acquired output distribution is given as,

$$C_{z,p} = - \sum z_i \log(y_d) \quad (3)$$

$$y_d = \zeta(c) = \frac{e^{c_d}}{\sum_{i=1}^D e^{c_i}} \text{ for } d = 1, \dots, D \quad (4)$$

Where function ζ takes a D -dimensional vector c as input and produces a vector y of probability within 0 and 1 as output.

$$C_{z,p} = - \log \left(\frac{e^{c_d}}{\sum_{i=1}^D e^{c_i}} \right) \text{ By Eqn (4)}$$

Where y_d is the softmax activation, which is combined with a categorical cross-entropy loss function. Equation 3 is optimized using root mean square propagation[11]. The optimizer updates the parameters recursively as,

$$E[G^2]_s = \beta E[G^2]_{s-1} + (1 - \beta) \left(\frac{\delta C_{z,p}}{\delta W} \right)^2 \quad (5)$$

Where $E[G^2]$ is the exponential average of square gradients, β is an average parameter, $\left(\frac{\delta C}{\delta W} \right)$ is the gradient of cost function C concerning weight W and η is the initial learning rate. Weights are further updated to avoid the drastic weight changes while learning,

$$W_s = W_{s-1} + \frac{\eta}{E[G^2]_s} \left(\frac{\delta C_{z,p}}{\delta W} \right)^2 \quad (6)$$

Regularizer term is slightly modified by adding a term \hat{I}_θ , which is a dissimilarity index of each chart images where $\hat{I}_\theta \in [0,1]$. Addition of dissimilarity index \hat{I}_θ with regularizer enables learning underlying dissimilarities within chart images belonging to different categories. It is used to avoid overfitting, which was caused by similar chart subtypes. It improves the learning accuracy of the hidden representation of different similar-dissimilar chart images, thus forcing the model to learn as per these variations. The cost function(Eqn 3) is modified by adding regularization terms λ to it. The updated cost functions of model,

$$\tilde{C}_{z,p} = C_{z,p} + \lambda \left\| \hat{I}_\theta \tilde{W} \right\| \quad (7)$$

where, \tilde{W} is weight matrix learned to minimize the loss, λ is the regularization parameter learned while training.

Computing Dissimilarity Index:

Similarity and dissimilarity within chart images might be more or less depending upon chart characteristics and features. The effect of these variations on chart images is not always consistent. Due to these inconsistencies, it was difficult to model the chart variations over similar structures. In order to extract the dissimilarity within targeted chart images, features from these charts are extracted. Associated features are essential for thorough analysis as they can expose significant relationships amongst them. Hence, correlation analysis[12] is incorporated

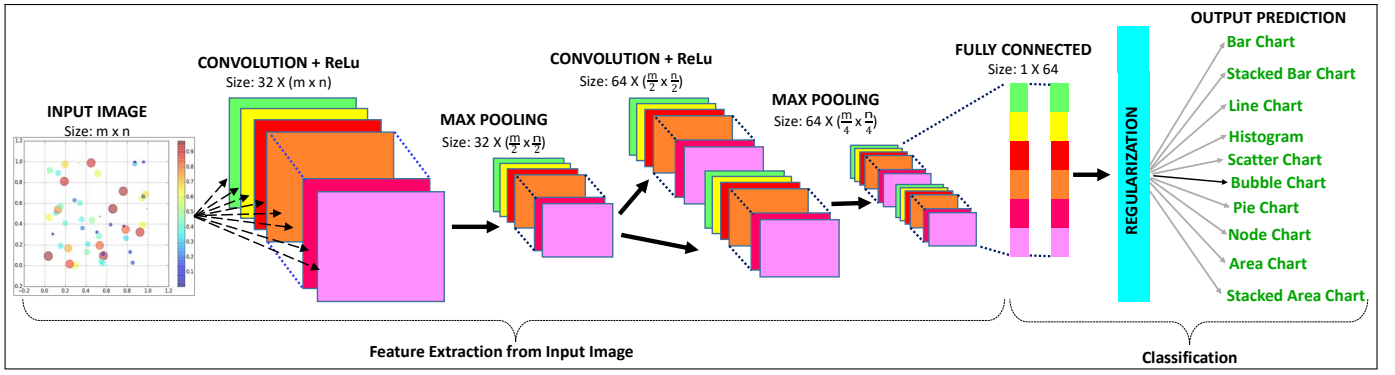


Fig. 1: Proposed model for chart classification

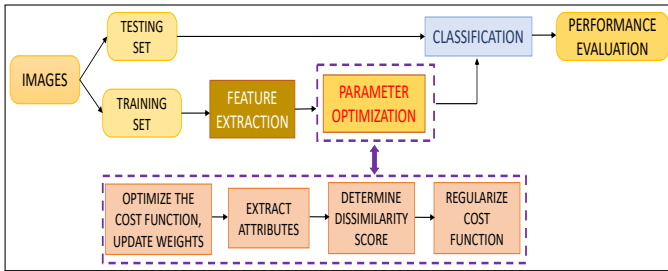


Fig. 2: Flowchart of the proposed method

to determine the set of associated features. The association affecting at a global level is determined using Pearsons correlation matrix[13]. Pearsons correlation coefficient R amid image under consideration T_1, T_2 and its mean \bar{T}_1, \bar{T}_2 respectively is given as,

$$R = \frac{\sum_{i=1}^n (T_{i1} - \bar{T}_1) (T_{i2} - \bar{T}_2)}{\sqrt{\sum_{i=1}^n (T_{i1} - \bar{T}_1)^2 (T_{i2} - \bar{T}_2)^2}} \quad (8)$$

$$R = \begin{cases} 1, & \text{if images are not associated} \\ -1, & \text{otherwise} \end{cases}$$

Every image has a region or a part of a region(sub-region) which has a dissimilar set of features while there are sub-regions that have a similar type of features. In other words, there are regions which have all the features as dissimilar or similar type in comparison to other sub-regions. For every such sub-region of interest ϑ_i , is computed separately depending upon the number of dissimilar regions found. Dissimilarity score for images, \hat{I}_ϑ could be given as

$$\hat{I}_\vartheta (T_1, T_2) = \sqrt{\sum_0^{N-1} d_i^1 - d_i^2} \quad (9)$$

where $d_i^2 = (t_i - \min(T_i)) / (\max(T_i) - \min(T_i))$

The global dissimilarity index matches whole images in a way that, variations in any of the images will affect the output dissimilarity. Given dissimilarity index acquires the set of different features that are useful in recognition. So accordingly

the optimizer will be,

$$E[G^2]_s = \beta E[G^2]_{s-1} + (1 - \beta) \left(\frac{\delta C_{z,p}}{\delta W} \right)^2 \quad (10)$$

Weight W is now given as,

$$W_s = W_{s-1} + \frac{\eta}{E[G^2]_s} \left(\frac{\delta C_{z,p}}{\delta W} \right)^2 \quad (11)$$

The regularization term is fused with a global dissimilarity index computed over each chart in order to find the differences amongst the target chart images. L_1, L_2 and L_{enet} regularizer are given as for input layer Y^l

- 1) L_1 (Lasso) regularization, $R(W_t) = \lambda_1 \|W_{t,j}^y\|$ is now stated as $R(W_t) = \lambda_1 \|\hat{I}_\vartheta W_{t,j}^y\|$
- 2) L_2 (Ridge) regularization, $R(W_t) = \lambda_2 \|W_{t,j}^y\|^2$ is now stated as $R(W_t) = \lambda_2 \|\hat{I}_\vartheta W_{t,j}^y\|^2$
- 3) L_{enet} (Elastic Net) regularization, $R(W_t) = (\lambda_1 \|W_{t,j}^y\| + \lambda_2 \|W_{t,j}^y\|^2)$ is now stated as $R(W_t) = \lambda_1 \|\hat{I}_\vartheta W_{t,j}^y\| + \lambda_2 \|\hat{I}_\vartheta W_{t,j}^y\|^2$

L_2 regularization is the summation of square of all weights, while compelling(reducing) the weights to be smallest (towards 0) but higher than 0. L_1 regularization is the absolute value of the weights, compelling parameters to be reduced to zero. L_{enet} seeks to combine both L_1 and L_2 regularization. It merges the penalties of both L_1 and L_2 norm, providing the best of both norms. The final cost function is now given as,

$$C_{z,p}^{L_1} = C_{z,p} + \lambda_1 \|\hat{I}_\vartheta W_{t,j}^y\| \quad (12)$$

$$C_{z,p}^{L_2} = C_{z,p} + \lambda_2 \|\hat{I}_\vartheta W_{t,j}^y\|^2 \quad (13)$$

$$C_{z,p}^{L_{enet}} = C_{z,p} + \left(\lambda_1 \|\hat{I}_\vartheta W_{t,j}^y\| + \lambda_2 \|\hat{I}_\vartheta W_{t,j}^y\|^2 \right) \quad (14)$$

III. EXPERIMENTAL RESULTS AND ANALYSIS

In this section, the proposed model is evaluated on the crafted dataset for chart image classification.

A. Dataset Preparation and Description

Due to the unavailability of any public dataset, multi-class chart images were handcrafted to quantitatively assess the proposed model. Each chart image has x-axis values, y-axis values, labels ticks, legends, chart titles, axis marks, and it was tried to retain all types of chart properties in the images. The dataset consists of a total of 5280 images out of which 3780 images are as the training set covering ten classes, namely bar chart, stacked bar chart, line chart, pie chart, scatter chart, histogram, area chart, nodal chart, bubble chart, stacked area chart. Apart from these training datasets, 1500 images were used to test the proposed model, which had near about 150 images belonging to each class. Samples of the learning dataset are shown in figure 3.

CHARTS	IMAGES	CHARTS	IMAGES
BAR CHART	300	LINE CHART	500
NODAL CHART	500	AREA CHART	480
SCATTER CHART	500	PIE CHART	300
HISTOGRAM	300	STACKED AREA CHART	300
STACKED BAR CHART	300	BUBBLE CHART	300

Fig. 3: Dataset classes

B. Result Analysis

Our handcrafted chart images were used to evaluate our model. Author [7],[14] suggested a model that effectively performed the work on chart images, so in order to compare the classification results with [7],[14], 1500 images were used which was drawn from 10 categories as stated in section III.A. The proposed model was tested on a handcrafted labeled dataset to evaluate its performance on a similar type or subtype of chart images. The experimental results of the proposed model in comparison to other benchmark systems are summarized in Table I, II and III. Hyperparameter choices for obtained results were- learning rate: 0.001, optimization algorithm: RMSProp, $\alpha=0.9$, Batch size: 25, Activation function: ReLu, $\lambda_1=0.0001$, $\lambda_2=0.0003$.

The reason behind variation in results was due to the similarity between chart types and subtypes. Table I shows the accuracy rate obtained for each chart class of our testing set, under different regularization norms. From table I, it is seen that model was able to predict scatter chart and bubble chart, but due to its similar structural property, the results generated are in the range of 92-97%, giving the highest accuracy rate of 96.66% and 97.33% respectively. On the other hand, the model was able to classify between histogram, stacked bar charts and bar charts appropriately, yielding the highest accuracy of 97.33%, 96.66%, and 97.80%, respectively. Dropout regularization yielded better results for the chart image dataset in comparison to other norms. Hence it can be stated that ridge, lasso, and elastic net with dropout norm was able to generate better results in comparison to a model without dropout norm under the same parameters

and criterias. Confusion matrix of bar chart for Elastic net+ Dropout is obtained as $\begin{bmatrix} TP = 125 & FP = 8 \\ FN = 25 & TN = 1342 \end{bmatrix}$

The model learns competently on heterogeneous data in terms of color intensity and structure orientations, as seen in figure 4. It was seen that the scatter chart, line chart, and bubble chart share a similar type of style, data representation, and orientation.

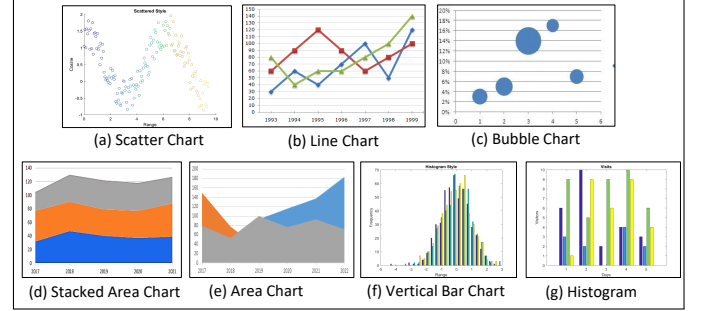


Fig. 4: Illustration of confusing chart images (a-b-c), (d-e), (g-f)

As per the accuracy level summarised in table II the model achieves the average accuracy (of all the categories) rate of 95.36% without regularization and 96.66% with regularization stating that proposed model with regularization has given slightly improved results in comparison to other methods as shown in table II. Author [7] classified images using SVM on low-level image features. Author [14] used a deep learning model for image classification and deep belief network for further reducing the size of the feature dimension. Table I provides detailed performance results over all the chart types.

As the proposed model targeted similar structured images, it was necessary to determine the true positive and false positive rate of chart images. Recall gave the correct positive prediction rate, while precision denotes the accurateness of the model to predict actual positive results over positive predicted results. The results of regularization show that the model has performed better in terms of prediction of similarly seemed chart images, thus giving better results in comparison to other methods.

Most effectually, the proposed model could very well discriminate between similarly seemed charts. Thus, regularization functions have benefitted the model to overcome the said challenges [5],[10] in a very efficient manner. Apart from this, the model was able to generate a high accuracy rate in recognizing various chart types.

Figure 5 concludes that the model with a gradual decrease/increase of loss/accuracy rate has given a better performance rate. While drastic fall/rise of loss/accuracy rate tends to overfit in training giving false results on testing. With dropout, there is a gradual increase in the training rate. While the loss rate with dropout is progressively decreasing with epochs.

Table III shows the accuracy rate of chart types with other benchmark systems [7],[14]. From the table, it can be seen that the proposed model has given good results for line charts

TABLE I: Accuracy Rate (%) of Chart Types from Our Dataset with Respect to Regularization Function

Chart Type	L_1 +DI	L_2 +DI	L_{enet} +DI	Dropout	L_1 +DI+Dropout	L_2 +DI+Dropout	L_{enet} +DI+Dropout
Bar Chart	94	93.33	94.66	96.66	95.33	94	97.80
Stacked Bar Chart	94.66	92	94	95.33	96.66	94.66	96.66
Line Chart	96.66	98.33	96.66	97.33	97.33	96	98.66
Pie Chart	95.33	90.66	96.66	98	97.33	93.33	98.33
Scatter Chart	93.33	90	94.66	96.66	96.66	92.66	96.66
Histogram	93.33	91.33	96	96.66	97.33	94	97.33
Area Chart	94	92	94	97.33	96	96	98
Nodal Chart	94.44	93.33	96.66	98.66	97.33	96	98.66
Bubble Chart	92	89.33	96	96	95.33	92	97.33
Stacked Area Chart	95.33	90	97.33	96	97.33	92.66	98.66

Abbreviations: DI-Dissimilarity Index

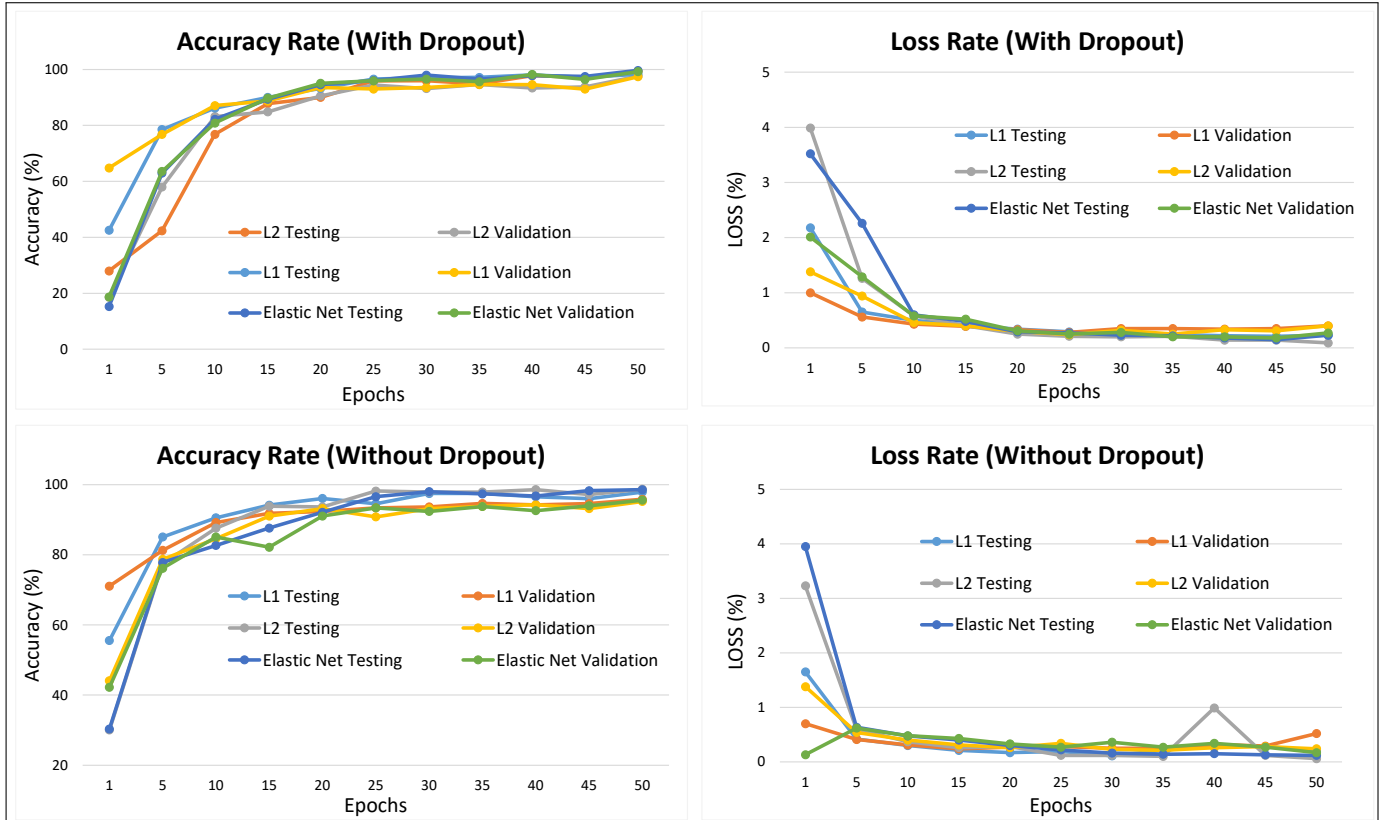


Fig. 5: Performance rate of the model (Training Vs. Validation)

TABLE II: Comparison with Existing Models

Model	Accuracy(%)
Tang Binbin, et al. [14]	72.6
Manolis Savva et al. [7]	90
Our model without DI	95.36
Our model with DI	96.66

giving maximum accuracy rate. The model was able to classify amongst histogram and bar chart, which has similar-looking structures. Though for a pie chart, the accuracy rate was 89.47% which is the lowest in comparison with other chart types, the reason for this might be the somewhat similar

structure features, i.e., different subsection having a different color intensity which is also the case of other chart types too for example bar chart and area charts.

From tables I and III, it can be noted that the model had classified various chart types correctly irrespective of their similar structures and characteristics. Norms with dropout have yielded better results than norms without dropout with an increase of near about 2% accuracy rate. For all chart types, minimum accuracy level was obtained by lasso norm in comparison to other norms, and the highest accuracy range was obtained by applying L_{enet} + dropout with the proposed model. The major obstacle of recognizing chart types [5],[10] and its variant was overcome by using chart dissimilarity based

TABLE III: Performance Rate Concerning Particular Chart Type

Chart Type	Proposed Model	Monolis Savva et al.[7]	Tang Binbin et al. [14]
Bar Chart	95.12	78	74.2
Line Chart	100	-	67.9
Pie Chart	89.47	79	59.4
Scatter Chart	91.63	79	84.2
Area Chart	96.13	88	-
Histogram	96.47	-	-
Nodal Bubble	98.24	-	-
Bubble Chart	96	-	-
Stacked Area Chart	96.28	-	-
Stacked Bar Chart	96	-	-

- : model does not consider particular chart type

regularization on the CNN model.



Fig. 6: Samples of the obtained results using the proposed model

IV. CONCLUSION

The paper proposed a chart dissimilarity based regularized model for learning hidden features and specifics which learns definite difference amongst a variety of chart images adjusting regularization method to tackle the challenges of chart image recognition due to its variations in the structure, orientation, and graphical components. Results depict that the proposed model has enabled adaptive learning of these variations of types and its subtypes, having a different structure, characteristics, and appearance.

REFERENCES

- [1] Poco, Jorge, and Jeffrey Heer, "Reverse- engineering visualizations: Recovering visual encodings from chart images.", in *Proceedings of Computer Graphics Forum*, vol. 36, no. 3, pp. 353-363. 2017.
- [2] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton, "Imagenet classification with deep convolutional neural networks", in *Advances in Neural Information Processing Systems*, pp. 1097-1105. 2012.
- [3] Muzammil Khan and Sarwar Shah Khan, "Data and Information Visualization Methods, and Interactive Mechanisms: A Survey", in *International Journal of Computer Applications*, 34(1):1-14, November 2011.
- [4] Tam, Gary KL, Vivek Kothari, and Min Chen, "An analysis of machine- and human-analytics in classification", in *IEEE Transactions on Visualization and Computer Graphics*, 23.1 (2016): 71-80.
- [5] Amara, Jihen, Pawandeep Kaur, Michael Owonibi, and Bassem Bouaziz, "Convolutional neural network based chart image classification.", 2017.
- [6] Huang, Weihua, Siqi Zong, and Chew Lim Tan, "Chart image classification using multiple-instance learning.", in *IEEE Workshop on Applications of Computer Vision (WACV'07)*, pp. 27-27. IEEE, 2007.
- [7] Savva, Manolis, Nicholas Kong, Arti Chhajta, Li Fei-Fei, Maneesh Agrawala, and Jeffrey Heer, "Revision: Automated classification, analysis and redesign of chart images.", in *Proceedings of the 24th annual ACM Symposium on User Interface Software and Technology*, pp. 393-402. 2011.

- [8] Chagas, Paulo, Rafael Akiyama, Aruanda Meiguins, Carlos Santos, Filipe Saraiva, Bianchi Meiguins, and Jefferson Morais. "Evaluation of convolutional neural network architectures for chart image classification." In *International Joint Conference on Neural Networks (IJCNN)*, pp. 1-8. IEEE, 2018.
- [9] Rodriguez, Pau, Jordi Gonzalez, Guillem Cucurull, Josep M. Gonfaus, and Xavier Roca, "Regularizing cnns with locally constrained decorrelations", arXiv preprint arXiv:1611.01967, 2016.
- [10] Jung, Daekyoung, Wonjae Kim, Hyunjo Song, Jeong-in Hwang, Bongshin Lee, Bohyoung Kim, and Jinwook Seo. "ChartSense: Interactive data extraction from chart images." in *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 6706-6717. 2017.
- [11] Geoffrey Hinton, "Neural Networks for machine learning", *Online course*. Published online: <https://www.coursera.org/learn/neural-networks/home/welcome>
- [12] Haghighat, Mohammad, Mohamed Abdel-Mottaleb, and Wadee Alhalabi. "Discriminant correlation analysis: Real-time feature level fusion for multimodal biometric recognition", in *IEEE Transactions on Information Forensics and Security*, 11.9 (2016): 1984-1996.
- [13] Vijay Kotu, Bala Deshpande, Classification, *Data Science (Second Edition)*, 2019.
- [14] Tang, Binbin, Xiao Liu, Jie Lei, Mingli Song, Dapeng Tao, Shuifa Sun, and Fangmin Dong. "Deepchart: Combining deep convolutional networks and deep belief networks in chart classification." in *Signal Processing*, 124 (2016): 156-161.
- [15] Chagas, Paulo, Rafael Akiyama, Aruanda Meiguins, Carlos Santos, Filipe Saraiva, Bianchi Meiguins, and Jefferson Morais, "Evaluation of Convolutional Neural Network Architectures for Chart Image Classification", In 2018 International Joint Conference on Neural Networks (IJCNN), pp. 1-8. IEEE, 2018.
- [16] Gao, Jinglun, Rafael E. Carrillo, and Kenneth E. Barner., "Image categorization for improving accessibility to information graphics", in *Proceedings of the 12th International ACM SIGACCESS Conference on Computers and Accessibility*, pp. 265-266. 2010.