
Interpreting Deep Learning: ISIC Melanoma

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
Matteo Russo
Princeton University
matteor@princeton.edu

Dale Lee
Princeton University
dalelee@princeton.edu

Andrew Griffin
Princeton University
andrewjg@princeton.edu

Abstract

Can an algorithm save a life? In the field of dermatology, an early diagnosis can make the difference between life and death. Skin cancer is the most prevalent kind of cancer in the United States, with 1 in 5 Americans developing skin cancer by the age of 70 [3]. Melanoma is the most fatal type of skin cancer, but can often be cured if detected early enough, with the 5 year survival rate for early-stage melanoma being 99% [4]. This unique aspect of this particular kind of cancer makes it ideal for a computer vision application since potential lesions are easy to identify, image, and visually diagnose, unlike other kinds of internal cancers that require a more vigorous biopsy procedure. In this project, we establish baseline classification techniques on melanoma image classification, data exploration and visualization via dimensionality reduction, as well as exploring Convolutional Neural Network approach, providing interpretability explanations of the predictions through three methodologies: Activation Maps Backpropagation, LIME and SHAP values. We found the highest performance with the CNN approach. We also found that the interpretable predictions truly shed light on the model and onto the degree to which we should trust the Deep Learning prediction, both of which are crucial in the field of bio-medicine.

1 Introduction

This report will show a natural progression through various machine learning techniques in order to complete the important task of classification of skin lesion images as benign or malignant. There will also be an emphasis on interpretation of models, especially given the fact that the classification task is currently being done solely by human doctors and a potential switch from humans to algorithms will not be reliable unless backed up by interpretations that help explain what the algorithm is picking up from the image data. Even without the possibility of replacement of doctors by algorithms, interpreting these models can give doctors a very powerful tool, and complement the biology and anatomy that has been studied on the topic. The progression towards effective classification of images will begin with the baseline classification techniques of random forest and logistic regression. These techniques are effective at classification of many types of data including text data for sentiment analysis, data on students for educational outcomes[6], etc. These techniques are general and give us something to compare to the results of our more complicated models. We see that the general classification techniques are not great, yielding only about 77% accuracy on 224x224 pixel images. This motivates us to use other techniques to analyze the data and attempt to gain a greater understanding for the underlying structure and patterns.

This will bring us to unsupervised learning in the form of dimension reduction techniques and visualizations. By performing PCA and visualizing the corresponding clusters with respect to the first

054 two principle components, we see that there is a large amount of class overlap. This confirms our
055 intuition that a linear classifier will not suffice for the data in question, and further enforces the
056 hypothesis that we need a better classifier.
057

058 Finally, we will implement a more complicated classifier using a Convolutional Neural Network
059 (CNN) in order to gain greater accuracy in classification. CNNs are great for image data as they take
060 advantage of local patterns and spatial considerations that exist in images. The one disadvantage
061 of using a complex nonlinear model such as CNN is the difficulty of interpretability. As previously
062 stated, this interpretability is key in the medical field where insights must be backed up by reliable
063 explanations. These interpretations will come in the form of visualizations of algorithms assigning
064 importance to different pixels and objects in an image. The effectiveness of the CNNs compared
065 to the baseline classifiers with the corresponding visual interpretations give us compelling evidence
066 that machine learning techniques *can*, in fact, be extremely useful for medical applications such as
067 melanoma detection.

068 We used the HAM10000 ISIC 2017 dataset from the International Skin Imaging Collaboration
069 Melanoma Project, which is a public archive of dermoscopic images of skin lesions [2]. Because
070 computation with RGB image data can be quite time complex, we used a reduced dataset of 1449
071 224x224px RGB malignant images and 1800 224x224 px RGB benign images via random sampling.
072 We scaled the original images which were 500x1000px to 224x224px, using using reverse mapping
073 from the destination to source image with the cv2 resample library. This method of image reduction
074 has found success in many implementations [1].
075

076 **2 Related Work**

077 Neural Networks are widely used in image classification—in particular, convolutional neural net-
078 works have been found to reduce computation time while reducing the chance of overfitting and
079 providing improved results on smaller datasets. Convolutional Neural Networks in the biomedical
080 imaging sphere have gained widespread success in classifying and segmenting images [9]. A re-
081 cent study published in Nature by Stanford researchers found that a CNN was able to outperform
082 board-certified dermatologists when using a pre-trained GoogleNet Inception v3 CNN [5].
083

084 There is also literature on fine-tuning an already trained network for image classification versus
085 creating one from scratch.[13] Because of issues with computational power and the need for large
086 data sets, many researchers use existing networks trained on other images for classification. We
087 believe melanoma images are sufficiently different from general image training data, therefore we
088 choose to fully train rather than fine tune an existing network.
089

090 **3 Methods**

091 **3.1 Baseline Classification Methods**

092 We use two standard classifiers to benchmark a baseline accuracy for the ISIC Melanoma Classifi-
093 cation Task:
094

- 095 • Random Forest with a max tree depth of 5 and number of estimator trees of 20.
- 096 • Logistic Regression with ℓ_2 regularization.
097

098 We found the above mentioned hyper parameters through 10 fold Grid-Search Cross Validation. We
099 searched over 5, 10, 20, 35, 50 number of trees in the random forest, with tree depths of 3 and 5. For
100 logistic regression, we searched over 11 and 12 penalty space.
101

102 **3.2 Dimensionality Reduction and Data Visualization**

103 Due to the large number of features we test on, our image classifier could suffer from the “curse
104 of dimensionality” and lead to overfitting of the training data due to the high number of irrelevant
105 parameters [11]. We decided to test this possibility by performing PCA with two principal compo-
106 nents, and visualizing the results to see what the principal components were in Figure 9. We did
107

not find a stark separation of clusters with PCA, however, so we decided to perform t-SNE projection to two dimensions, and found a sharper divide between malignant and benign dimensionally reduced lesions. We recognize the fact that t-SNE performs well with low-dimensional samples and that is the reason why we decided to first reduce the dataset through PCA and then apply t-SNE to the linearly reduced dataset. Even in this case, nonetheless, we acknowledge that this pipeline of dimension reduction has no theoretical ground but it gives similar performances in the visualization to the ones t-SNE provides. We could visualize t-SNE reduction in Fig.1, PCA reduction in Fig.9 and PCA followed by t-SNE in Fig. 11 and Fig. 12.

3.3 Convolutional Neural Network (CNN)

A Convolutional Neural Network is a machine learning algorithm meant to capture the importance of different aspects of images to their classification. It performs remarkably well with image data as it takes into account spatial consideration of matrices (pixels) by applying filters to subsets of the pixels, then pooling together the output of those filters in order to reduce dimensions and capture the most important features.

3.4 CNN Visualizations

We will see in Section 4 that convolutional neural networks perform much better at the classification of benign and malignant lesions in the ISIC dataset. However, because the neural network is very complex and highly nonlinear, it can be difficult to interpret. Specific to a medical dataset like ISIC, this interpretation is paramount. It is unlikely that professionals in the medical field will be comfortable with trusting the output of a model where nobody can understand exactly why classifications are being made. The interpretability of a complex model, in general, is as important if not more important than the model prediction itself. It is often the case that there exists a true dichotomy between the complexity of the model, and thus, the hardness of the patterns it is able to capture and its interpretability. As is the case with many machine learning techniques, visualizations can shed light on the most important aspects of the model. The visualizations that are used in this report are class activation maps as well as LIME and SHAP visualizations.

3.4.1 Class Activation Maps

Class Activation maps allow us to see which areas in the image are relevant to the classification of a specific class. For example, CAMs are used in order to show the importance of faces and limbs in the classification of humans along or the importance of ears and noses in the classification of dogs and other animals.[14] As highlighted in Section 3.4, Neural Networks are constructed in convolutional layers created by different filtering of images by a kernel. Each convolutional layer extracts high-level features such as edges or patterns in the input image, and by analyzing CAMs at each layer we can get an appreciation for how each layer extracts a different feature. Some layers clearly extract the shape of the objects in the image, other layers clearly extract color, and yet others extract combinations of features that are not as easily recognizable. We will look at these CAMs with respect to properly classified samples as well as misclassified samples to see exactly how the CNN model assigns importance to areas in the picture for classification. In other words, the Activation Map Analysis we have carried out relies on a white-box knowledge of the network itself, meaning that we have perfect knowledge over all the components. As a consequence, albeit their strong visual interpretability, Activation Maps do not scale well to other models.

3.4.2 LIME

Urging to abstract oneself from the perfect knowledge of the network, Ribeiro et al.'s LIME (Local Interpretable Model-agnostic Explanations [10]) performs a black-box locally interpretable model evaluation by perturbing, within a neighborhood, specific samples' feature values and measuring the resulting impact in the classification. Formally, let $g \in \mathcal{G}$ be an interpretable model from the set of interpretable models \mathcal{G} and let f be the original globally unknown model. Further, let $\Omega(g)$ be a measure of model g complexity (such as Radamacher's Complexity) and let π_x be a definition of locality (e.g. neighborhood) for point x . Last, let $\mathcal{L}(g, f, \pi_x)$ be the measure of mistrust we have on model g in approximating f within locality π_x . Then, LIME solves the following minimization

162 problem:

163

$$\lambda(\mathbf{x}) = \arg \min_{g \in \mathcal{G}} \mathcal{L}(g, f, \pi_{\mathbf{x}}) + \Omega(g) \quad (1)$$

164
165
166 The notion of locality $\pi_{\mathbf{x}}$ in \mathcal{L} allows us not to worry about the overall model f but to only consider
167 a neighborhood approximation of it.

168
169 **3.4.3 SHAP**

170 As much as LIME, with the aim of visual interpretability, rather than explaining the complex model
171 f in its entirety, Lundeberg and Lee in [8] propose to study the model f itself when evaluated at a
172 single data point, as per the following equation, where g represents the additively binary simplified
173 version of f , ϕ the vector of feature coefficients for each of the covariate components and \mathbb{I} is the
174 indicator which is equal to 0 if the feature is not present and 1 if it is:

175
176
$$g = \phi^T(\mathbf{x})\mathbb{I}(\mathbf{x}) \quad (2)$$

177 The only question that remains unanswered regards how to choose such that the following four
178 notions of fairness are fulfilled for any sample \mathbf{x} :

- 179
180 1. *Perfect redistribution*: $g(\mathbf{x}) = f(\mathbf{x})$.
- 181 2. *Equality*: features contributing equally to the classification should have the same SHAP
182 value.
- 183 3. *Neutrality*: features not contributing should have null SHAP value.
- 184 4. *Additivity*: $g(\mathbf{x} + \mathbf{y}) = g(\mathbf{x}) + g(\mathbf{y})$.

185 The following expression gives us the feature coefficients that would satisfy the four properties of
186 above:

187
188
$$\phi_i = \sum_{S \subseteq \mathcal{F} \setminus \{i\}} \frac{|S|!(|\mathcal{F}| - |S| - 1)!}{|\mathcal{F}|!} [f_{S \cup \{i\}} - f_S] \quad (3)$$

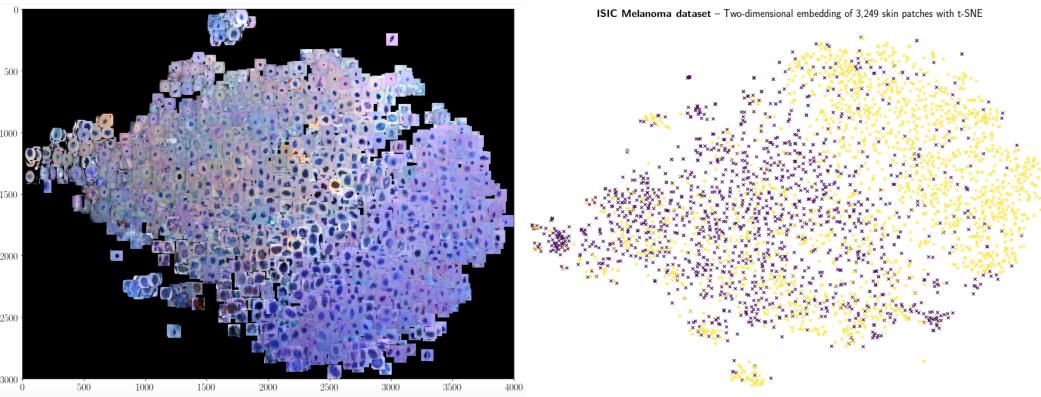
189 Hereby, \mathcal{F} represents the set of all features, i represents the i th component of feature map ϕ and
190 S any subset of $\mathcal{F} \setminus \{i\}$. In other words, SHAP calculates the importance of each feature (in our
191 case, pixel), in any possible order so to avoid unfair order bias in features comparison. As we may
192 observe from Eq.2, the problem of finding every possible subset $S \subseteq \mathcal{F} \setminus \{i\}$ and retraining on
193 each of them explodes combinatorially, and, thus, we need convex approximation techniques to gain
194 computational efficiency. For a Deep Neural Network model f , as thoroughly explained in [8],
195 SHAP is replaced by a technique named Deep SHAP. The latter links Shapley values computed in
196 small local portions of the DNN to the ones for the whole network, via a recursive backpropagation
197 of DeepLIFT multipliers ([12]).

200 **4 Results**

201
202 **4.1 Exploratory Analysis: PCA and t-SNE Feature Reduction**

203 When observing the PCA transformed dataset pictured in Figure 9, we found that the two directions
204 of maximum variance appeared to be color and size—images that were more blue appeared to cluster
205 higher in the plot, while larger lesion images tended towards the right side of the plot, and smaller
206 lesions tended towards the left. However, PCA by itself did not provide a very clear separation
207 between the malignant and benign lesions. Because the blue tinted images seem to correspond
208 most with the cluster of benign images, we considered the fact that this could be an artifact of the
209 data, however, we found that matplotlib was modifying the background color of the images, but the
210 correctly colored images were being input to the classifiers.

211 We found improved separation between malignant and benign lesions through first applying PCA,
212 then t-SNE to the dataset, as shown in Figure 11 and Figure 12 in the Appendix. PCA is only taking
213 the principal components of the images, while t-SNE is perhaps a better application to nonlinear data.
214 This was confirmed when we applied t-SNE by itself to the dataset, finding that this dimensionality
215 reduction methodology offered the clearest separation between the two kinds of lesions.



216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
Figure 1: t-SNE of benign and malignant images. Yellow datapoints are benign, and black datapoints are malignant.

4.2 Baseline Classification Results

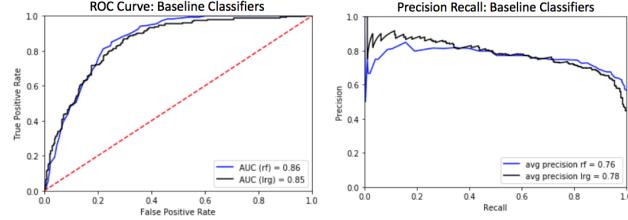


Figure 2: ROC curves with AUC in legend. Precision-Recall curves, with average precision scores in legend.

To evaluate the performance of the classifiers on our test data set, we first computed the average accuracy, average precision, and false negative rate, as shown in Table 1. We generated Receiver Operating Characteristic (ROC) curves and precision recall curves, as shown in Fig.6. The accuracy we received from simple baseline classification was surprisingly high—we found accuracy to be **76.3%** for random forest and **78.5%** for logistic regression. However, we found our false negative rate to be **9.2%** for random forest and **12.5%** for logistic regression. While we were initially impressed with the accuracy of these models, we realized that a false negative rate of around **10%** means that for one in ten patients, we would falsely diagnose that their melanoma was benign. This, as can be expected, can be disastrous for any sort of biomedical application and oncological in particular. Thus, our models urge to focus on decreasing this false negative rate, while also increasing accuracy. Table 1 summarizes all the information about classifiers performances.

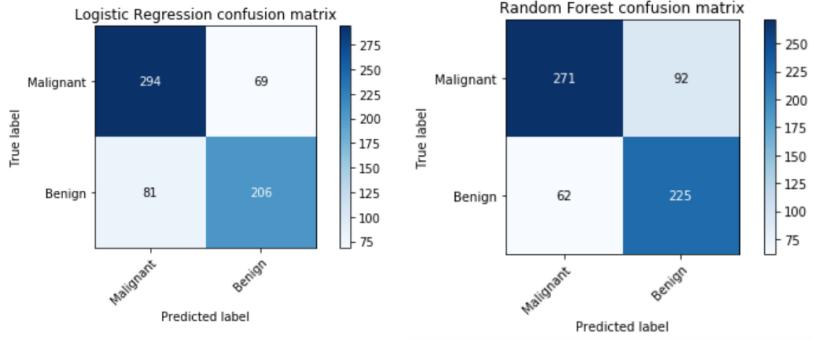
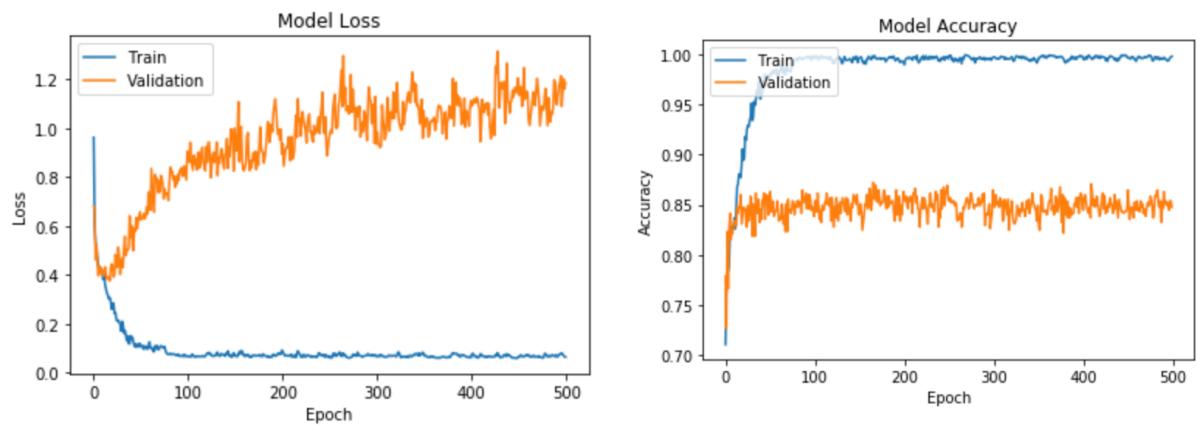


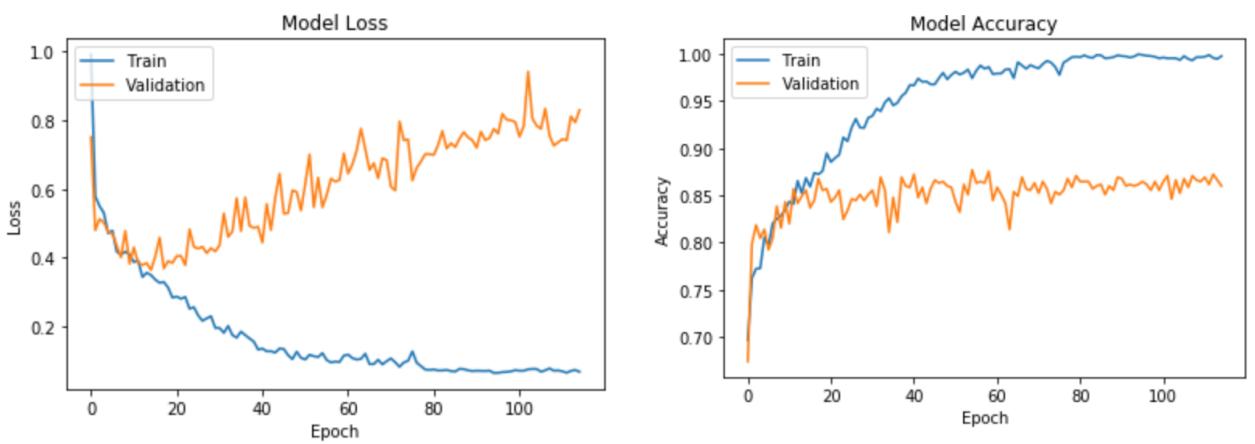
Figure 3: Confusion Matrices for Logistic Regression and Random Forest Classifiers

270 **4.3 Convolutional Neural Network Classification Results**
 271

272 As per [7], we have implemented a ResNet, whose full specification is outlined in the Appendix
 273 at 2. We have chosen this model because it is not too deep and, given the fact that our dataset is
 274 relatively small when compared to the cardinality of DNN parameters, ResNet guarantees high
 275 performance rates as explained later in this section. The neural network described above was
 276 trained using the training set of 2600 images from the ISIC dataset. We began by training the CNN
 277 using 500 epochs as shown in Fig.4 and obtained an accuracy of **81%**, higher than either of the two
 278 baselines (**76%**). When looking at the plot of model accuracy, we noticed the model overfitting to
 279 the data as shown by the validation loss increasing with an increase in the number of epochs. This
 280 lead us to look into early stop regularization, not allowing the CNN to continue if the validation loss
 281 continues to rise. As we see in Fig.5, this reduced the number of epochs to slightly over 100 and
 282 resulted in an improved accuracy of **87%**.
 283



298 Figure 4: Model Loss and Model Accuracy for Convolutional Neural Network run with 500 epochs
 299 and no early stop regularization.
 300



318 Figure 5: Model Loss and Model Accuracy for Convolutional Neural Network run with 100 epochs
 319 completed with early stop regularization
 320

321 At the same time, as we can observe from the confusion matrix of below, the ResNet FNR (the
 322 fraction of melanoma erroneous diagnoses) is **7.1%**, which is about a **5%** improvement from
 323 Logistic Regression and a **2%** one from Random Forest.

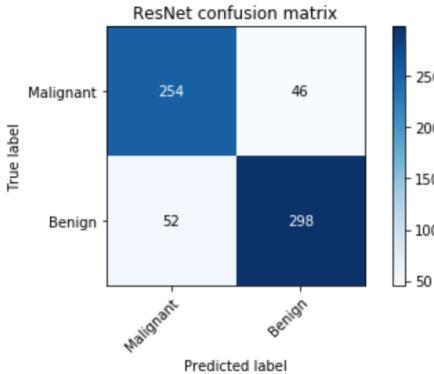


Figure 6: Confusion Matrix for ResNet Classifier.

As mentioned before, Table 1 summarizes all the information about classifiers performances.

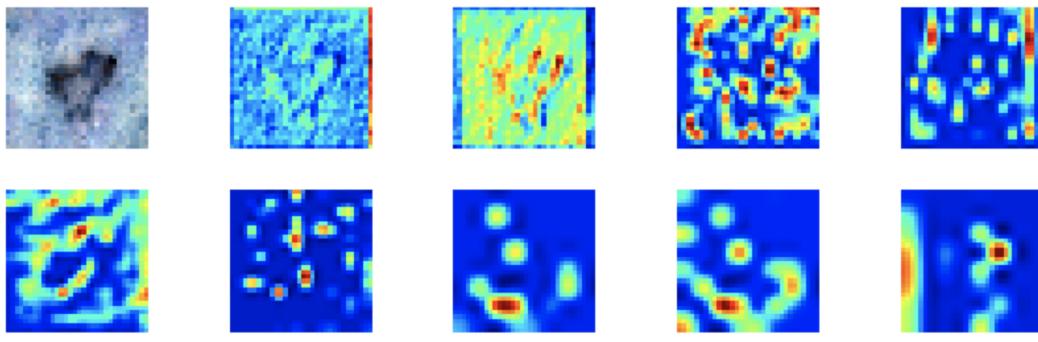
Classifier	Accuracy	Precision	FPR	FNR
Logistic Regression	0.769	0.785	0.106	0.125
Random Forest	0.763	0.764	0.144	0.092
ResNet	0.872	0.764	0.080	0.071

Table 1: Results from classifiers on test set with 80%-20% train-test split (2599 samples in training set and 650 in testing set).

This gives us confidence that the neural network is indeed a more powerful classifier than the baseline, however we are still left with the task of interpretation of the model as previously stated. For this we look at a correctly classified malignant image alongside the activation maps, LIME, and SHAP visualizations.

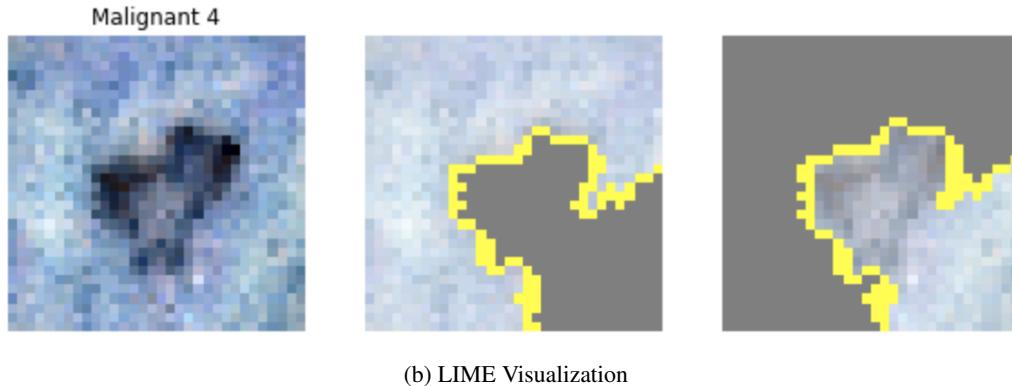
4.4 Comparison of CNN Visualizations

Fig. 7 shows each of the neural net visualizations on a correctly classified benign image from the testing set. Looking at the Class Activation Map, we can see the evolution of the convolutions and how they pick up on different important features of the image as described in Section 3.4. The first few layers look to be focusing on general features like shape and surroundings of the main center of mass. As the layers progress, areas of the lesion which are darker near the edge of the lesion are assigned more importance. This makes sense biologically as we will see in Section 4.5. Further confirmation of the oncological interpretation comes from the LIME and SHAP visualizations.



(a) Class Activation Map by Layer

378
 379 Fig. 7b shows in its second and third image shows which part of the image explains respectively the
 380 benign and the malignant class. According to LIME, in fact, the lesion explains the malignant class
 381 which is what we would expect from a biomedical interpretation. Indeed, the image is correctly
 382 classified by the ResNet as malignant.
 383



393
 394 On the other hand, according to SHAP, in Fig.7c) second image, we see that all the "hot" pixel
 395 values (in red) reside outside the lesion region of the image, further confirming that the benign class
 396 is explained by outside pixels: this implies that the boundaries are a really discriminatory feature (the
 397 SHAP last two images are one the color inverse of the other given that they explain complementary
 398 classes).
 399

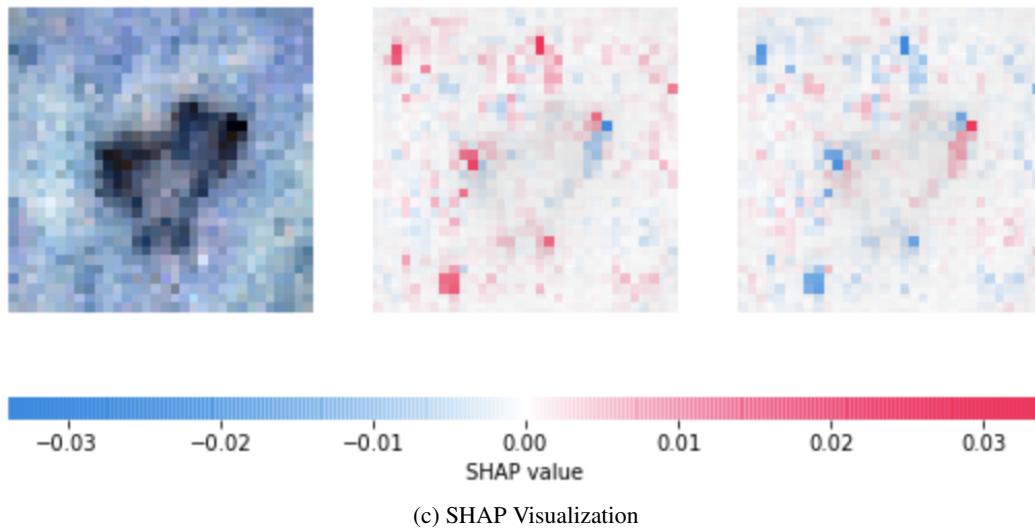


Figure 7: Correctly Classified Image

4.5 Biological Discussion

423 We were interested in comparing the areas of greatest weight in the CNNs classification with a
 424 visual diagnosis a melanoma lesion and see if there was significant overlap. The five key signs in
 425 visually diagnosing melanoma follow the ABCDE acronym: Asymmetry (half of the lesion does not
 426 match the other half), Border irregularity, Color (not uniform or featuring multiple colors), Diameter
 427 greater than 3mm, and Evolution of size, shape or color.
 428



Figure 8: Full resolution image of correctly classified Melanoma image from Fig.7. Circles indicate areas in the original image with high SHAP values.

Interestingly, it appeared that there was not overwhelming overlap—the only major aspect that we noticed was that SHAP values appeared marginally higher around the lesion border, in the red circles displayed in Fig. 8. Its notable that the points with highest SHAP values also corresponded to areas of color contrast within the images, which could indicate the importance of color in the CNN’s interpretation.

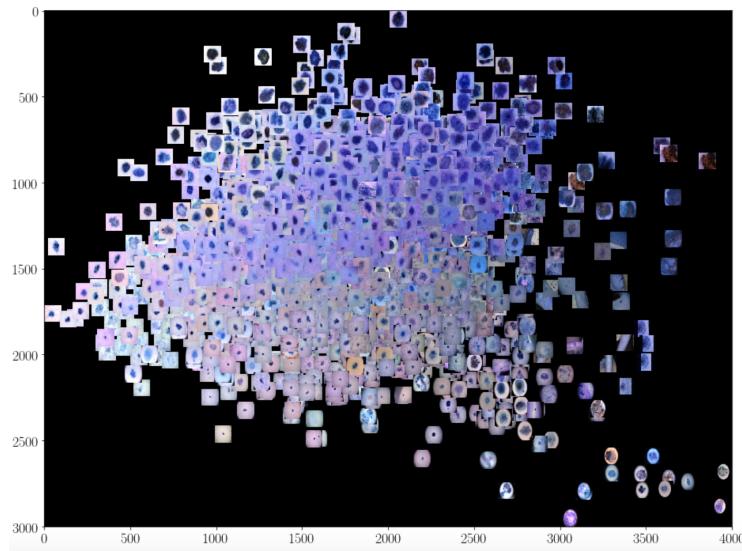
5 Discussion and Conclusion

From this project, we find reasonable accuracy from baseline classification methods of random forest and logistic regression, but find a much higher accuracy rate of **87%** with a Convolutional Neural Network. As we have mentioned, this accuracy comes at a bit of a cost with respect to interpretability given the non-linearity of the classifier. For this reason, we created visualizations of the model for our image data (Activation Maps, LIME and SHAP) in order to understand a few key features of the data. We see that the parameters of greatest variance are size and color of the lesion—reasonable findings given these two characteristics are two out of the five most important characteristics used to diagnose melanoma by human doctors.

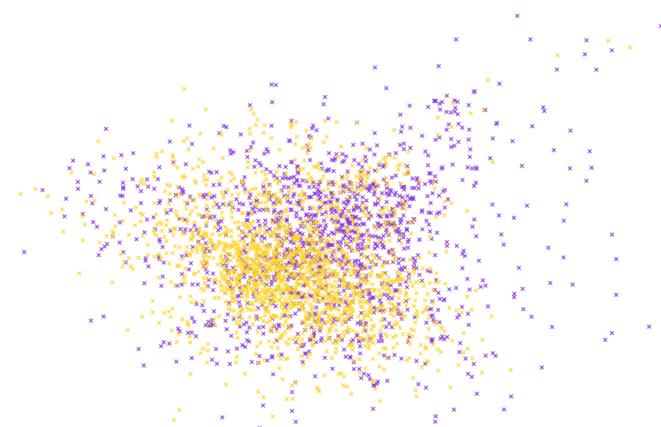
This ultimately reinforces the desired relationship between machine learning and medical diagnosis. Deep learning techniques can be used for diagnosis with a very high accuracy, however the complementary function for human experts that comes along with the visualization may indeed prove more important—especially while deep learning techniques are still young as they are now. We can be optimistic that these methods will continue to evolve to improve medical diagnosis along with our understanding of medical processes at a deeper level.

486
487
488
489
490 **6 Appendix**
491
492
493
494
495

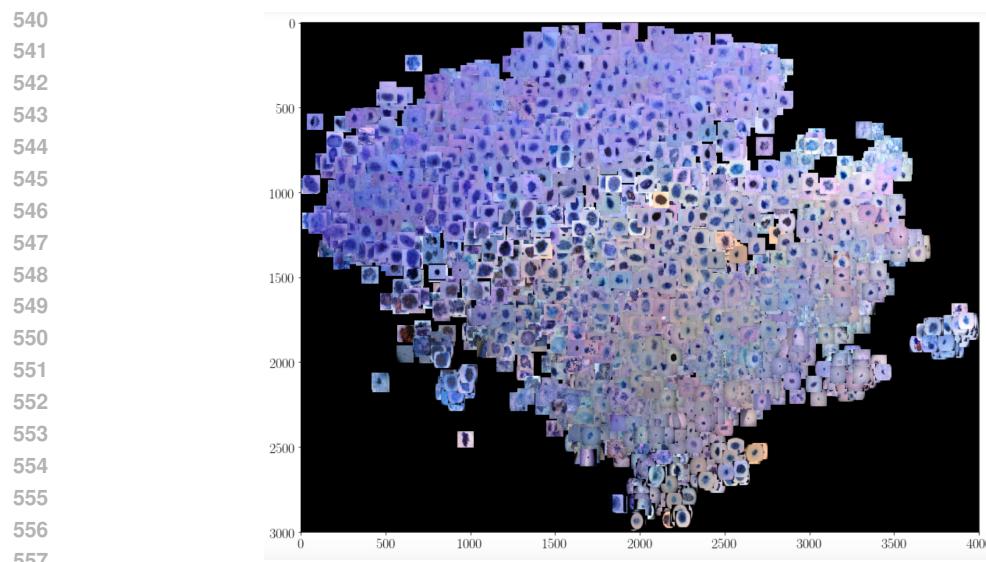
496 **6.1 Dimensionality Reduction**
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513



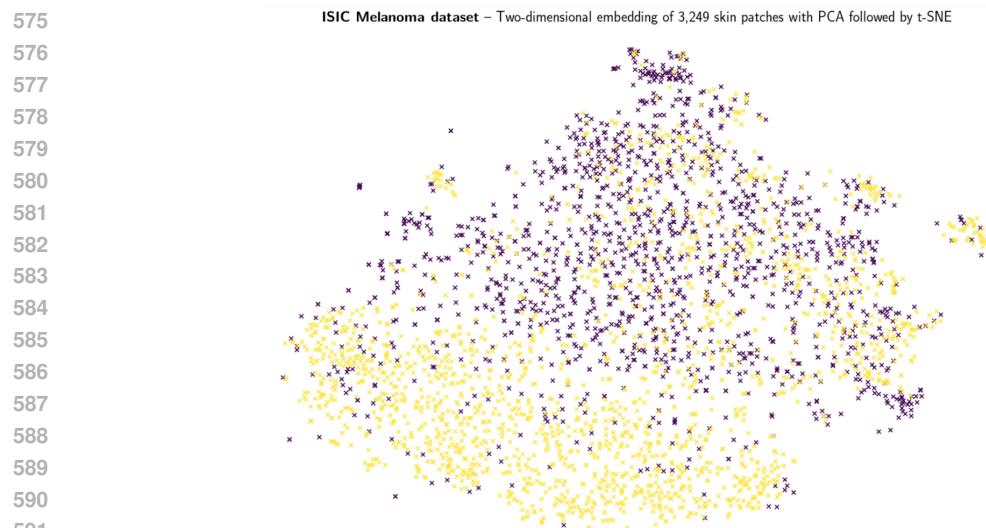
514 Figure 9: PCA of benign and malignant images.
515
516
517
518
519
520
521
522 ISIC Melanoma dataset – Two-dimensional embedding of 3,249 skin patches with PCA
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538



539 Figure 10: PCA on benign and malignant images. Yellow datapoints are benign, and black data-
points are malignant.



558 Figure 11: Visualization of class separation after performing t-SNE on PCA reduced dataset.
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574

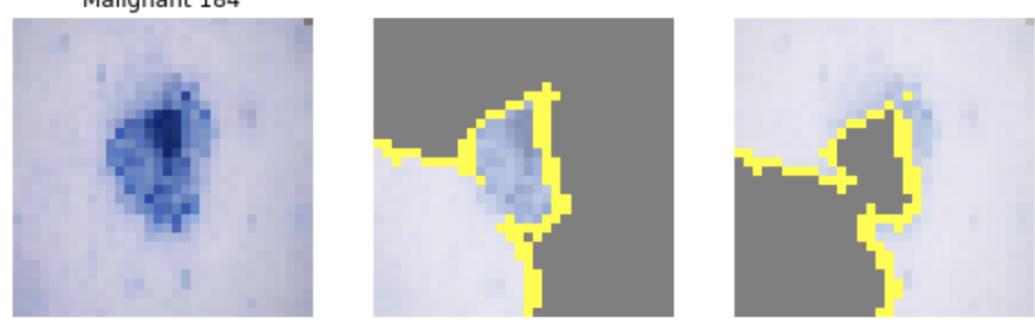


592 Figure 12: t-SNE on PCA reduced dataset of benign and malignant images. Yellow datapoints are
593 benign, and black datapoints are malignant.

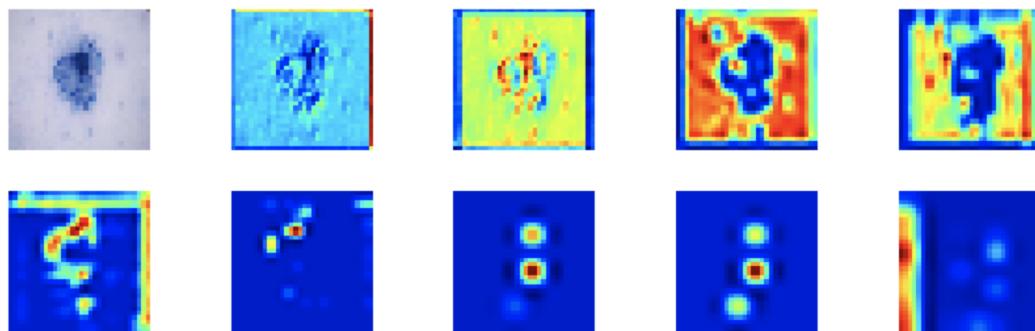
594
 595 **6.2 Deep Convolutional Neural Network Architecture**
 596
 597

Layer	Type	Output Shape	Parameters
conv2d_1	Conv2D	(None, 32, 32, 32)	896
activation_1	Activation	(None, 32, 32, 32)	0
batch_normalization_1	Batch Norm	(None, 32, 32, 32)	128
conv2d_2	Conv2D	(None, 32, 32, 32)	9248
activation_2	Activation	(None, 32, 32, 32)	0
batch_normalization_2	Batch Norm	(None, 32, 32, 32)	128
max_pooling2d_1	Max Pooling	(None, 16, 16, 32)	0
dropout_1	Dropout	(None, 16, 16, 32)	0
conv2d_3	Conv2D	(None, 16, 16, 64)	18496
activation_3	Activation	(None, 16, 16, 64)	0
batch_normalization_3	Batch Norm	(None, 16, 16, 64)	256
conv2d_4	Conv2D	(None, 16, 16, 64)	36928
activation_4	Activation	(None, 16, 16, 64)	0
batch_normalization_4	Batch Norm	(None, 16, 16, 64)	256
max_pooling2d_2	Max Pooling	(None, 8, 8, 64)	0
dropout_2	Dropout	(None, 8, 8, 64)	0
conv2d_5	Conv2D	(None, 8, 8, 128)	73856
activation_5	Activation	(None, 8, 8, 128)	0
batch_normalization_5	Batch Norm	(None, 8, 8, 128)	512
conv2d_6	Conv2D	(None, 8, 8, 128)	147584
activation_6	Activation	(None, 8, 8, 128)	0
batch_normalization_6	Batch Norm	(None, 8, 8, 128)	512
max_pooling2d_3	Max Pooling	(None, 4, 4, 128)	0
dropout_3	Dropout	(None, 4, 4, 128)	0
flatten_1	Flatten	(None, 2048)	0
dense_1	Dense	(None, 2)	4098

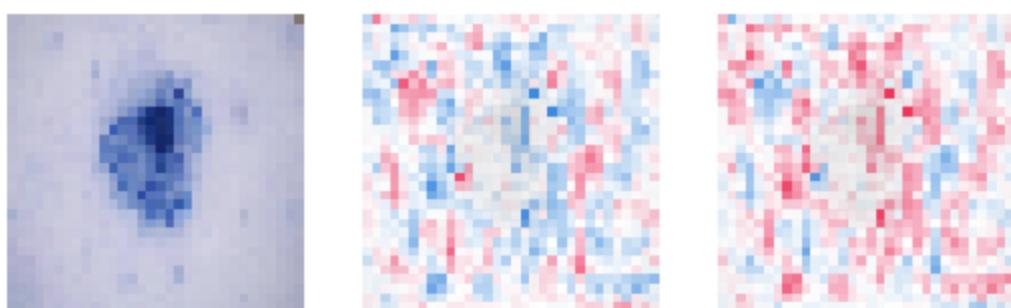
624 Table 2: Full ResNet architecture specification
 625
 626
 627
 628
 629
 630 Thus, we have that:
 631
 632 • Total number of parameters: 292,898
 633
 634 • Trainable number of parameters: 292,002
 635
 636 • Non-trainable number of parameters: 896
 637
 638
 639
 640 **6.3 Interpretable Results in High-Confidence Misclassification Settings**
 641
 642 We conduct the same type of analysis as in Section 4.4 but with respect to the following kind of
 643 image: the image has been misclassified with high degree of confidence. In practical terms, the
 644 model has predicted a benign tumour with high confidence, when that patient had, in fact, melanoma.
 645 We can observe how the activation maps across layers are not capable of putting under clear marking
 646 a melanoma lesion. We can also see how the "hot" zones in the third activation map are at the
 647 exterior, meaning that not only it has detected the boundaries, but also is fooled into believing that
 648 the distinctive characteristic of that sample lie in the surroundings of the melanoma center of mass.



(a) LIME Visualization



(b) Class Activation Map by Layer



-0.0010 -0.0005 0.0000 0.0005 0.0010
SHAP value

(c) SHAP Visualization

Figure 13: Misclassified Image

702
703 **References**

- 704 [1] Different methods of image mapping, its advantages and disadvantages.
705 [2] International skin imaging collaboration: Melanoma project. 2017.
706 [3] Skin cancer. american academy of dermatology association.
707 [4] Key statistics for melanoma. *American Cancer Society*, 2018.
708 [5] Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau,
709 and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural
710 networks. *Nature*, 542, 2017.
711 [6] Dorina Kabakchieva. Predicting student performance by using data mining methods for clas-
712 sification. *Cybernetics and information technologies*, 13(1):61–72, 2013.
713 [7] Abhijeet Kumar. Object-recognition-cifar-10. [https://github.com/abhijeet3922/
714 Object-recognition-CIFAR-10](https://github.com/abhijeet3922/Object-recognition-CIFAR-10), 2018.
715 [8] Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *CoRR*,
716 abs/1705.07874, 2017.
717 [9] Philipp Fischer Olaf Ronneberger and Thomas Brox. U-net: Convolutional networks for
718 biomedical image segmentation. *Springer Link*, 2015.
719 [10] Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explain-
720 ing the predictions of any classifier. *CoRR*, abs/1602.04938, 2016.
721 [11] Chiara Sabatti. Basic statistical references for the design and analysis of gene-chips experi-
722 ments.
723 [12] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features
724 through propagating activation differences. *CoRR*, abs/1704.02685, 2017.
725 [13] Nima Tajbakhsh, Jae Y Shin, Suryakanth R Gurudu, R Todd Hurst, Christopher B Kendall,
726 Michael B Gotway, and Jianming Liang. Convolutional neural networks for medical image
727 analysis: Full training or fine tuning? *IEEE transactions on medical imaging*, 35(5):1299–
728 1312, 2016.
729 [14] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning
730 deep features for discriminative localization. In *Proceedings of the IEEE conference on com-
731 puter vision and pattern recognition*, pages 2921–2929, 2016.
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755