Interactive Deep Learning for Sorting Plant Images by Visual Phenotypes

Huimin Han^{1,1}, Ritvik Prabhu^{1,1}, Timothy Smith^{2,2}, Kshitiz Dhakal^{1,1}, Xing Wei^{1,1}, Song Li^{1,1}, and Chris North^{1,1}

 $^{1}\mathrm{Virginia}$ Tech $^{2}\mathrm{Virginia}$ Techh

November 30, 2022

Abstract

This paper proposes an interactive system called Andromeda that enables users to interact with machine learning models by sorting images in a reduced dimension plot. In our system, a dimension reduction algorithm projects the images into a 2D space representing similarities between the images based on visual features extracted by a deep neural network. With Andromeda, users can alter the projection by dragging a subset of the images into groups according to their domain expertise. The underlying machine learning model learns the new projection by optimizing a weighted distance function in the feature space, and the model re-projects the images accordingly. The users can explore multiple custom projections, and can export a model for future classification tasks. Our approach incorporates user preferences into machine learning model construction and allows reuse of pre-trained image processing models to accomplish new tasks based on user inputs. Using edamame pod images as an example, we transferred a maturity based model into a model that can classify number of seeds per pod to demonstrate the utility of our system.

Interactive Deep Learning for Exploratory Sorting of Plant Images by Visual Phenotypes

Huimin Han^a, Ritvik Prabhu^a, Timothy Smith^a, Kshitiz Dhakal^b, Xing Wei^b, Song Li^b, and Chris North^a

^aDepartment of Computer Science, Virginia Tech ^bSchool of Plant and Environmental Sciences, Virginia Tech

ABSTRACT

This paper proposes an interactive system called Andromeda¹ that enables users to interact with machine learning models to allow for exploratory sorting of images through a cognitive approach that uses a reduced dimension plot. In our system, a dimension reduction algorithm projects the images into a 2D space representing similarities between the images based on visual features extracted by a deep neural network. With Andromeda, users can alter the projection by dragging a subset of the images into groups according to their domain expertise. The underlying machine learning model learns the new projection by optimizing a weighted distance function in the feature space, and the model re-projects the images accordingly. The users can explore multiple custom projections to learn about the visual support for different groupings based on explainable-AI feedback. Our approach incorporates user preferences into machine learning model construction and allows transfer learning from pre-trained image processing models to accomplish new tasks based on user inputs. Using edamame pod images as an example, we interactively re-project the images into different groupings based on maturity and disease, and identify important visual features from the pixels highlighted by the model.

Keywords: Deep learning feature visualization, dimension reduction, interactive visual analytics

1. INTRODUCTION

In recent years, computer vision and artificial intelligence (AI) have played crucial roles in automating imagebased decision processes in agriculture research and production. Many AI models have been developed to diagnose plant diseases,² determine plant species,^{3,4} and assess plant product quality.⁵ Most published models follow a simple workflow consisting of a label-train-test-release cycle.

However, researchers sometimes need a more exploratory approach in which they interactively explore many alternative labelings based on their domain expertise. Thus, the goal of this paper is to create an algorithm that aids users in performing exploratory sorting on the fly, as they might sort physical seed pods on a tabletop, augmented with interactive machine learning that incorporates user perception feedback.

To dynamically incorporate user's exploratory feedback into the machine learning processes, we developed an interactive machine learning platform within a computational notebook (Jupyter) called Andromeda. Users can provide feedback to the machine learning models through regrouping the projected images. First, we exploit transfer learning⁶ to extract image features using a popular pre-trained image-processing neural network. Second, the user is provided with a 2-D interactive projection where they can visualize similarities and re-group the images. Third, a machine-learning algorithm for inverse projection learns the revised similarities between the user-specified groups and renders a new projection of the images. It uses the interactive feedback from the user to learn a novel distance model that weights the visual features extracted by the network. Fourth, the explainable-AI highlights the pixels in the images corresponding to the up-weighted image features, thus providing visual justification for the user's grouping. In this manuscript, using images of edamame pods as our model system, we demonstrate the functionality and utility of our interactive machine learning approach.

Further author information: (Send correspondence to Chris North and Song Li)

Chris North: E-mail: north@vt.edu, Telephone: +1 540 231 2458

Song Li: E-mail: songli@vt.edu, Telephone: +1 540 231 2756

2. MATERIALS AND METHODS

2.1 Dataset and Preprocessing

Images used in this paper were collected by the Li Lab of Applied Machine Learning in Genomics and Phenomics at Virginia Tech.⁷ This dataset comprises ready-to-harvest, late-to-harvest, and diseased pod images (100 images with 10-20 pods in each image). Figure 1 shows the sample raw data and image pre-processing results. We used an improved vegetation index, Excess Green minus Excess Red (ExG- ExR),⁸ to identify pods for our data sets. ExR was subtracted from ExG with a zero threshold to create the ExG-ExR binary image. After computing a binary image from vegetation indices, we applied several morphological transformations.^{9,10} We used dilation to increase the object area and closing and opening, which cleaned background noise by imputing missing pixel values. Finally, after vegetation indices and morphological transformations, we obtained a binary image mask with pods as white and background as black. Pods were detected by finding the contours of these masks.



Figure 1. Sample raw data and preprocessing results for a diseased pod

2.2 Feature Extraction and Visual Back Propagation

We use a convolutional neural network (CNN) to convert images into meaningful quantitative representations.



Figure 2. Resnet-18 Image Feature Extractor and Visual Back Propagation for Explainable-AI

- CNN Representations: Convolutional neural networks are a powerful tool to learn representations of image data with multiple levels of abstraction. In particular, convolutional neural network (CNN) models are widely used in computer vision.^{11, 12} Although any pre-trained CNN model could be applied, in this paper we use the pre-trained ResNet-18 model from ImageNet.¹³ ResNet was introduced in 2015 and won several competitions in computer vision¹⁴, and is one of the widely used CNN models to extract features from images. We use the last convolutional layer to extract 512 features from each image.¹⁵
- Visual Back Propagation To visualize each feature extracted from the ResNet-18 model, we utilize a modified version of the visual back propagation¹⁶ method to visualize sets of pixels of the input image that contribute most to each feature. Starting from the 512 feature map from the last convolutional layer, we back propagate each feature and average the feature maps after each ReLU layer. The averaged feature

map of the last convolutional layer is scaled-up via deconvolution and multiplied by the averaged feature map from the previous layer. The resulting intermediate mask is again scaled-up and multiplied. This process is repeated until we reach the input image layer. We initiate visual back propagation using feature weights from the interactive dimension reduction model, thus enabling explainability of the projection.

2.3 Dimension Reduction and Visual Analytics

To visualize similarities between images, we use a dimension reduction (DR) algorithm to project the 512dimensional data into a 2-dimensional plot. To allow users to drag images and form new projections, we use an interactive framework called Andromeda.¹ A weighted Multi-Dimensional Scaling (MDS) algorithm with a weighted distance metric enables both forward and inverse projection. Although any dimension reduction algorithm could be applied (such as PCA or t-sne), MDS closely matches the user's cognitive task of sorting images by mapping high-dimensional image similarities to 2D distances. Proximal images in the projection are similar in the weighted high-dimensional image feature space ¹. MDS also easily adapts to different distance functions for experimentation. While MDS projects the high-dimensional data to a 2D scatter plot,¹⁷ a weighted distance function with user-specified weights on each dimension enables alternative projections that emphasize different dimensions. After interactive sorting, an inverse dimension reduction algorithm learns distance function weights for the user-modified projections of the images. Figure 3 shows the system design.



Figure 3. Andromeda System Design

3. INTERACTIVE CASE STUDY

These case studies validate the hypothesis that our interactive system with features extracted by the deep neural network can dynamically capture the novel abstractions interactively specified by the users when dragging images in the plot, and that the visual back propagation method provides explainability to uncover the important visual features identified through the interactive learning process that support the user's abstractions.

3.1 Pods Based on Maturity Stage

The maturity stage of each pod as either diseased, late-to-harvest, or ready-to-harvest is a phenotype that can be determined by trained observers. Here we test if a user can sort the images according to these phenotypes with the help of the machine learning. The image data for 30 randomly chosen edamame pods are displayed on the 2D projection as shown in figure 4 (a); note that the color coding was added to show that the default plot does not capture the desired phenotypes as clear clusters. The initial weights for each feature are equal. In figure 4 (b) the user interactively drags 15 pods highlighted in green in order to group them into 3 clusters according to desired phenotype categories. Figure 4 (c) shows the updated projection, which produced three main clusters of pods according to their maturity stage. The red cluster shows the pods that are too late to harvest, the blue cluster are diseased and the green cluster are ready to harvest. This indicates that the desired phenotypes of each pod were effectively captured by the weighted features and successfully learned by Andromeda.

Furthermore, the explainable-AI visualizations of specific pods depict the most important visual features in our re-grouping as learned by the interactive model. In figure 4 (d) we see that one of the more important visual features learned to determine disease phenotype is a salient discolored spot. Similarly, in figure 4 (e,f), areas of each pod correlating to important features are highlighted. This provides us with insight into which parts of the pod are important for visually discerning a diseased, late-to-harvest, or ready-to-harvest product.



3.2 Pods Based on Seed Numbers

The number of seeds per pod is an important phenotype that potentially affects consumer acceptance of the product. However, the images were not originally collected to determine the number of seeds. The number of seeds is a novel visual feature that can be observed directly by the end users but is not initially clustered in the default projection. Images of 30 edamame pods are displayed in the 2D plot in Figure 5 (a) with equal weights applied for each feature. The user interactively drags 15 pods highlighted in green to group them into 3 clusters according to number of seeds, as shown in 5 (b). Figure 5 (c) shows the updated projection with clear clustering. We find that the "number of seeds" phenotype is well captured by the weighted features learned by Andromeda.



To explain the feature space with updated weights, we select the features with higher weights as an example. In Figure 5 (d,e,f), the most relevant CNN features mainly captures the overall shape of the pod to differentiate pods with different numbers of seeds.

These case study results indicate that Andromeda with features extracted by CNN deep neural networks can indeed enable interactive sorting of pod images according to various human-guided visual phenotypes. In future work, the resulting weighted feature models could then be used as classifiers for larger collections of images. We plan to extend these methods to more complex input scenarios, such as images of pods on live plants captured in the field with mobile phones.

REFERENCES

- Zeitz, J., Dowling, M., Wenskovitch, J., Crandell, I., Wang, M., House, L., Leman, S., and North, C., "Observation-level and parametric interaction for high-dimensional data analysis," ACM Trans. Interact. Intell. Syst. 8 (June 2018).
- [2] Mohanty, S. P., Hughes, D. P., and Salathé, M., "Using deep learning for image-based plant disease detection," Front. Plant Sci. 7, 1419 (Sept. 2016).
- [3] Barré, P., Stöver, B. C., Müller, K. F., and Steinhage, V., "LeafNet: A computer vision system for automatic plant species identification," *Ecol. Inform.* 40, 50–56 (July 2017).
- [4] Tan, J. W., Chang, S.-W., Abdul-Kareem, S., Yap, H. J., and Yong, K.-T., "Deep learning for plant species classification using leaf vein morphometric," *IEEE/ACM Trans. Comput. Biol. Bioinform.* 17, 82–90 (Jan. 2020).
- [5] Yu, D., Lord, N., Polk, J., Dhakal, K., Li, S., Yin, Y., Duncan, S. E., Wang, H., Zhang, B., and Huang, H., "Physical and chemical properties of edamame during bean development and application of spectroscopybased machine learning methods to predict optimal harvest time," Food Chem. 368, 130799 (Jan. 2022).
- [6] Zhuang, F., Qi, Z., Duan, K., Xi, D., Zhu, Y., Zhu, H., Xiong, H., and He, Q., "A comprehensive survey on transfer learning," *Proceedings of the IEEE* 109(1), 43–76 (2021).
- [7] Li, S., "Li lab of applied machine learning in genomics and phenomics." https://lilabatvt.github.io/.
- [8] Meyer, G. and Camargo Neto, J., "Verification of color vegetation indices for automated crop imaging applications," *Computers and Electronics in Agriculture* 63, 282–293 (10 2008).
- [9] Raid, A., Khedr, W., El-dosuky, M., and Aoud, M., "Image restoration based on morphological operations," International Journal of Computer Science, Engineering and Information Technology 4, 9–21 (07 2014).
- [10] Gil, J. and Kimmel, R., "Efficient dilation, erosion, opening and closing algorithms," in [ISMM], (2000).
- [11] Denker, J. S., Gardner, W. R., Graf, H. P., Henderson, D., Howard, R. E., Hubbard, W., Jackel, L. D., Baird, H. S., and Guyon, I., "Neural network recognizer for hand-written zip code digits," in [Proceedings of the 1st International Conference on Neural Information Processing Systems], NIPS'88, 323–331, MIT Press, Cambridge, MA, USA (1988).
- [12] LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D., "Backpropagation applied to handwritten zip code recognition," *Neural Computation* 1(4), 541–551 (1989).
- [13] He, K., Zhang, X., Ren, S., and Sun, J., "Deep residual learning for image recognition," CoRR abs/1512.03385 (2015).
- [14] Schaetti, N., "Character-based convolutional neural network and resnet18 for twitter author profiling : Notebook for pan at clef 2018," (09 2018).
- [15] Bian, Y., Wenskovitch, J., and North, C., "Deepva: Bridging cognition and computation through semantic interaction and deep learning," in [Proceedings of the IEEE VIS Workshop MLUI 2019: Machine Learning from User Interactions for Visualization and Analytics], (2019).
- [16] Bojarski, M., Choromanska, A., Choromanski, K., Firner, B., Ackel, L. J., Muller, U., Yeres, P., and Zieba, K., "Visualbackprop: Efficient visualization of cnns for autonomous driving," in [2018 IEEE International Conference on Robotics and Automation (ICRA)], 4701–4708 (2018).
- [17] Cox, M. A. A. and Cox, T. F., [Multidimensional Scaling], 315–347, Springer Berlin Heidelberg, Berlin, Heidelberg (2008).

Response to the Reviews and Decision

Title: Interactive Deep Learning for Investigative Sorting of Plant Images by Visual Phenotypes

DOI: 10.1002/essoar.10508768.1

> Authors: Huimin Han Ritvik Prabhu Timothy Smith Kshitiz Dhakal Xing Wei Song Li Chris North

Date: January 26, 2022

Message from the Authors

Dear Editors and Reviewers,

We thank you for your constructive comments, which are valuable and very helpful for revising and improving our paper. We have addressed the comments and incorporated your valuable suggestions in the revised version. The updated contents are colored in red in the revised paper to differentiate with contents in the original one.

We address each comment separately in the following detailed response. The comments we received are boxed, and our responses are written following each comment. We've added more clarification about the relationship between our approach and transfer learning, the validation for our approach, and modelling choices. We hope that you find the revised version satisfactory.

Sincerely,

Huimin Han, Ritvik Prabhu, Timothy Smith, Kshitiz Dhakal, Xing Wei, Song Li, Chris North.

"The writing is clear although some statements are too general and/or no justification is provided. For example, the authors claim that they can change the classification task of a DNN and reuse a previous DNN as the 'majority of parameters in the lower levels of the neural networks do not need to be retrained' but provide no evidence or citation to support this. They appear to be suggesting an approach to transfer learning but never mention this term or related work in this area. They provide no classification results so it is unclear whether the extracted features are useful. Yes, input from the user can suggest new feature weights, but how this impacts classification accuracy is not discussed."

Response

We appreciate your careful review and detailed feedback. Our focus in the revised version was to clearly state the relationship between our approach and transfer learning and the clarification for the usefulness of the extracted features. We hope that you find the following response satisfactory.

Reviewer Comment

"They appear to be suggesting an approach to transfer learning but never mention this term or related work in this area."

Response

Thanks for your review and yes we used a pre-trained model as a feature extractor, which is one of the three major Transfer Learning scenarios¹, we've added more reference papers about using pre-trained model as a feature extractor in the revised manuscript.

Reviewer Comment

"They provide no classification results so it is unclear whether the extracted features are useful. Yes, input from the user can suggest new feature weights, but how this impacts classification accuracy is not discussed."

Response

Thank you for the comment. Talking about those extracted features, the main process we did is to take a CNN pretrained on ImageNet, remove the last fully-connected layer (this layer's outputs are the 1000 class scores for a different task like ImageNet), then treat the rest of the CNN as a fixed feature extractor for the pods dataset. CNN features are more generic in early layers and more

 $^{^{1}} https://cs231 n.github.io/transfer-learning/tf$

original-dataset-specific in later layers, as you've mentioned that it's hard to interpret single feature extracted by a neural network (whether it is efficient or useful for downstream classification tasks), hence we use the visual backpropogation method to interpret those features, we use Andromeda as a tool to enable users to form different projections based on the pods' Phenotypes, then the colormap generated by the visual backpropogation could illustrate whether those features could help distinguish data points or not by applying the learned weights. In future work, the resulting weighted feature models could then be used as classifiers for larger collections of images. We plan to extend these methods to more complex input scenarios, such as images of pods on live plants captured in the field with mobile phones.

:

"The idea of 're-use' of a DNN is interesting and has application in phenomics. The authors did a proof-of-concept of DNN reuse with a study of edamame pods. The reuse concept could have been explored further, including how much is reusable and how much effort is saved vs. starting from scratch. The DNN is a basic pre-trained algorithm, but it is unclear how this model was trained as most information about the model is missing or assumed known by the reader. They show that the model adapts to changes in the input data, which is as expected. Most modeling choices are unexplained/unmotivated (e.g., why this DNN? why MDS and why reduction to 2D?). None of the classification results are provided (accuracy?) nor are the percentage of pods in each class in the data. Validation?"

Response

Thanks for the constructive comments. We discussed more about the modeling choices in our revised papers. We hope that you find the following response satisfactory.

Reviewer Comment

"The reuse concept could have been explored further, including how much is reusable and how much effort is saved vs. starting from scratch."

Response

Thanks for the comments. How much is reusable and how much effort is saved are popular topics that computer vision researchers are doing nowadays for transfer learning. For this paper, we aim to provide investigative image sorting based on one fixed feature extractor. We aim to provide proofs that those features extracted by the pre-trained CNN model could capture human concept like "disease" or "number of seeds".

Reviewer Comment

"The DNN is a basic pre-trained algorithm, but it is unclear how this model was trained as most information about the model is missing or assumed known by the reader."

Response

Thanks for the comments. We have add more reference talking about the pre-trained model that we chose for this paper.

Reviewer Comment

"Most modeling choices are unexplained/unmotivated (e.g., why this DNN? why MDS and why reduction to 2D?)."

Response

We chose ResNet18 as the Neural Network of choice because of the consideration of model complexity as cited in the original paper. As supported by citations included in the original paper, ResNet18 is a powerful neural network model widely used in computer vision tasks. Furthermore, we can confidently claim that using ResNet18 over traditional convectional neural networks is certainly more desirable due to the maximum threshold depth in the traditional convolutional neural networks i.e. the training and testing error percentage rises if the threshold is crossed. For dimension reduction algorithms, MDS is easy to interpret and includes parameters that enable multiple forms of interaction. This is due to WMDS spatializations using distance between observations1 that reflect relative similarity; two points close to each other in a weighted MDS low-dimensional spatialization are considered more similar to each other in the high-dimensional space than are two points far from each other. We reduce the dimension into 2D to enable users implement human sorting, it is not able to do interactive tasks in multi-dimensional space. We cited more papers to support our modelling choices in the revised paper.

Reviewer Comment

"None of the classification results are provided (accuracy?) nor are the percentage of pods in each class in the data. Validation?"

Response

Thanks for the comments, the validation for our approach is by observation level human perception. After the underlying model updating the feature weights, the users could see the new projection, the distance between data points indicate the similarities. Whether the data points could be grouped together by applying the learned weights is how we evaluate whether the CNN model capture the human concepts.

"The approach appears novel and useful. A key weakness is that it is not clear how different methods are evaluated, or how this method performs relative to alternatives. Further, while a software system is described, there is no evidence that the software is available."

Response

We would like to thank you for you positive feedback. Your detailed comments have considerably helped with improving the clarity of the revised manuscript.

Reviewer Comment

"A key weakness is that it is not clear how different methods are evaluated, or how this method performs relative to alternatives."

Response

Thank you for pointing out a valid concern. The reason there has not been any comparisons to alternative methods is because there is no other tool that has been implemented to allow investigative sorting of images. With that being said, we have specifically used this approach as it focuses on the user's cognitive spacial cues. In other words, such an intuitive approach would allow the user to draw a clear distinction between different groups projected on the Andromeda.

Reviewer Comment

"Further, while a software system is described, there is no evidence that the software is available."

Response

Thank you for your feedback. We will be releasing an open-source computational notebook (Jupyter) through Binder that will contain the version of the Andromeda talked about in the manuscript.

This is a proof-of-concept study of how a DNN trained for one classification task on visual images may be re-used for a similar classification task from the same images. The idea of getting input from users during the modeling process is very interesting and may save modeling time and resources. However, it is difficult to determine how classification accuracy changes with user input to the feature space, or how much user input is required for successful transfer of classification accuracy. We see image subspaces that are important to classification, but it is unknown how these compare pre- and post- user input. The amount of data shown is VERY small, particularly for CNNs, and no validation is provided.

Response

Thanks for the constructive comments. We hope that you find the following responses satisfactory.

Reviewer Comment

"However, it is difficult to determine how classification accuracy changes with user input to the feature space, or how much user input is required for successful transfer of classification accuracy."

Response

Preservation of accuracy and how much effort is saved are popular topics that computer vision researchers are exploring for transfer learning. For this paper, we aim to provide investigative image sorting based on one fixed feature extractor. We would like to provide proofs that those features extracted by the pre-trained CNN model could capture human concept like "disease" or "number of pods" through an interactive, cyclical process of human interaction with the data samples.

Reviewer Comment

"We see image subspaces that are important to classification, but it is unknown how these compare pre- and post- user input."

Response

In a system that allows investigative image sorting, highlighting image sub-spaces that are important to classification allows human users to better understand the features that the underlying model might use to represent certain distinguishing characteristics that they find in the data samples. This helps the human users to better understand if the underlying model produces a feature space which can represent characteristics of data samples that the users identify. For use cases in biology, the system might be used to determine if a model can represent a certain phenotype that researchers identify. Alternatively, researchers could interact with data samples in the system to investigate undiscovered phenotypes represented in the model's feature space. The paper has been updated with important image subspaces pre- and post- user input to demonstrate the value added through the investigative process.

"The paper points to an interesting and potentially useful approach, but lacks sufficient information to adopt the approach and lacks quantitative results or comparison with related approaches."

Response

We would like to thank you for you positive feedback. Your detailed comments have considerably helped with improving the clarity of the revised manuscript.

Reviewer Comment

"The paper points to an interesting and potentially useful approach, but lacks sufficient information to adopt the approach and lacks quantitative results or comparison with related approaches."

Response

Thank you for pointing out a valid concern. We have specifically used this approach as it focuses on the user's cognitive spacial cues. In other words, such an intuitive approach would allow the user to draw a clear distinction between different groups projected on the Andromeda.

"Need reference for claim that 'majority of parameters in the lower levels of the neural networks do not need to be retrained'. Relationship to transfer learning?"

Response

We would like to thank you for you positive feedback. Your detailed comments have considerably helped with improving the clarity of the revised manuscript.

Reviewer Comment

"Need reference for claim that 'majority of parameters in the lower levels of the neural networks do not need to be retrained'. Relationship to transfer learning?"

Response

Thank you for your response. We indeed use the process of Transfer Learning in this project. The ResNet-18 is an 18 layer neural network trained on over a million images from the ImageNet database. Due to the high accuracy of the model, we decided to re-purpose a few higher layers of the neural network to make it relevant to our data set on hand. In fact, we re-purposed the neural network differently for the two data sets using transfer learning. In our research, our approach is to use multidimensional reduction to provide a 2D visualization for humans to provide interactive input that, through inverse dimensionality reduction, is imparted on the model through adjusted feature weights. Those features are features captured in the top-most layer(s) of the original model.