

Classifying skin lesions

By Mohamed Elghetany, Megan Martin, Mike Khor

Background & Objective

- Skin cancer is the most common cancer worldwide¹
- *Early and correct* classification impacts treatment
- Lesion appearance is main driver of biopsy/treatment

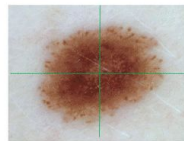
Objective:

Train a model to classify skin lesion images for different types of cancers and benign lesions

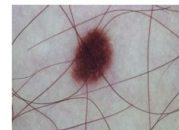
Benign

Asymmetry

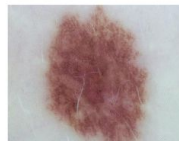
Both sides match other



Border: Regular Edges

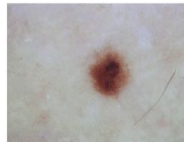


Colour: Consistent Shades



Diameter

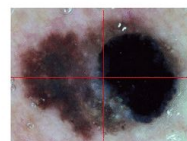
Lesion is smaller than 6mm



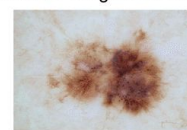
Melanoma

Asymmetry

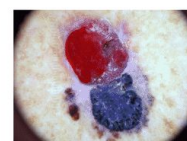
One side does not match other



Border: Irregular or Blurred

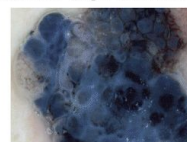


Colour: Different Shades



Diameter

Lesion is larger than 6mm



About the Data

- Dataset obtained from International Skin Imaging Collaboration (ISIC)³
- Utilizing ISIC 2019 dataset which contains dermoscopic images for **8** diagnostic categories
 - Patient metadata for each image
- Dataset contains **25,331** images
 - 21,311 images with complete metadata
 - **Unbalanced** dataset - most are labeled as benign

	image	age_approx	anatom_site_general	lesion_id	sex
0	ISIC_0000000	55.0	anterior torso	NaN	female
1	ISIC_0000001	30.0	anterior torso	NaN	female
2	ISIC_0000002	60.0	upper extremity	NaN	female
3	ISIC_0000003	30.0	upper extremity	NaN	male
4	ISIC_0000004	80.0	posterior torso	NaN	male
...
25326	ISIC_0073247	85.0	head/neck	BCN_0003925	female
25327	ISIC_0073248	65.0	anterior torso	BCN_0001819	male
25328	ISIC_0073249	70.0	lower extremity	BCN_0001085	male
25329	ISIC_0073251	55.0	palms/soles	BCN_0002083	female
25330	ISIC_0073254	50.0	upper extremity	BCN_0001079	male
25331 rows × 5 columns					

Metadata for ISIC 2019 dataset.

Addressing Class Imbalance

- >50% images in a single class (benign moles)
- Downsampled or upsampled classes to 2,000 images each, resulting in 16,000 total images for balanced dataset
- Augmentation including random flips, rotations, saturation, contrast, brightness changes were utilized to upsample

MEL	NV	BCC	AK	BKL	DF	VASC	SCC
4522.0	12875.0	3323.0	867.0	2624.0	239.0	253.0	628.0

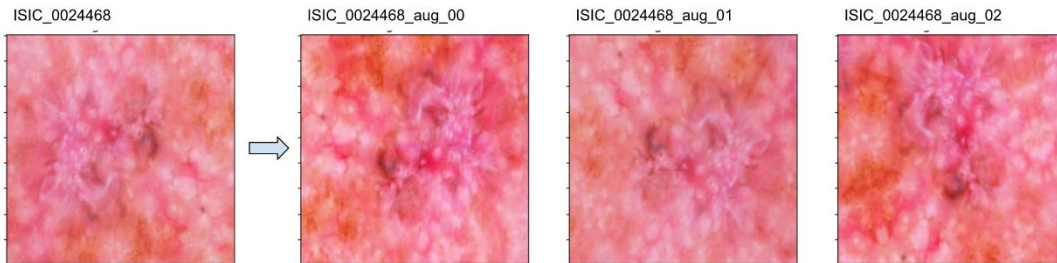
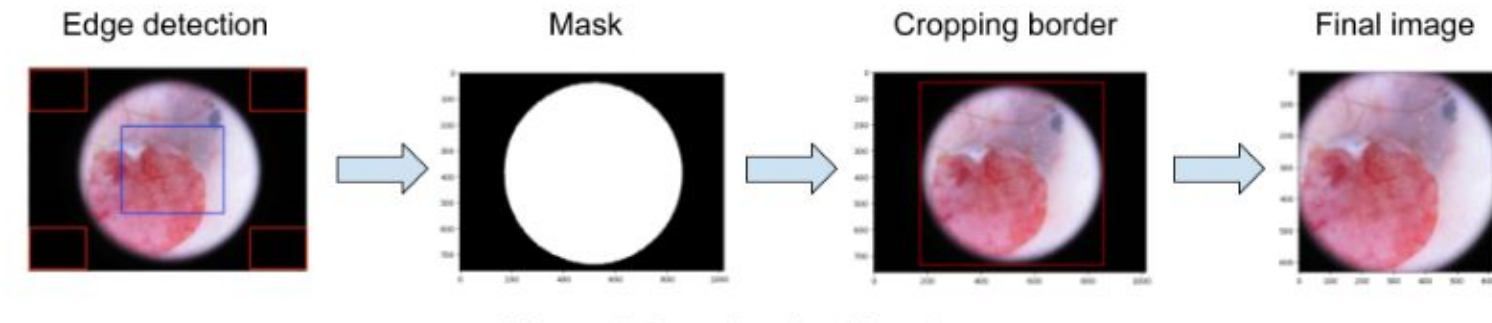
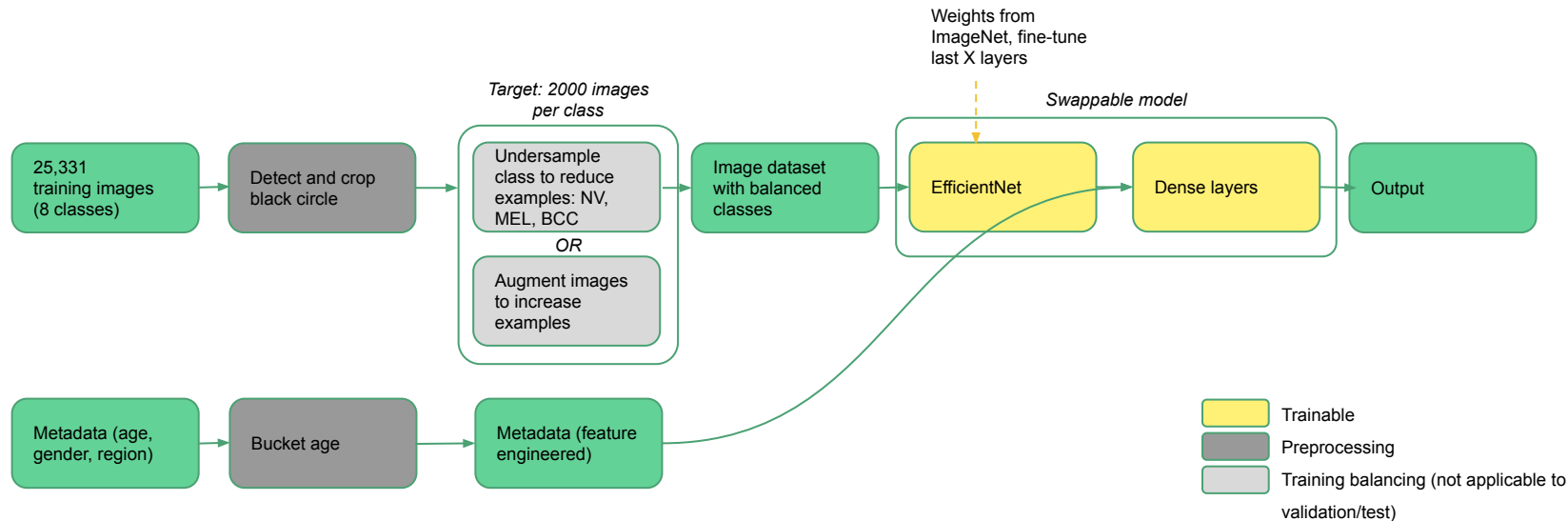


Image pre-processing

- Scaled down varied image sizes to 224 x224
- 25% of images had black regions surrounding lesion
 - *Detection*: does this image need cropping? Use value difference between corner and center.
 - *Mask*: find all pixels that are above a threshold value – “light region”
 - *Determine crop border*: min/max position of “light region”
 - *Crop*: produce final image



Block Diagram



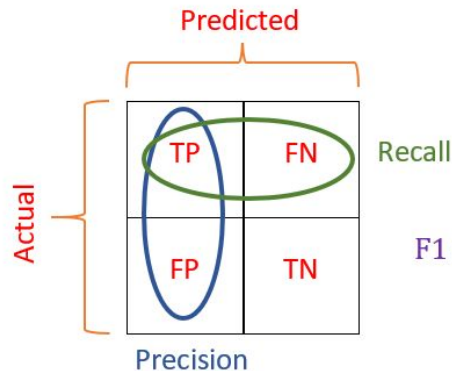
Tech Stack

- Jupyter Lab in Amazon Sagemaker (varying EC2 sizes)
- Python 3 image: Tensorflow 2.6, Python 3.8 GPU optimized
- Scikit-Learn
- Keras
- Numpy, Matplotlib, Pandas

Evaluation and Success Parameters

Success defined as:

- Recall: > 60%
- Precision: >50%
- ROC AUC: >60%
- Accuracy: >55%



$$F1 \text{ score} = \frac{2 * (\text{precision} * \text{recall})}{\text{precision} + \text{recall}}$$

Models and Methods

1. Logistic Regression
2. K-Nearest Neighbors (KNN)
3. Random Forest
4. Convolutional Neural Network (CNN)
5. CNN + hyperparameter tuning using Optuna
6. CNN + transfer learning using EfficientNet
7. CNN + Optuna + EfficientNet
8. CNN + XGBoost
9. CNN+ EfficientNet + Metadata (mixed image and structured data)
10. Clustering of Meta data using K-means and K-prototypes

Logistic Regression, KNN, Random Forest

First, we tried standard algorithms- results were not promising

- Logistic regression: noisy validation results
- K-nearest neighbor: many false positives for MEL
- Random forest: $n=3$ (best accuracy), also false positives for MEL

Figure: Training history of logistic regression

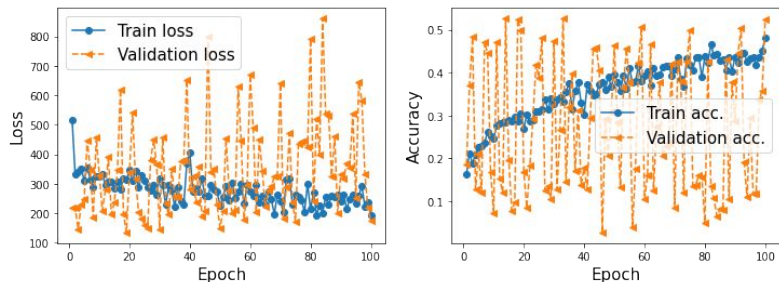


Figure: KNN accuracies across k's

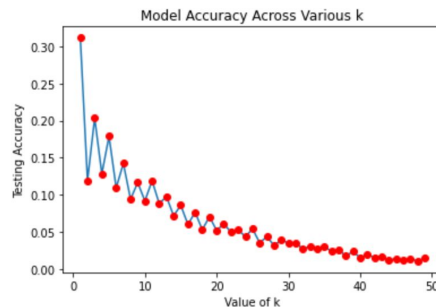
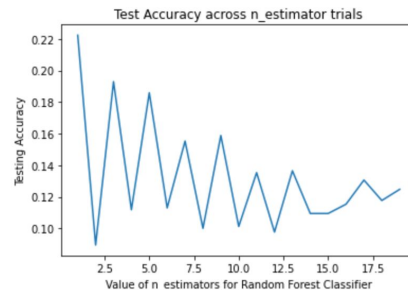


Figure: Random forests across number of trees



CNN

Brief description of approach

1. Convolutional Neural Network (CNN)
2. CNN + hyperparameter tuning using Optuna
3. CNN + transfer learning using EfficientNet
4. CNN + Optuna + EfficientNet
5. CNN + XGBoost
6. CNN+ EfficientNet + Metadata (mixed image and structured data)



Figure 5. Hyperparameter tuning of CNN structure

1. **A** sets of: (options: 2, 3)
 - Conv2D with:
 - **B** filters (options: 5 to 16)
 - **C** kernel size (options: 3 to 5)
 - **D** setting of padding/no padding (options: True/False)
 - activation: relu
 - MaxPooling
 - **E** pool size (options: 2, 3)
 - **F** stride (options: 2, 3)
 - Dropout: 0.1
2. **G** Flatten
3. **H** layers of Dense (options: 0 to 3)
 - **H** units (options: 5 to 20)
 - Activation: relu
 - Dropout: 0.1
4. One output layer of Dense
 - output units: 8
 - Activation: softmax

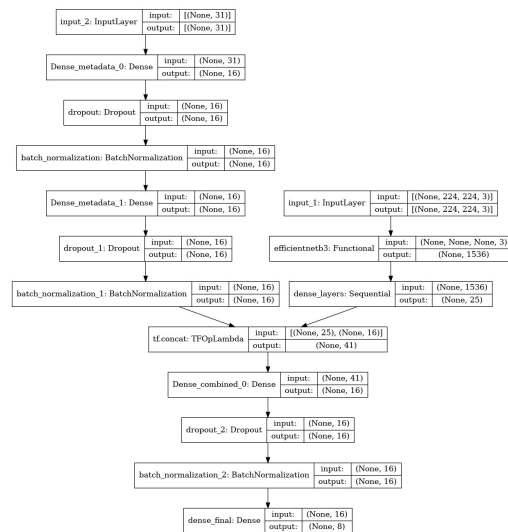
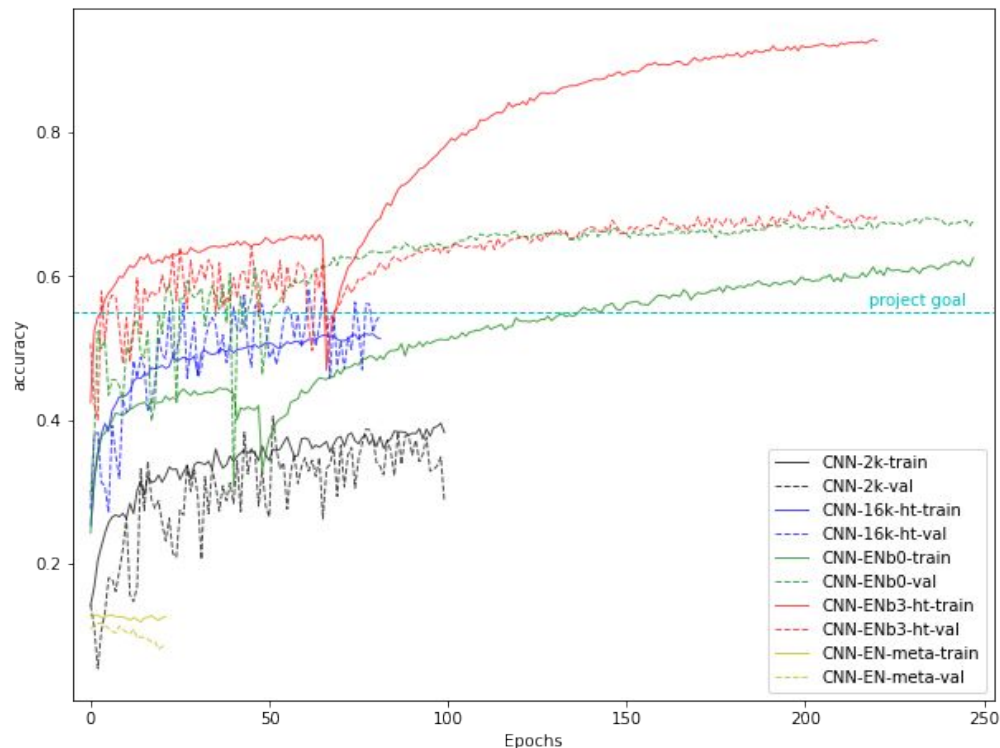


Figure: CNN + metadata input (normal dense layers; left) and image data input (CNN; right), which gets concatenated and put through a dense layer before classification.

Accuracy: top CNN models

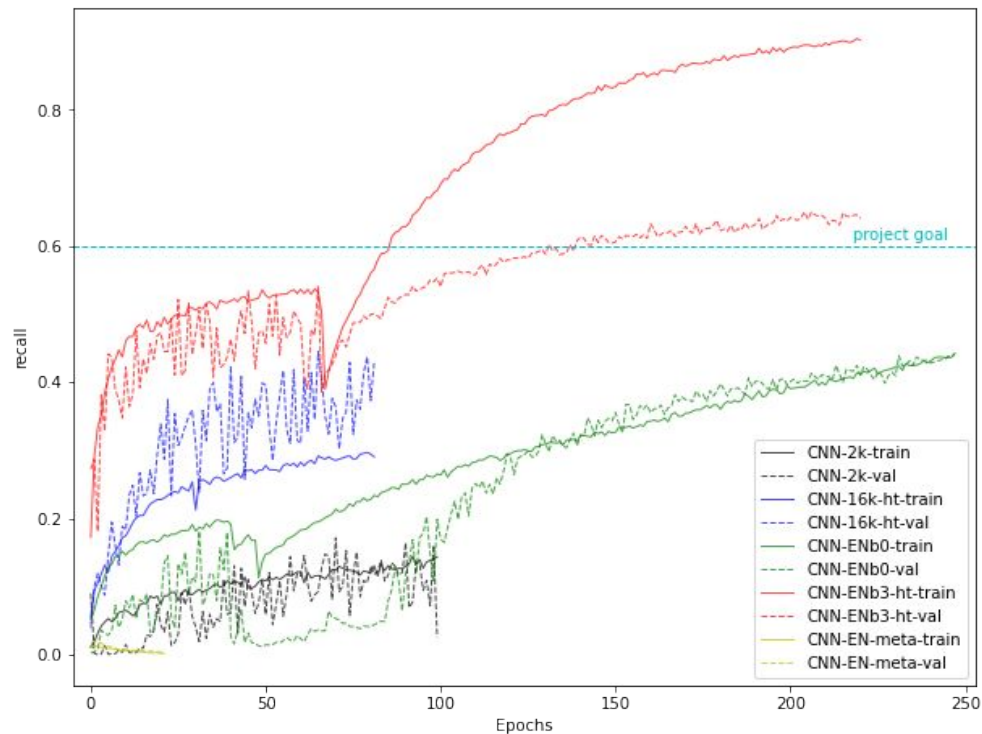


Note: these are categorical accuracies (based on model output probabilities)

Model abbreviations:

- **CNN-2k**: small dataset
- **CNN-16k-ht**: big dataset Optuna-tuned
- **CNN-ENb0**: EfficientNetB0
- **CNN-ENb3-ht**: EfficientNetB3 with Optuna hyperparameter-tuning
- **CNN-EN-meta**: EfficientNet with metadata

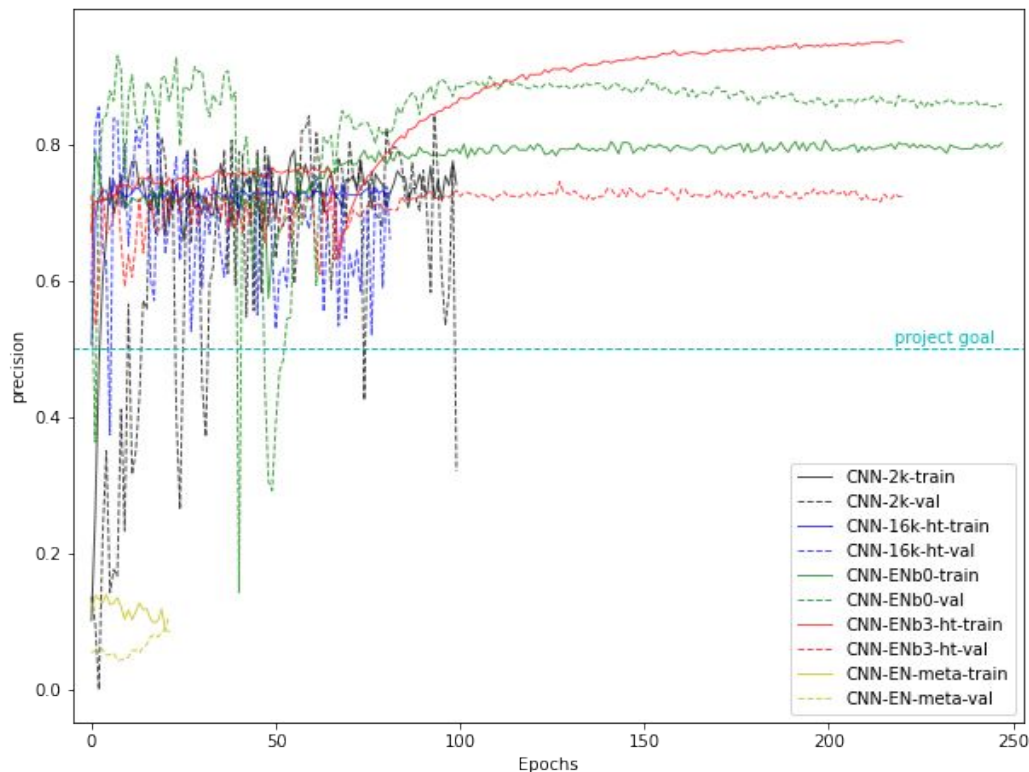
Recall: top CNN models



Model abbreviations:

- CNN-2k: small dataset
- CNN-16k-ht: big dataset Optuna-tuned
- CNN-ENb0: EfficientNetB0
- CNN-ENb3-ht: EfficientNetB3 with Optuna hyperparameter-tuning
- CNN-EN-meta: EfficientNet with metadata

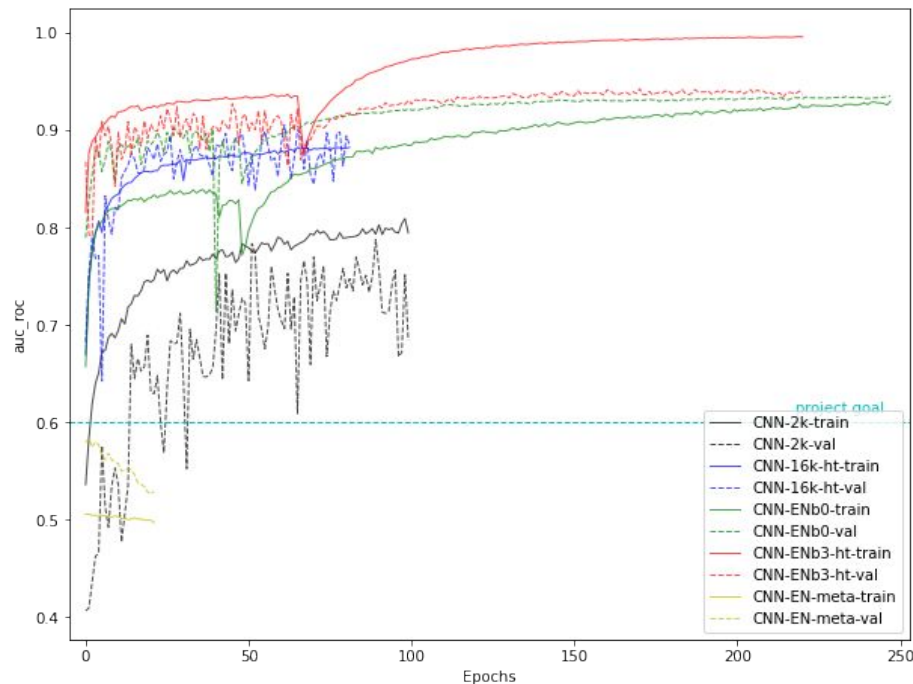
Precision: top CNN models



Model abbreviations:

- *CNN-2k*: small dataset
- *CNN-16k-ht*: big dataset Optuna-tuned
- *CNN-ENb0*: EfficientNetB0
- *CNN-ENb3-ht*: EfficientNetB3 with Optuna hyperparameter-tuning
- *CNN-EN-meta*: EfficientNet with metadata

AUC-ROC: top CNN models



Model abbreviations:

- *CNN-2k*: small dataset
- *CNN-16k-ht*: big dataset Optuna-tuned
- *CNN-ENb0*: EfficientNetB0
- *CNN-ENb3-ht*: EfficientNetB3 with Optuna hyperparameter-tuning
- *CNN-EN-meta*: EfficientNet with metadata

Metadata Clustering

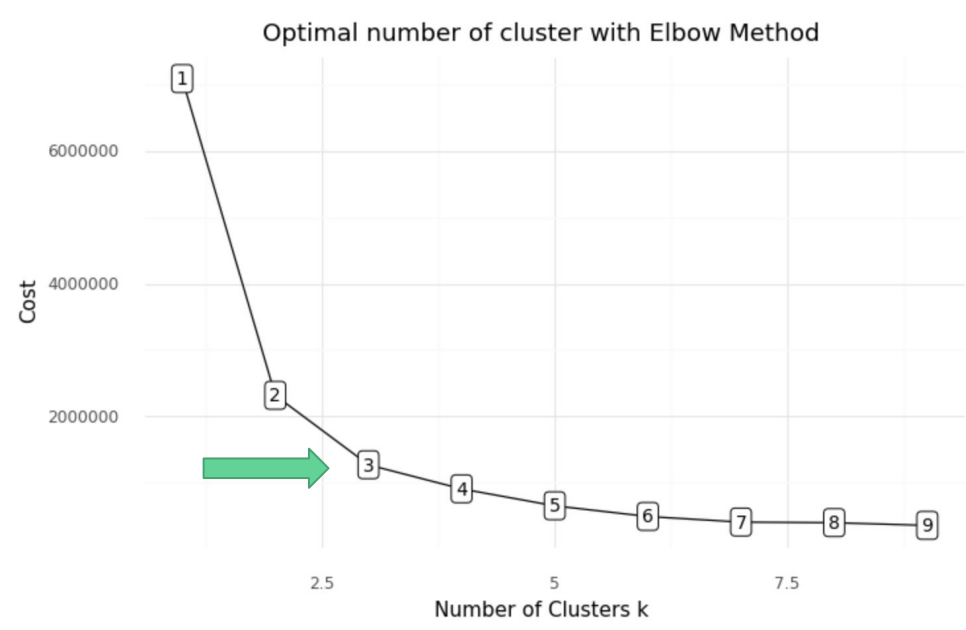
K-Means vs K-Prototypes

- Our metadata contains a mix of numerical and categorical data
- The standard K-means algorithm isn't directly applicable to categorical data
 - Categorical data is discrete
 - A Euclidean distance isn't really meaningful
- k-prototypes is a combination of k-means and k-modes.
 - Can handle both numerical and categorical features simultaneously

	image	age_approx	anatom_site_general	lesion_id	sex
0	ISIC_0000000	55.0	anterior torso	NaN	female
1	ISIC_0000001	30.0	anterior torso	NaN	female
2	ISIC_0000002	60.0	upper extremity	NaN	female
3	ISIC_0000003	30.0	upper extremity	NaN	male
4	ISIC_0000004	80.0	posterior torso	NaN	male
...
25326	ISIC_0073247	85.0	head/neck	BCN_0003925	female
25327	ISIC_0073248	65.0	anterior torso	BCN_0001819	male
25328	ISIC_0073249	70.0	lower extremity	BCN_0001085	male
25329	ISIC_0073251	55.0	palms/soles	BCN_0002083	female
25330	ISIC_0073254	50.0	upper extremity	BCN_0001079	male
25331 rows × 5 columns					

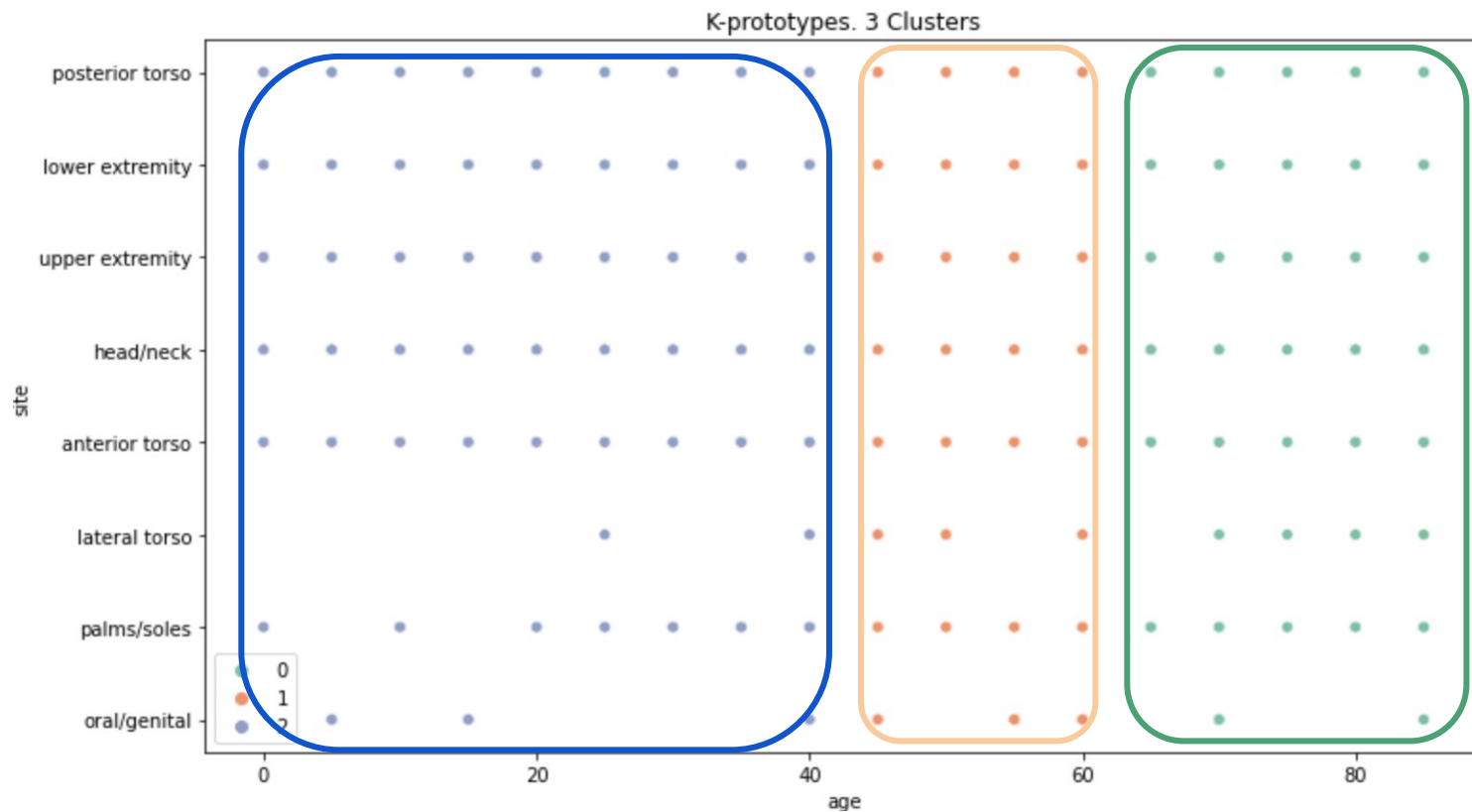
K-prototypes: Choosing K

Elbow method [1:10]

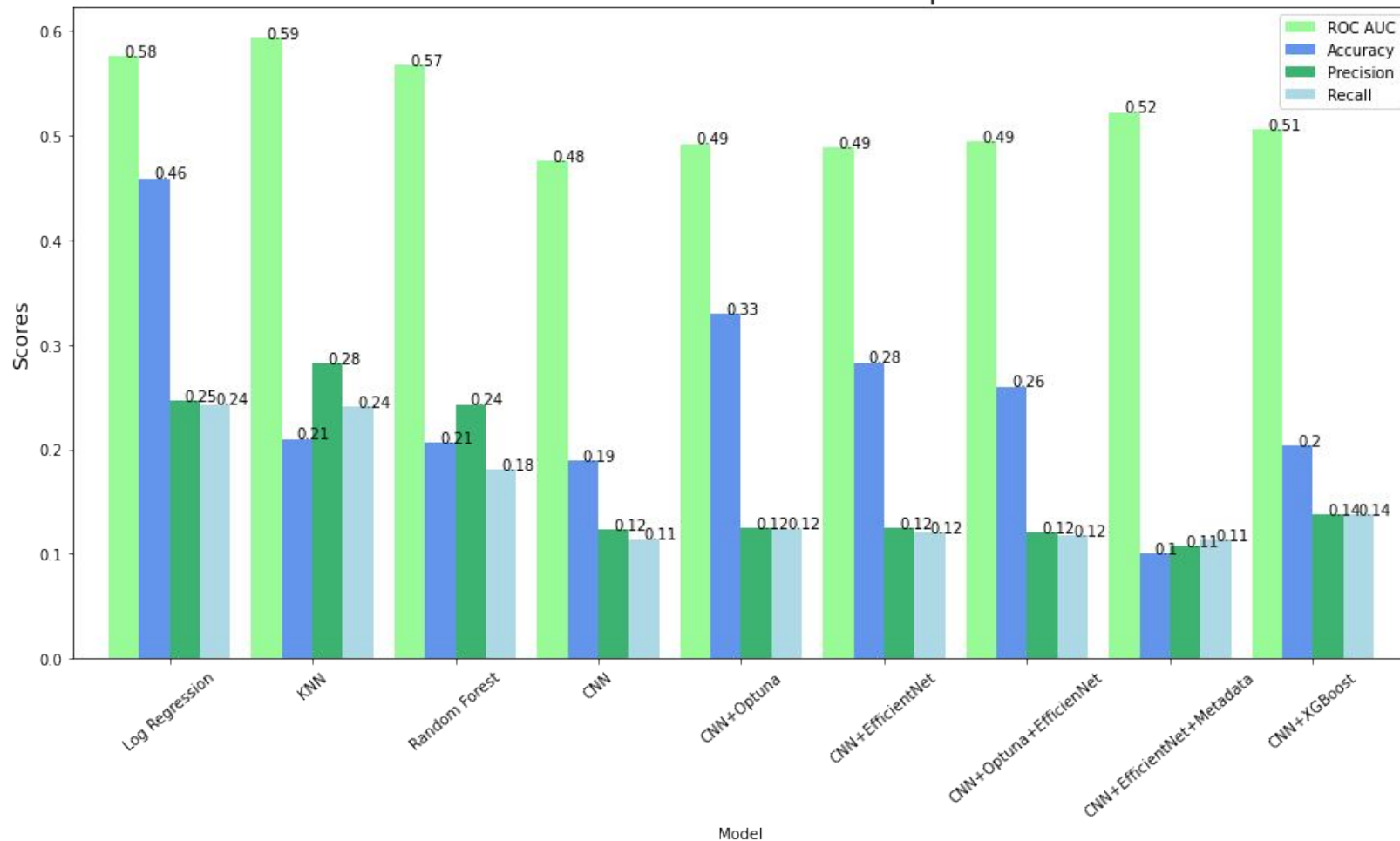


K-prototypes

K=3

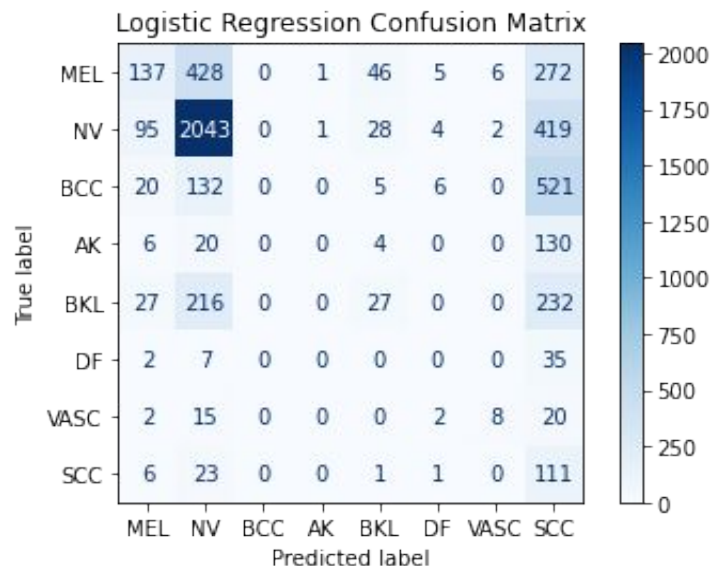


Model Performance Metrics Comparison

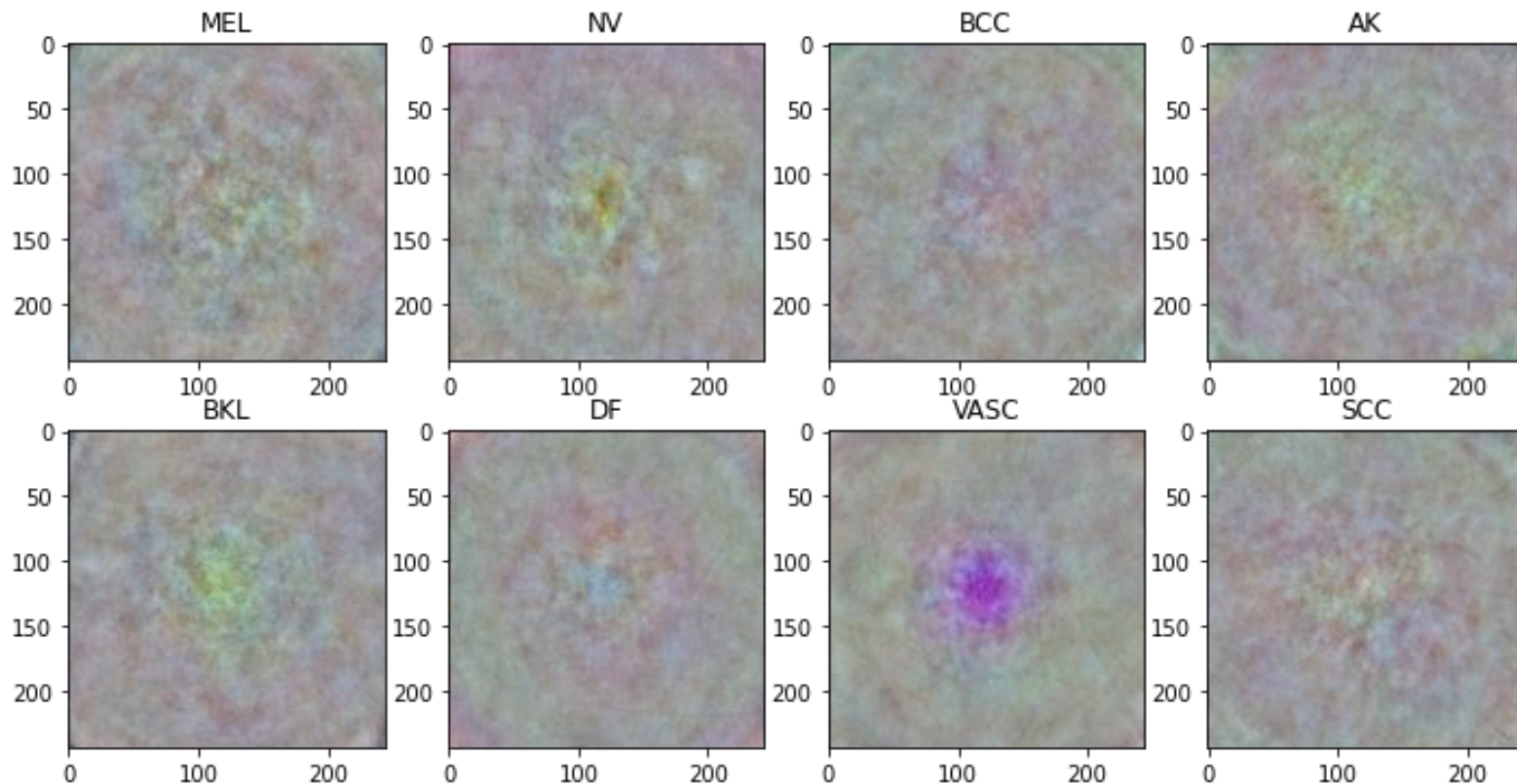


Model Conclusions

- Overall, model performance was poor
- None of our models met performance thresholds
- All models had difficulty with NV vs. MEL classification
 - CNN and Logistic Regression had best recall score for melanoma
 - CNN + tuning or transfer learning improved performance more for benign classes
- Performance not sufficient for clinical use



Logistic regression weights for each skin class



Strengths

- Real-world dataset representing real variability
- Image Augmentation implemented for class balance
- Cropping method implemented
- Optimization utilized across models
 - Optimized K in KNN and n_estimators in random forest
 - Elbow method for selecting number of clusters in K-prototypes
 - Tested performance across different EfficientNet versions

Constraints and Limitations

- Smaller training sets due to computational limitations
- Reduced generalizability for different races
- Future directions:
 - Principal component analysis
 - Further image augmentation
 - Assessing impact of presence of rulers, marks, hair, etc for certain classes



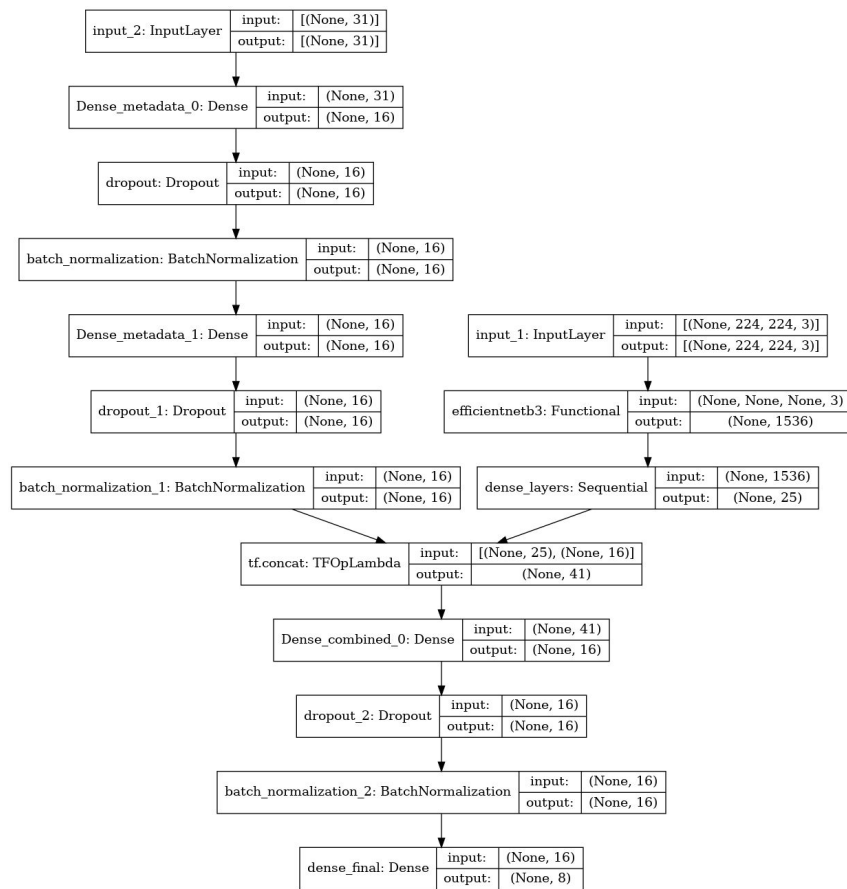
References

- [1] International Skin Imaging Collaboration (ISIC) on Kaggle, 2020: <https://www.kaggle.com/competitions/siim-isic-melanoma-classification/overview>
- [2] Cancer.net, February 2022. Melanoma: Statistics. <https://www.cancer.net/cancer-types/melanoma/statistics#:~:text=It%20is%20estimated%20that%207%2C650,people%20worldwide%20died%20from%20melanoma.>
- [3] ISIC Archive : <https://www.isic-archive.com/>
- [4] Tschandl P., Rosendahl C. & Kittler H. The HAM10000 dataset, a large collection of multi-source dermoscopic images of common pigmented skin lesions. Sci. Data 5, 180161 doi:10.1038/sdata.2018.161 (2018)
- [5] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, Allan Halpern: "Skin Lesion Analysis Toward Melanoma Detection: A Challenge at the 2017 International Symposium on Biomedical Imaging (ISBI), Hosted by the International Skin Imaging Collaboration (ISIC)", 2017; arXiv: 1710.05006.
- [6] Marc Combalia, Noel C. F. Codella, Veronica Rotemberg, Brian Helba, Veronica Vilaplana, Ofer Reiter, Allan C. Halpern, Susana Puig, Josep Malvehy: "BCN20000: Dermoscopic Lesions in the Wild", 2019; arXiv:1908.02288.
- [7] Haenssle, H. A., et al. (2018). Man against machine: diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists. Annals of oncology : official journal of the European Society for Medical Oncology, 29(8), 1836–1842. <https://doi.org/10.1093/annonc/mdy166>
- [8] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. (2019). Optuna: A Next-generation Hyperparameter Optimization Framework. In KDD. <https://optuna.org/>
- [9] Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. arXiv. <https://doi.org/10.48550/arXiv.1905.11946>
- [10] dmlc XGBoost read-the-docs website (2022). <https://xgboost.readthedocs.io/en/stable/index.html>
- [11] Ren, X., Guo, H., Li, S., Wang, S., Li, J. (2017). A Novel Image Classification Method with CNN-XGBoost Model. In: Kraetzer, C., Shi, YQ., Dittmann, J., Kim, H. (eds) Digital Forensics and Watermarking. IWDW 2017. Lecture Notes in Computer Science(), vol 10431. Springer, Cham. https://doi.org/10.1007/978-3-319-64185-0_28
- [12] Chollet, F (2020) Keras, Transfer learning & fine-tuning. https://keras.io/guides/transfer_learning/

Appendix

CNN + Metadata structure

Example structure of model using metadata input (normal dense layers; left) and image data input (CNN; right), which gets concatenated and put through a dense layer before classification.



Optuna Hyperparameter Tuning

Hyperparameter tuning of CNN structure

1. **A** sets of: (options: 2, 3)
 - Conv2D with:
 - **B** filters (options: 5 to 16)
 - **C** kernel size (options: 3 to 5)
 - **D** setting of padding/no padding (options: True/False)
 - activation: relu
 - MaxPooling
 - **E** pool size (options: 2, 3)
 - **F** stride (options: 2, 3)
 - Dropout: 0.1
2. Flatten
3. **G** layers of Dense (options: 0 to 3)
 - **H** units (options: 5 to 20)
 - Activation: relu
 - Dropout: 0.1
4. One output layer of Dense
 - output units: 8
 - Activation: softmax

Optimization History Plot

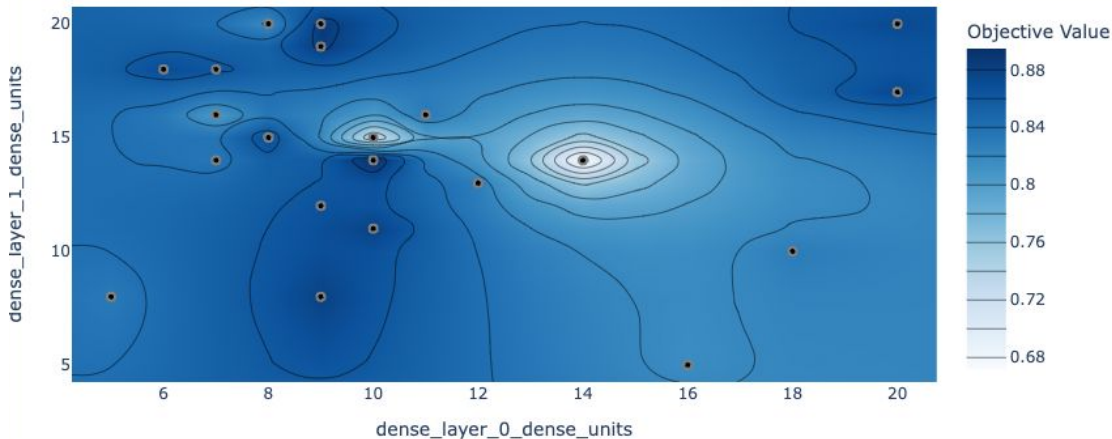
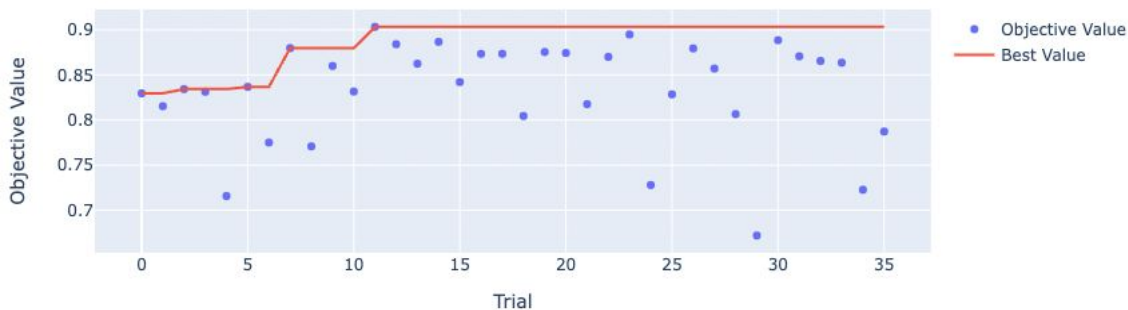
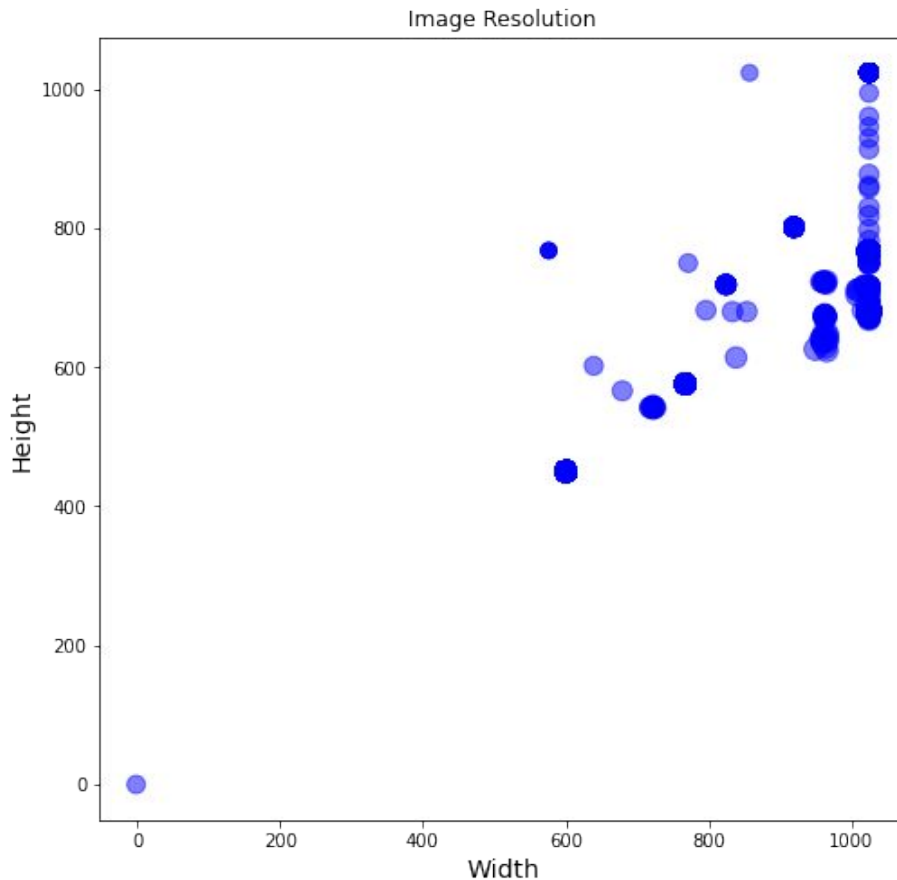


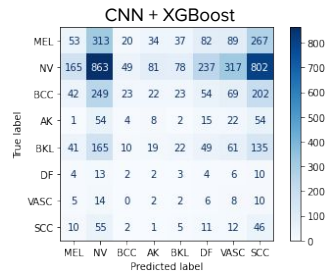
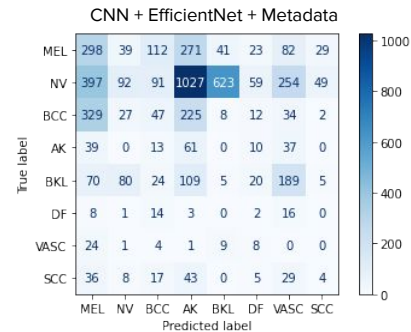
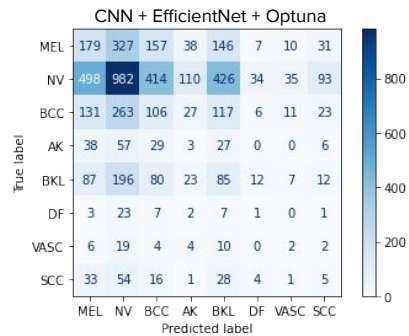
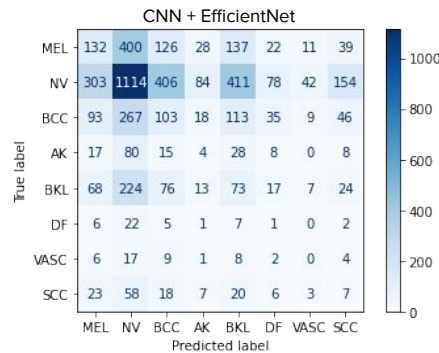
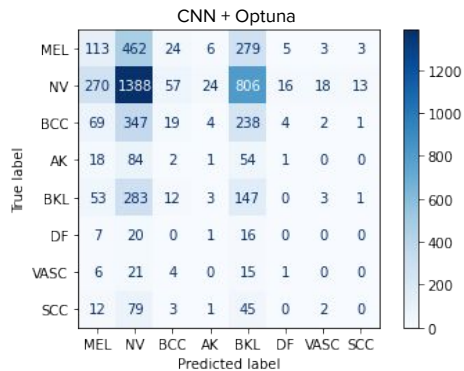
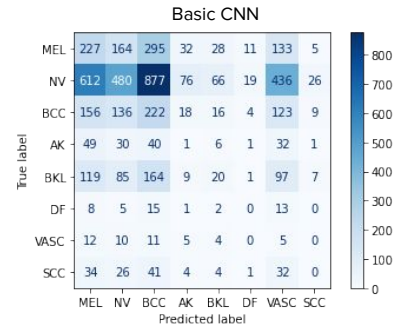
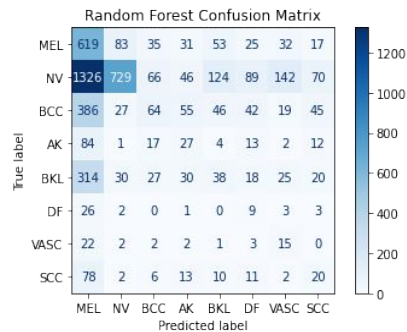
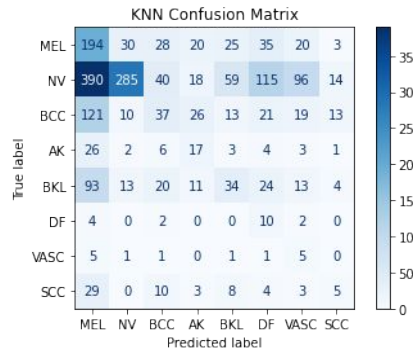
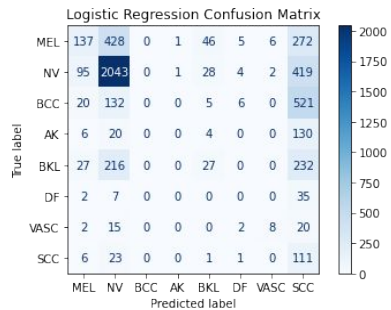
Image size distributions

Variable, but two main groups:

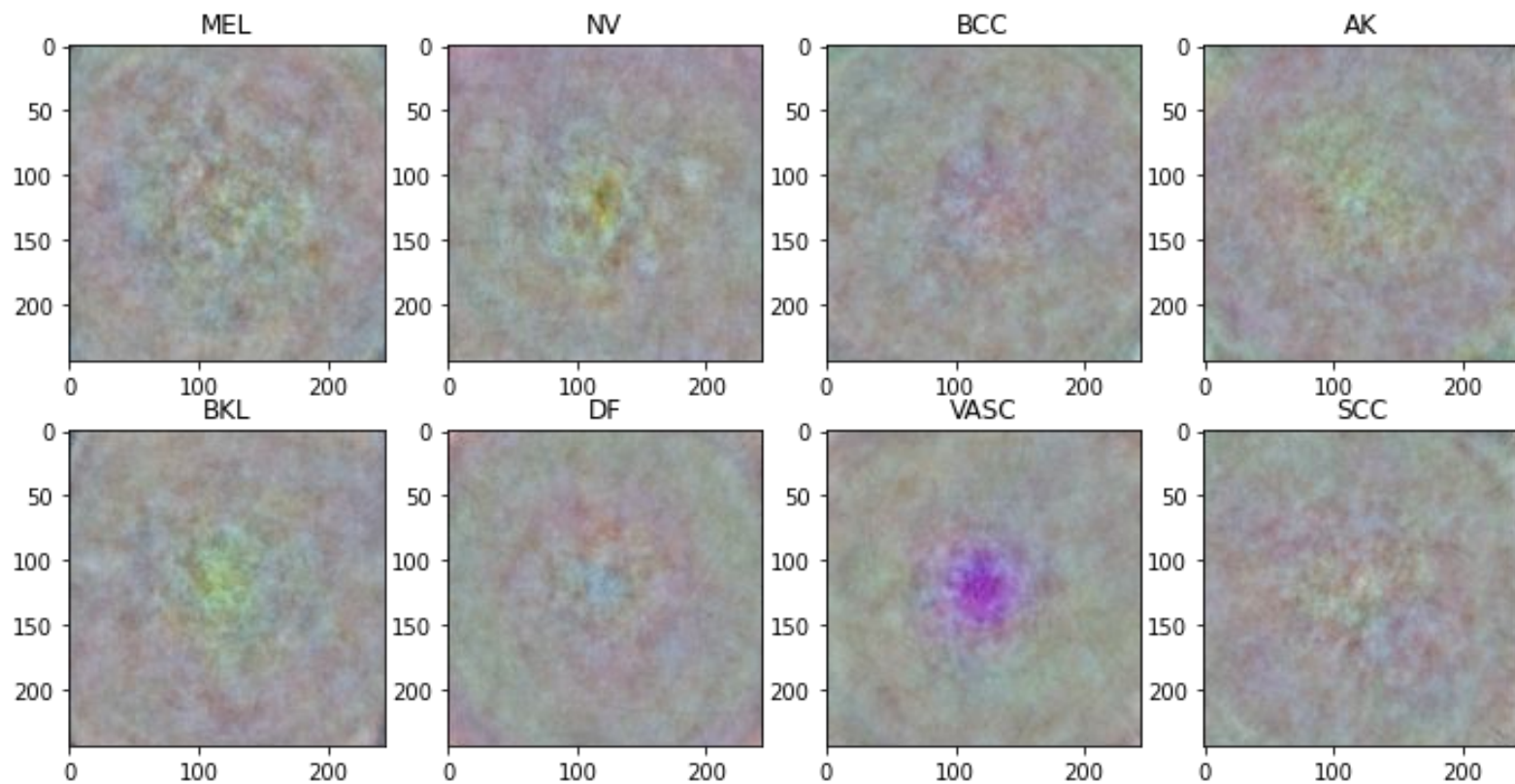
10,014 images of 600 x 450

12,410 images of 1024 x 1024





Visualizing logistic regression weights for each skin class



Logistic Regression Classification Report					KNN Classification Report					Random Forest Classification Report					CNN Classification Report				
	precision	recall	f1-score	support		precision	recall	f1-score	support		precision	recall	f1-score	support		precision	recall	f1-score	support
MEL	0.464	0.153	0.230	895	MEL	0.540	0.076	0.133	355	MEL	0.336	0.164	0.221	895	MEL	0.187	0.254	0.215	895
NV	0.708	0.788	0.746	2592	NV	0.836	0.280	0.420	1017	NV	0.832	0.281	0.420	2592	NV	0.513	0.185	0.272	2592
BCC	0.000	0.000	0.000	684	BCC	0.257	0.142	0.183	260	BCC	0.295	0.094	0.142	684	BCC	0.133	0.325	0.189	684
AK	0.000	0.000	0.000	160	AK	0.179	0.274	0.217	62	AK	0.132	0.169	0.148	160	AK	0.007	0.006	0.007	160
BKL	0.243	0.054	0.088	502	BKL	0.238	0.160	0.192	212	BKL	0.138	0.076	0.098	502	BKL	0.137	0.040	0.062	502
DF	0.000	0.000	0.000	44	DF	0.047	0.556	0.086	18	DF	0.043	0.205	0.071	44	DF	0.000	0.000	0.000	44
VASC	0.500	0.170	0.254	47	VASC	0.031	0.357	0.057	14	VASC	0.062	0.319	0.105	47	VASC	0.006	0.106	0.011	47
SCC	0.064	0.782	0.118	142	SCC	0.125	0.081	0.098	62	SCC	0.107	0.141	0.122	142	SCC	0.000	0.000	0.000	142
accuracy			0.459	5066	micro avg	0.354	0.210	0.263	2000	micro avg	0.396	0.207	0.272	5066	accuracy			0.189	5066
macro avg	0.247	0.243	0.180	5066	macro avg	0.282	0.241	0.173	2000	macro avg	0.243	0.181	0.166	5066	macro avg	0.123	0.114	0.094	5066
weighted avg	0.475	0.459	0.437	5066	weighted avg	0.590	0.210	0.292	2000	weighted avg	0.547	0.207	0.293	5066	weighted avg	0.327	0.189	0.209	5066
					samples avg	0.210	0.210	0.210	2000	samples avg	0.207	0.207	0.207	5066					
Accuracy: 0.4591					Accuracy: 0.21					Accuracy: 0.20706671930517173					Accuracy: 0.1885				
ROC AUC: 0.5760					ROC AUC: 0.5933327707205982					ROC AUC: 0.5672968987988957					ROC AUC: 0.4760				

CNN + Optuna Classification Report					CNN + EfficientNet Classification Report					CNN+Optuna+EfficientNet Classification Report				
	precision	recall	f1-score	support		precision	recall	f1-score	support		precision	recall	f1-score	support
MEL	0.206	0.126	0.157	895	MEL	0.204	0.147	0.171	895	MEL	0.171	0.187	0.179	895
NV	0.517	0.535	0.526	2592	NV	0.511	0.430	0.467	2592	NV	0.500	0.371	0.426	2592
BCC	0.157	0.028	0.047	684	BCC	0.136	0.151	0.143	684	BCC	0.118	0.140	0.128	684
AK	0.025	0.006	0.010	160	AK	0.026	0.025	0.025	160	AK	0.024	0.031	0.027	160
BKL	0.092	0.293	0.140	502	BKL	0.092	0.145	0.112	502	BKL	0.092	0.155	0.116	502
DF	0.000	0.000	0.000	44	DF	0.006	0.023	0.009	44	DF	0.000	0.000	0.000	44
VASC	0.000	0.000	0.000	47	VASC	0.000	0.000	0.000	47	VASC	0.000	0.000	0.000	47
SCC	0.000	0.000	0.000	142	SCC	0.025	0.049	0.033	142	SCC	0.052	0.063	0.057	142
accuracy			0.329	5066	accuracy			0.283	5066	accuracy			0.260	5066
macro avg	0.125	0.124	0.110	5066	macro avg	0.125	0.121	0.120	5066	macro avg	0.120	0.118	0.117	5066
weighted avg	0.332	0.329	0.317	5066	weighted avg	0.326	0.283	0.301	5066	weighted avg	0.314	0.260	0.281	5066
Accuracy: 0.3293					Accuracy: 0.2831					Accuracy: 0.2598				
ROC AUC: 0.4918					ROC AUC: 0.4891					ROC AUC: 0.4937				