

# SKIN JOB

## *Data Model Dermatology*

Name: Gong Qi Chen, Huihuang Liu

Course: DSE I2100

Instructor: Michael Grossberg



# Introduction

- Project Overview
- Research Purpose
- Why We Cares
- Data introduction
- Exploratory Data Analysis
- Methods
- Evaluation
- Conclusion

# Project Overview

In this project, we incorporated a series of machine learning techniques to help us build a model that can distinguish skin cancers from other tumors. A convolution neural network model was selected and built with the accuracy of matching human benchmark. Alongside, we also identified few key features that affect the model Training.

# Research Purpose

- Objective:
  - Use image data to classify skin lesions
  - Emphasize on recall and precision rate of skin cancers
- Benchmark:
  - Human accuracy: 67-75%
- Our goal
  - Match or better human benchmark
  - Achieve higher F1 score for skin cancers



# Why we care?

Limited screening methods

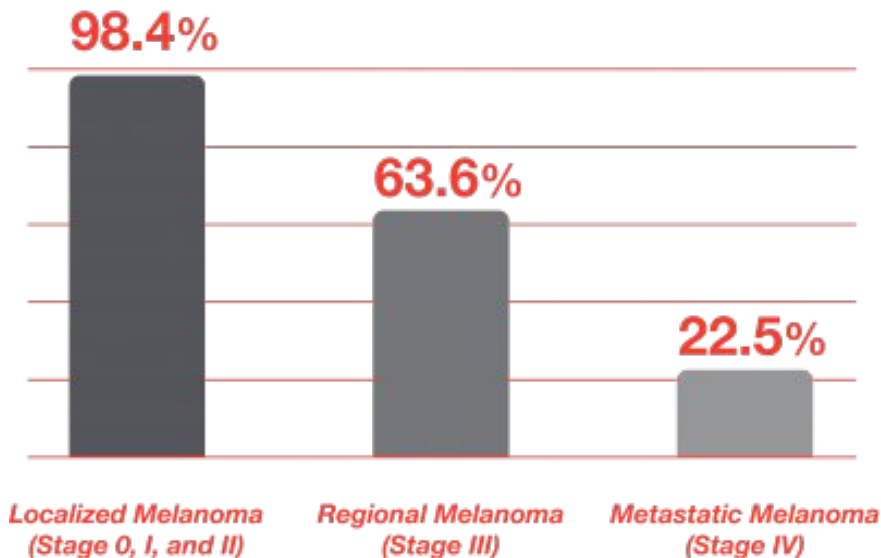
Once malignancy is confirmed, usually late stage

Accurate skin cancer detection at earlier stage is the key:

1. Improved survival
2. Improved clinical outcomes
3. Improved quality of life

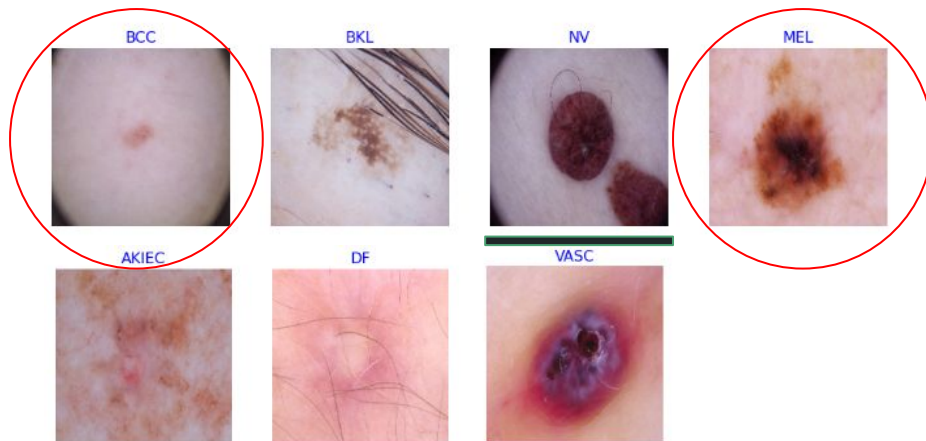


**Five-Year Survival Rate by Melanoma Stage**



# Data Introduction

- Dataset: HAM10000: Human Against Machine with 10000 training images
- What:
  - 10015 dermoscopic images
  - Types of label: 5 benign + 2 malignant
    - Benign: (AKIEC, BKL, DF,NV, VASC), **Malignant: (BCC, MEL)**
    - NV most common label
  - Variables:
    - lesion\_id, image\_id, dx, dx\_type, age and localization



	lesion_id	image_id	dx	dx_type	age	sex	localization
0	HAM_0000118	ISIC_0027419	bkl	histo	80.0	male	scalp
1	HAM_0000118	ISIC_0025030	bkl	histo	80.0	male	scalp
2	HAM_0002730	ISIC_0026769	bkl	histo	80.0	male	scalp
3	HAM_0002730	ISIC_0025661	bkl	histo	80.0	male	scalp
4	HAM_0001466	ISIC_0031633	bkl	histo	75.0	male	ear

## Data Introduction

- Who:
  - International Skin Imaging Collaboration (ISIC) archive
  - Standard source for dermoscopic image analysis research
- When: 2018 challenges
- Why: Enhance the diagnostic accuracy for distinguishing melanoma from other tumors

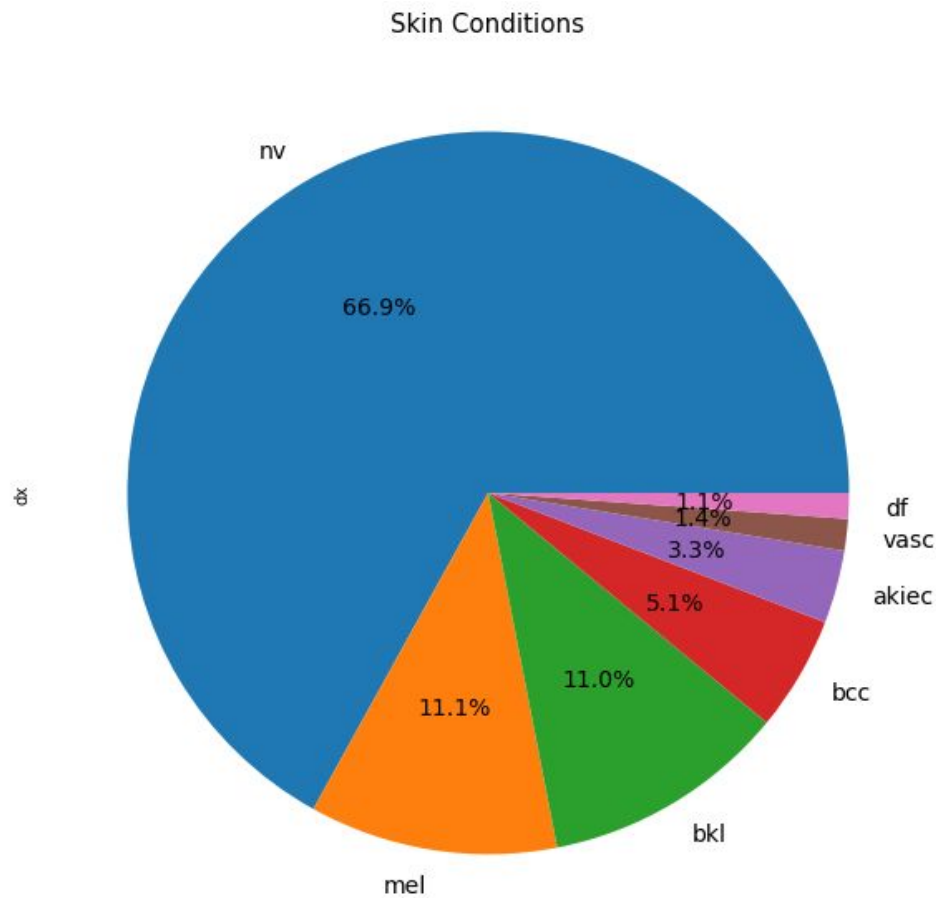
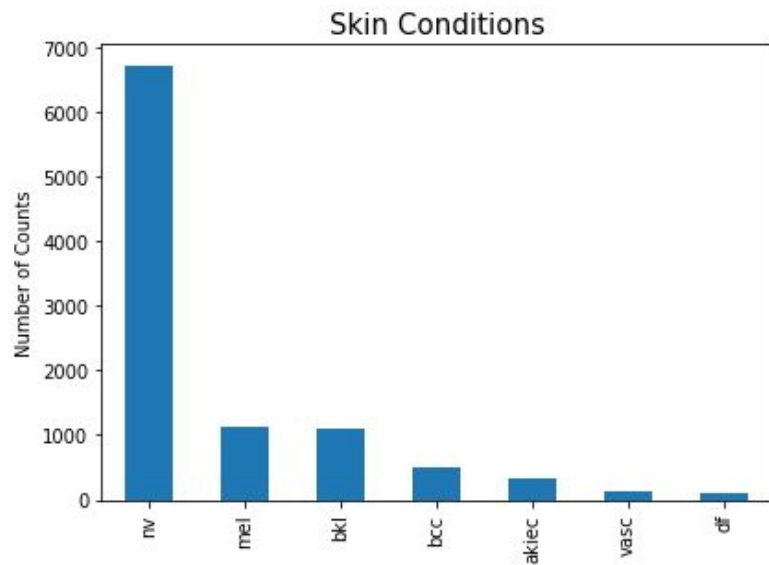
# Exploratory Data Analysis

- Missing Data in Age
  - Replaced with average mean

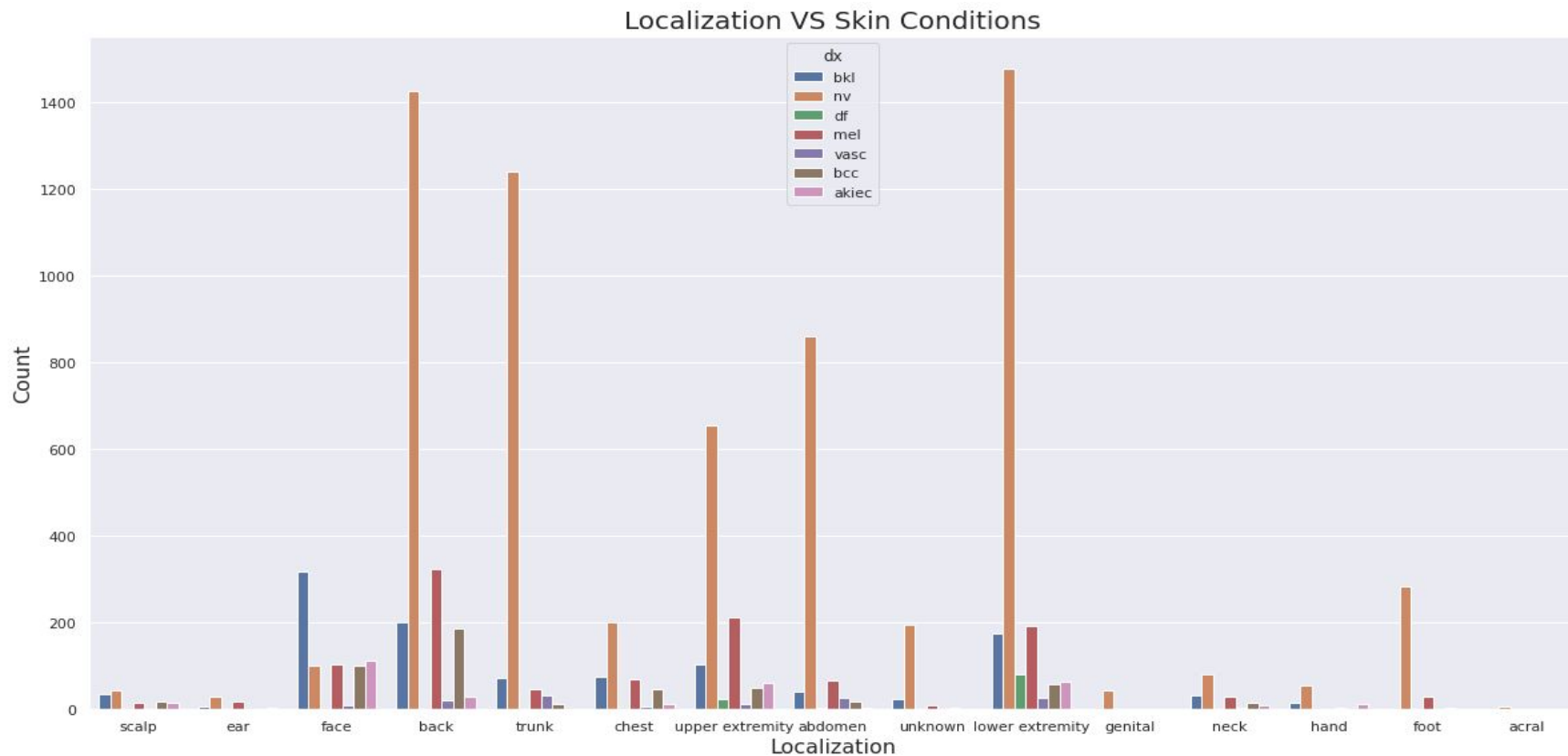
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10015 entries, 0 to 10014
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   lesion_id       10015 non-null  object
1   image_id        10015 non-null  object
2   dx              10015 non-null  object
3   dx_type         10015 non-null  object
4   age             9958 non-null   float64
5   sex             10015 non-null  object
6   localization    10015 non-null  object
dtypes: float64(1), object(6)
memory usage: 547.8+ KB
```



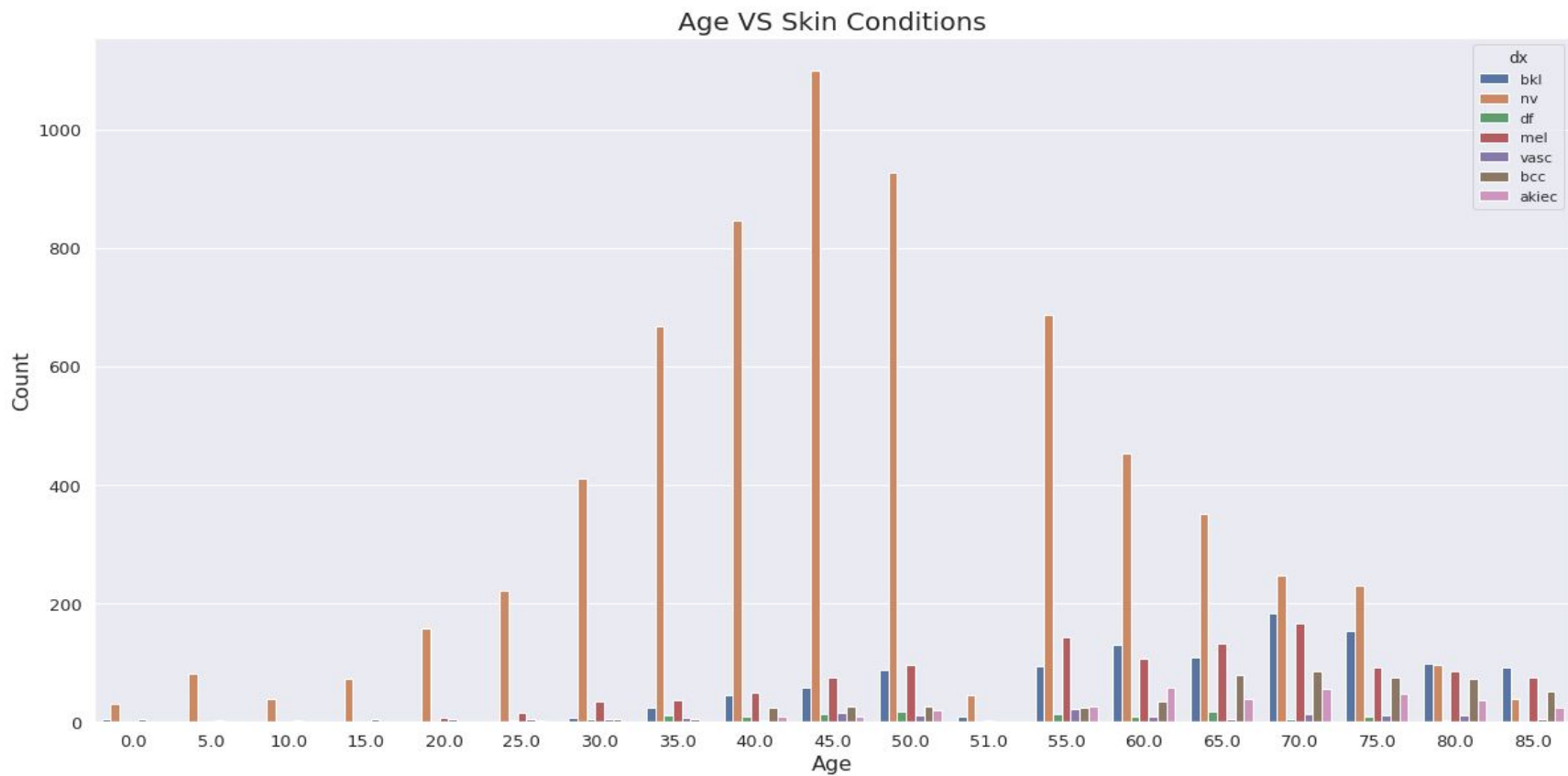
# EDA - Imbalance Data



# EDA - Bivariate Analysis



# EDA - Bivariate Analysis



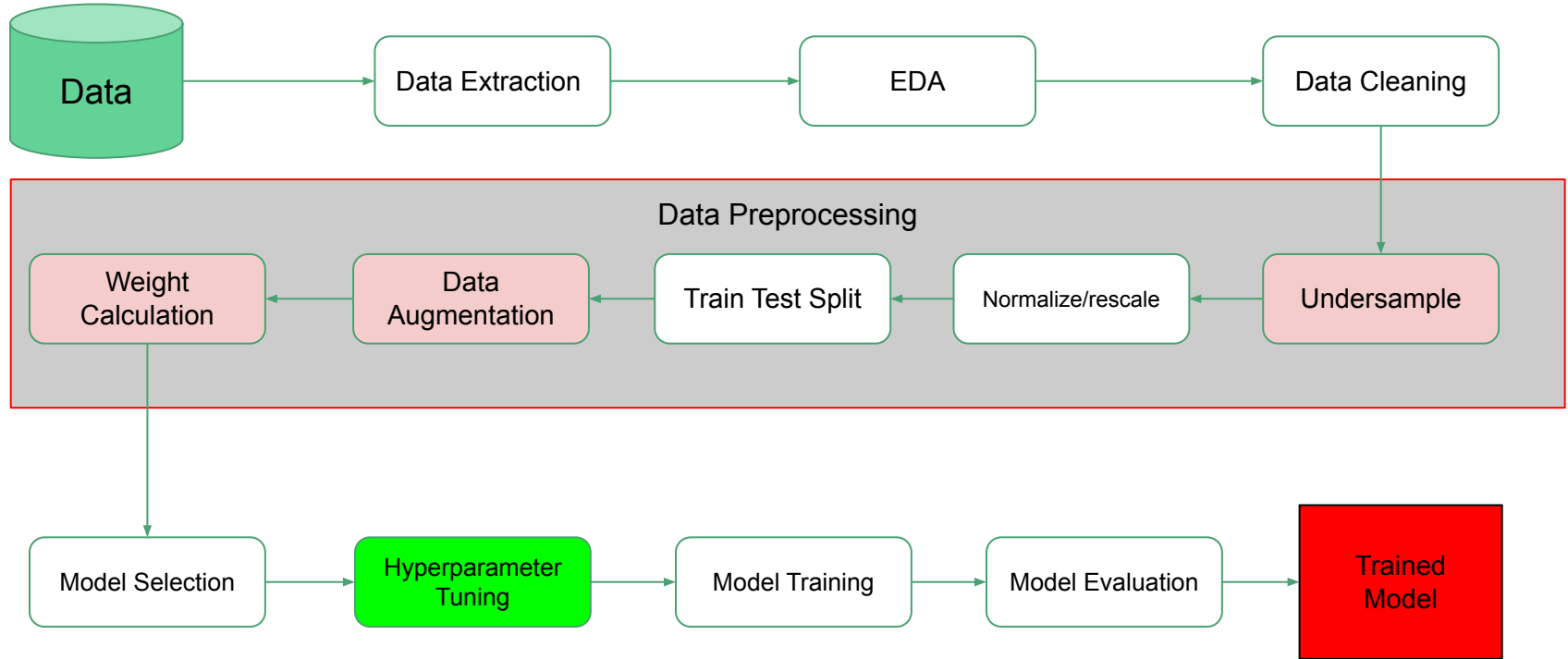
# Feature Selection

- Image Data
  - 90x90
- Age
- Localization

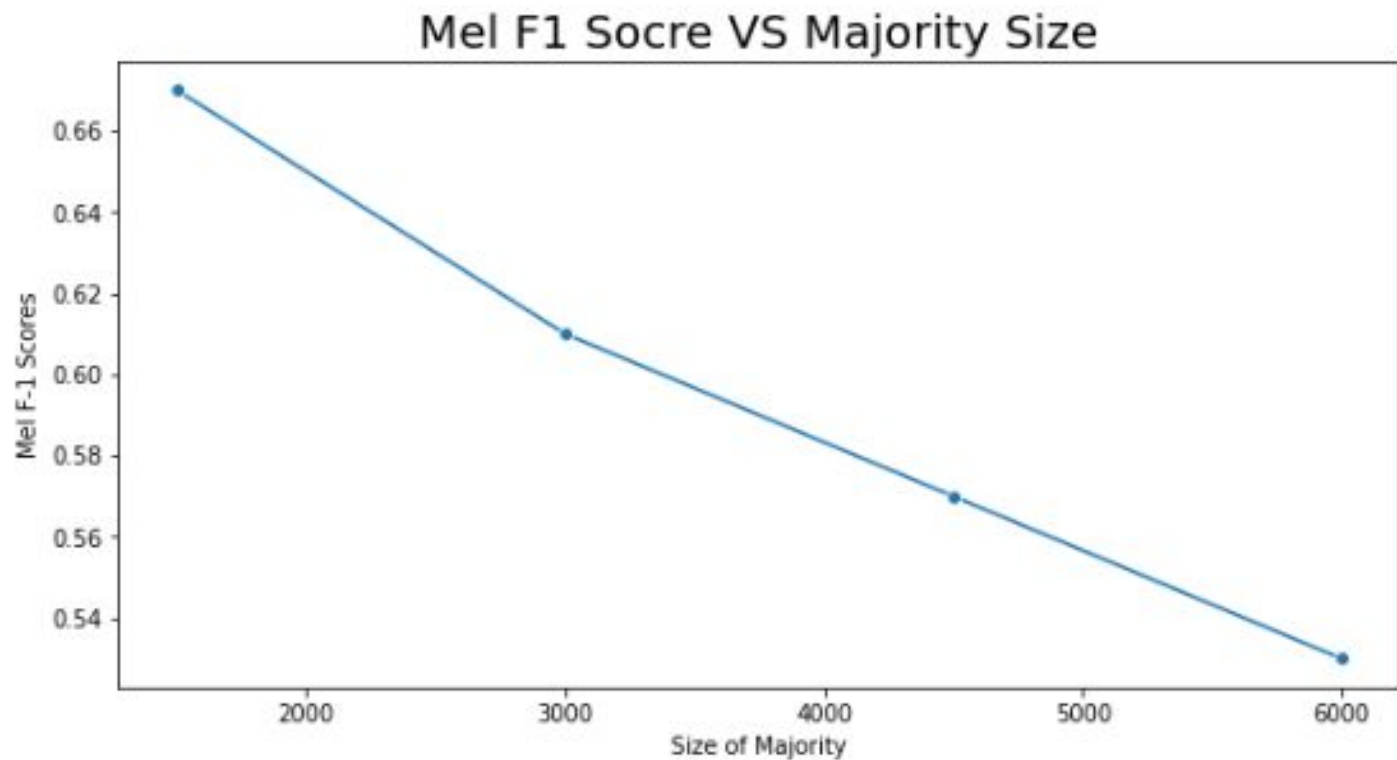
## ***Model Improvement with Age and Localization:***

Models	Val accuracy	Mel F1 Score	Bcc F1 Score
Ridge Classifier	9.30%	0.11	0.1
Logistic Classifier	11.81%	0.16	0.04
SVC	9.02%	0.05	0.12
Random Forest	7.51%	0.04	0.05
DNN	3.65%	0.15	0.03

# Pipeline Flow Chart



# Undersample



# Model Selection

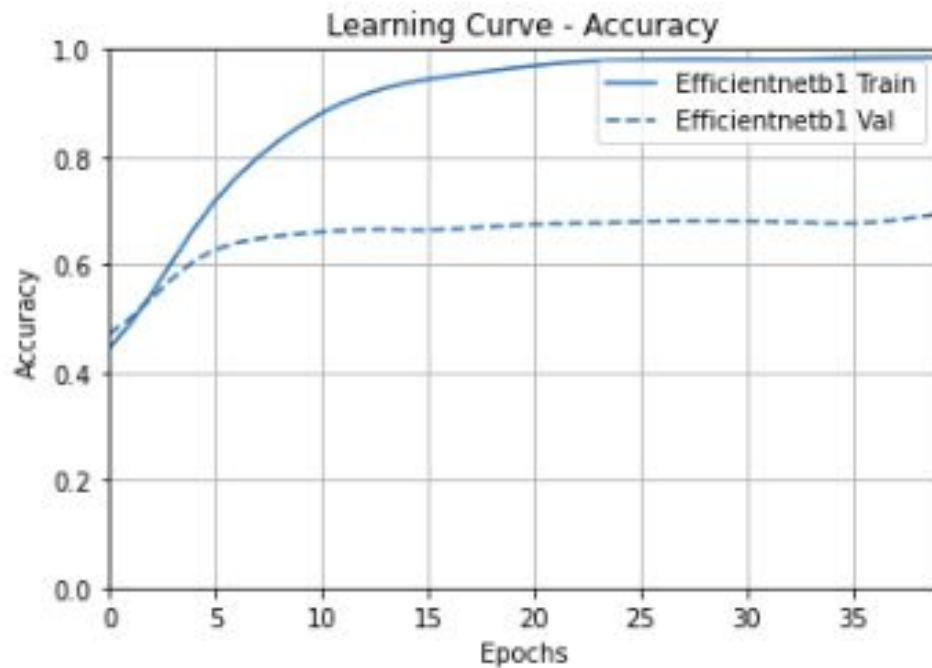
- Sklearn Models
  - Ridge Classifier
  - Logistic Classifier
  - SVC
  - Random Forest
- Keras Models
  - DNN (Dense Neural Network)
  - CNN (Convolutional Neural Network)
  - CNN with EfficientNetB1

# Model Selection

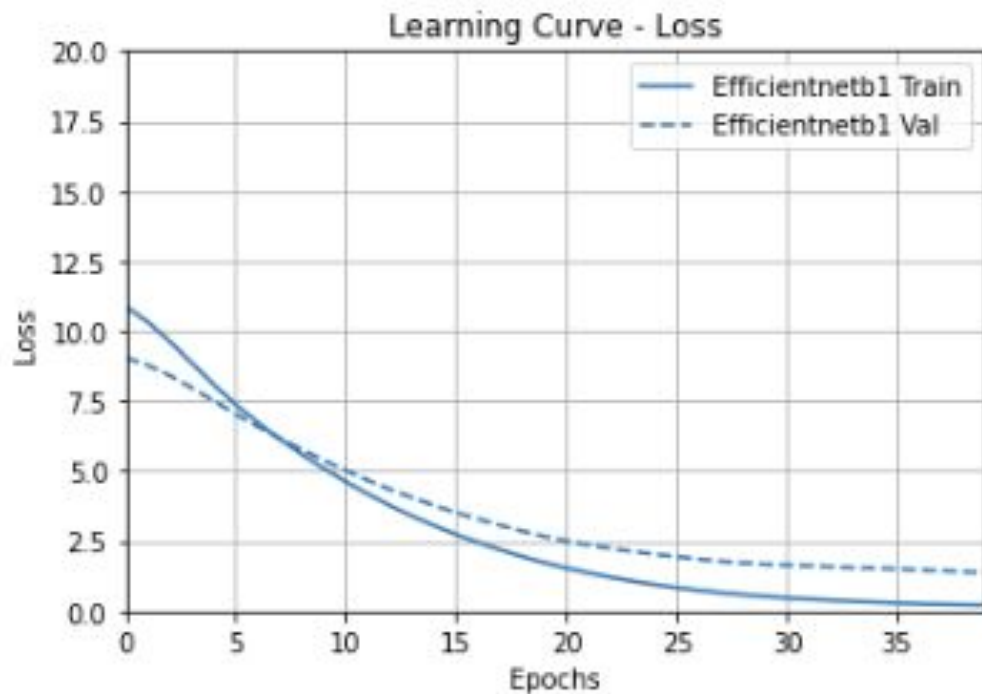
Models	Val accuracy	Mel F1 Score	Bcc F1 Score	Train - Val
Ridge Classifier	62.06%	0.45	0.33	8.14%
Logistic Classifier	67.57%	0.50	0.43	11.36%
SVC	72.72%	0.49	0.55	16.88%
Random Forest	68.93%	0.48	0.35	11.57%
DNN	71.65%	0.42	0.34	2.39%
CNN	75.23%	0.45	0.48	7.59%
CNN (EfficientNet B1)	81.53%	0.55	0.62	15.42%



# Learning Curve - Accuracy



# Learning Curve - Loss



# Confusion Matrix

Confusion Matrix

Actual \ Predicted	akiec	bcc	bkl	df	mel	nv	vasc
akiec	33	6	5	1	3	3	0
bcc	10	57	4	1	4	2	1
bkl	10	9	120	0	19	19	0
df	0	3	0	9	1	3	0
mel	3	3	28	0	97	21	1
nv	5	3	17	0	21	177	0
vasc	0	3	1	0	0	0	19

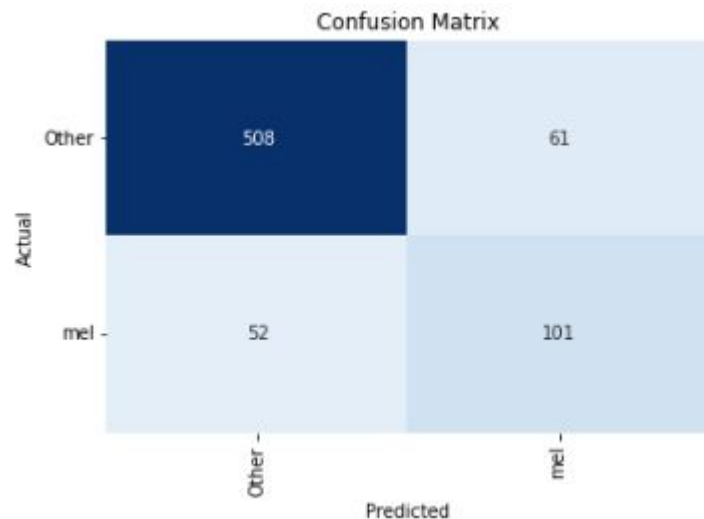
- Mel: 97 out of 153
- Bcc: 57 out of 79

# Classification Report

Classification Report:

	precision	recall	f1-score	support
akiec	0.54	0.65	0.59	51
bcc	0.68	0.72	0.70	79
bkl	0.69	0.68	0.68	177
df	0.82	0.56	0.67	16
mel	0.67	0.63	0.65	153
nv	0.79	0.79	0.79	223
vasc	0.90	0.83	0.86	23
accuracy			0.71	722
macro avg	0.73	0.69	0.71	722
weighted avg	0.71	0.71	0.71	722

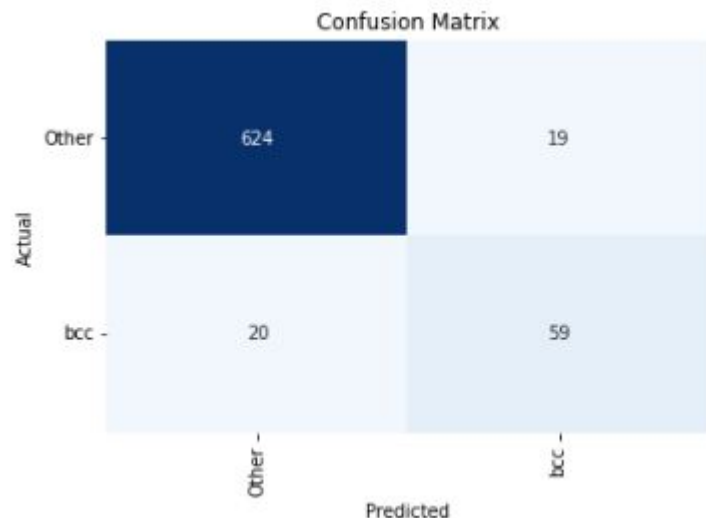
# One vs Rest - Mel



## Classification Report:

	precision	recall	f1-score	support
other	0.91	0.89	0.90	569
mel	0.62	0.66	0.64	153
accuracy			0.84	722
macro avg	0.77	0.78	0.77	722
weighted avg	0.85	0.84	0.85	722

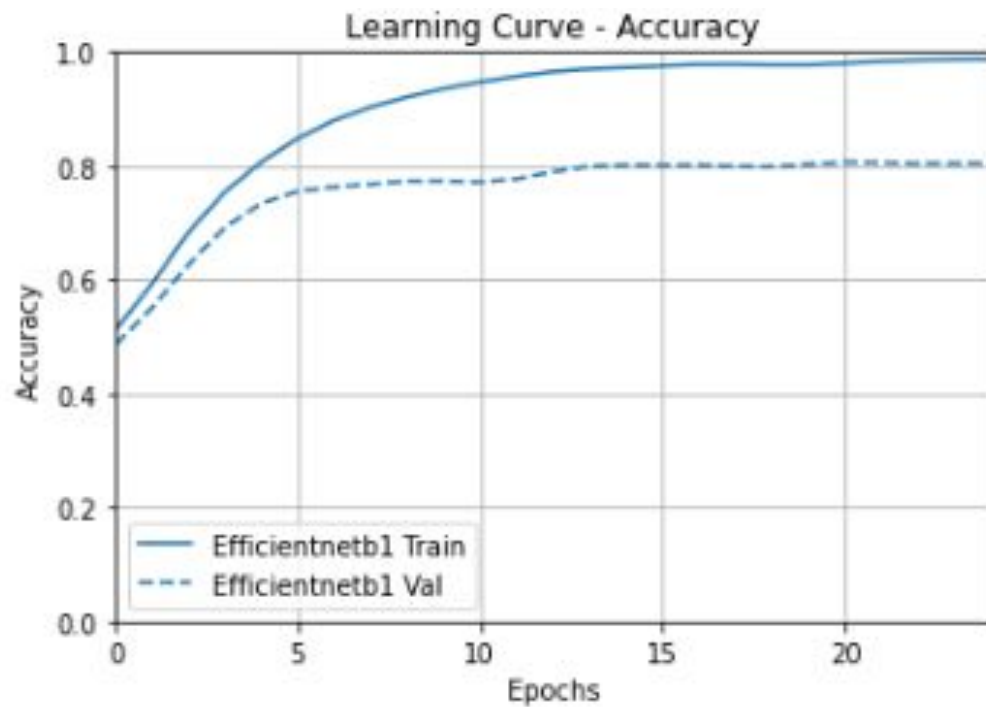
# One vs Rest - Bcc



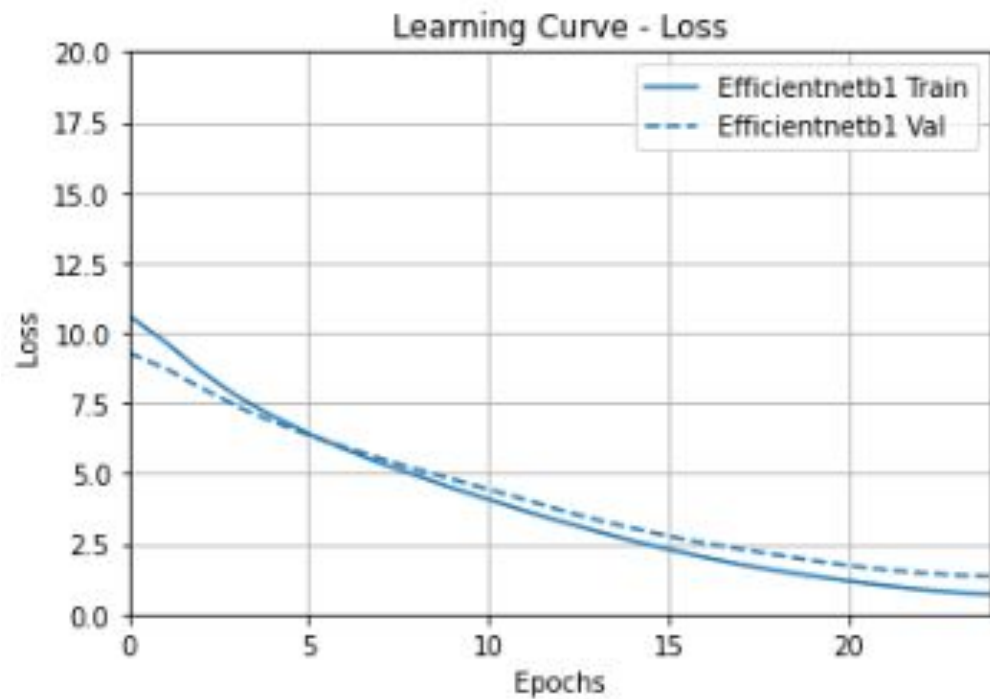
## Classification Report:

	precision	recall	f1-score	support
other	0.97	0.97	0.97	643
bcc	0.76	0.75	0.75	79
accuracy			0.95	722
macro avg	0.86	0.86	0.86	722
weighted avg	0.95	0.95	0.95	722

## On High Resolution



# On High Resolution





# On High Resolution

Confusion Matrix

Actual \ Predicted	akiec	bcc	bkl	df	mel	nv	vasc
akiec	40	2	3	1	3	2	0
bcc	5	62	4	3	4	1	0
bkl	8	7	134	1	16	10	1
df	1	2	0	12	0	1	0
mel	4	2	16	1	112	17	1
nv	2	7	8	1	17	188	0
vasc	0	1	1	0	1	1	19

- Mel: 122 out of 153
- Bcc: 62 out of 79

# On High Resolution

## Classification Report:

	precision	recall	f1-score	support
akiec	0.67	0.78	0.72	51
bcc	0.75	0.78	0.77	79
bk1	0.81	0.76	0.78	177
df	0.63	0.75	0.69	16
mel	0.73	0.73	0.73	153
nv	0.85	0.84	0.85	223
vasc	0.90	0.83	0.86	23
accuracy			0.79	722
macro avg	0.76	0.78	0.77	722
weighted avg	0.79	0.79	0.79	722

# Conclusion

- Model:
  - Data preprocessing technique is effective
  - Our model is matching with human benchmark
- Findings:
  - Age and localization are very effective in model training
  - One vs Rest is effective for Bcc, but not Mel
- Limitations:
  - Dataset suffer from biases: light skins
  - More risk factor features will improve model accuracy: family history, pre-existing condition

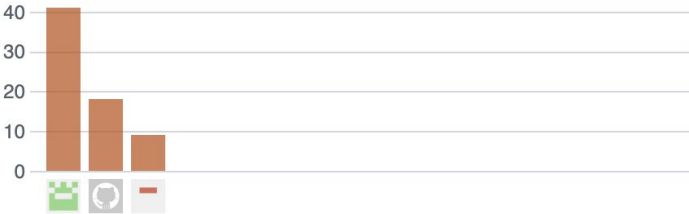
# Attribute

April 17, 2022 – May 17, 2022

Period: 1 month ▾












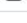




Overview			
<div><div></div></div> <div>23 Active pull requests</div>		<div><div></div></div> <div>0 Active issues</div>	
<div><div></div>23</div> <div>Merged pull requests</div>	<div><div></div>0</div> <div>Open pull requests</div>	<div><div></div>0</div> <div>Closed issues</div>	<div><div></div>0</div> <div>New issues</div>









Excluding merges, **3 authors** have pushed **68 commits** to main and **68 commits** to all branches. On main, **0 files** have changed and there have been **0 additions** and **0 deletions**.



23 Pull requests merged by 2 people

# Notebook List








 albert6051 Add files via upload		
..		
 1.0-gqc-initial-EDA.ipynb	Move Notebook to Individual folders	
 2.0-gqc-Linear_and_Logistic_Model_test.ipynb	Move Notebook to Individual folders	
 2.1-gqc-Linear_and_Logistic_Model_Sklearn.ipynb	Update to Sklearn Model	
 3.0-gqc-DNN_and_CNN_Model_test.ipynb	Move Notebook to Individual folders	
 3.1-gqc-CNN_with_regularization.ipynb	Add files via upload	
 3.2-gqc-DNN.ipynb	Add files via upload	
 4.0-gqc-balanced_image_numpy_convertor.ipynb	Move Notebook to Individual folders	
 4.0-gqc-image_numpy_convertor.ipynb	Move Notebook to Individual folders	
 5.0-gqc-pretrained_models_MobileNetV2.ipynb	Move Notebook to Individual folders	
 5.1-gqc-pretrained_models_EfficientNetB1.ipynb	Add files via upload	
 5.2-gqc-pretrained_models_EfficientNetB1_Data_Augme...	Move Notebook to Individual folders	
 7.0-gpc-Resample_for_Balancing_Data.ipynb	Move Notebook to Individual folders	
 8.0-gqc-preprocessing_pipeline.ipynb	Add files via upload	
 9.0-gqc-Sklearn_Models.ipynb	Sklearn Models test	
 9.1-gqc-Sklearn_Models_pca.ipynb	PCA model test	

HuiHuang Liu and HuiHuang Liu sklearn model update with baggingclassifier		
..		
 6.0 SVM.ipynb	Move Notebook to Individual folders	
 Initial_EDA+Descriptions.ipynb	EDA with description	
 Linear+Logistic_model_with_keras.ipynb	update on the previous version	
 SVM_tf_Update.ipynb	updatesvm model with multicategorical classification	
 keras_initial.ipynb	keras models update	
 keras_models.ipynb	keras models update	
 sklearn_models.ipynb	sklearn models update	
 sklearn_update_with_BaggingClassifier.ipynb	sklearn model update with baggingclassifier	

# Reports

HuiHuang Liu and HuiHuang Liu update the version

..

 Reports before merging	merging everything to final_report and created new folder
 figures	Setup Cookie Cutter
 .gitkeep	Setup Cookie Cutter
 Final_Report.ipynb	update the version
 Report_status.ipynb	updated data section with description of each categories of tumor
 Status Report#1.pptx	Status Report 5/3
 research interest.docx	Setup Cookie Cutter

# Reference

1. Rosendahl, C., Tschandl, P., Cameron, A. & Kittler, H. Diagnostic accuracy of dermoscopy for melanocytic and nonmelanocytic pigmented lesions. *J Am Acad Dermatol* 64, 1068–1073 (2011).
2. Bechelli S, Delhommelle J. Machine Learning and Deep Learning Algorithms for Skin Cancer Classification from Dermoscopic Images. *Bioengineering* (Basel). 2022;9(3):97. Published 2022 Feb 27.
3. Binder, M. et al. Application of an artificial neural network in epiluminescence microscopy pattern analysis of pigmented skin lesions: a pilot study. *Br J Dermatol* 130, 460–465 (1994).
4. Codella, N. C. F. et al. 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). Preprint at <https://arxiv.org/abs/1710.05006> (2017).
5. Deng, J. et al. ImageNet: A large-scale hierarchical image database, 2009 IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, 2009, pp. 248–255 (2009).
6. 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 (2018).
7. Dreiseitl, S., Binder, M., Hable, K. & Kittler, H. Computer versus human diagnosis of melanoma: evaluation of the feasibility of an automated diagnostic system in a prospective clinical trial. *Melanoma Res* 19, 180–184 (2009).
8. Kharazmi, P., Kalia, S., Lui, H., Wang, Z. J. & Lee, T. K. A feature fusion system for basal cell carcinoma detection through data-driven feature learning and patient profile. *Skin Res Technol* 24, 256–264 (2017).
9. Sinz, C. et al. Accuracy of dermoscopy for the diagnosis of nonpigmented cancers of the skin. *J Am Acad Dermatol* 77, 1100–1109 (2017).
10. Esteva, A. et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* 542, 115–118 (2017).
11. Han, S. S. et al. Classification of the clinical images for benign and malignant cutaneous tumors using a deep learning algorithm. *J Invest Dermatol*, Preprint at <https://doi.org/10.1016/j.jid.2018.01.028> (2018).