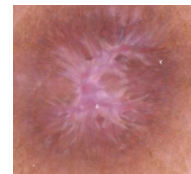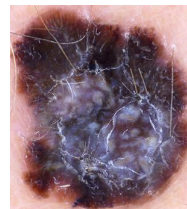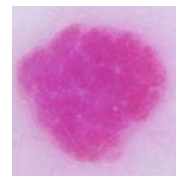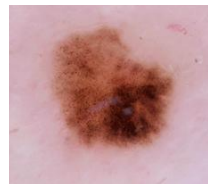# Capstone 3:
# Image Classification

Jarom Hatch

# Introduction

# Problem Statement

- Machine learning is evolving and more and more becoming part of our everyday lives
  - Predicting stock prices
  - Giving tv show and product suggestions
  - Helping self driving cars identify people, vehicles, and obstacles
- Machine learning can and is also used in the medical world
- Can we use images, as well as some specific metadata associated with those images, to accurately diagnose skin lesions?
  - Beneficial for both patients and doctors to speed up diagnoses and treatment plans
  - Potentially easier and less expensive than biopsies
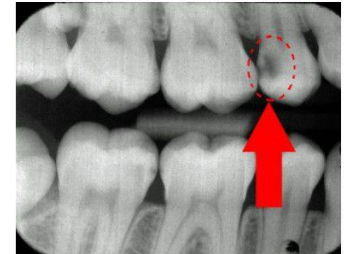
# Target Audience

- Dermatologists and Patients with increased risk of skin disease
    - Although, the benefits from this learning can be applied to all aspects of medicine where visual identification of issues is present
- Secondary audience:
    - Physicians as a whole, Pharmaceuticals

Chickenpox rash

Smallpox rash

# Data Overview

# Initial Datasets

- **2** main types provided:
  - Images
  - Metadata
- **Images**:
  - 10,015 color images comprised of 7 different lesions
  - 450 x 600 pixels
- **Metadata**:
  - Lesion ID, Image ID, Diagnosis, Diagnosis method, age, sex, and location of lesions
    - 7 different diagnoses
    - 4 different diagnosis methods
    - Ages between 0 and 85
    - 14 locations
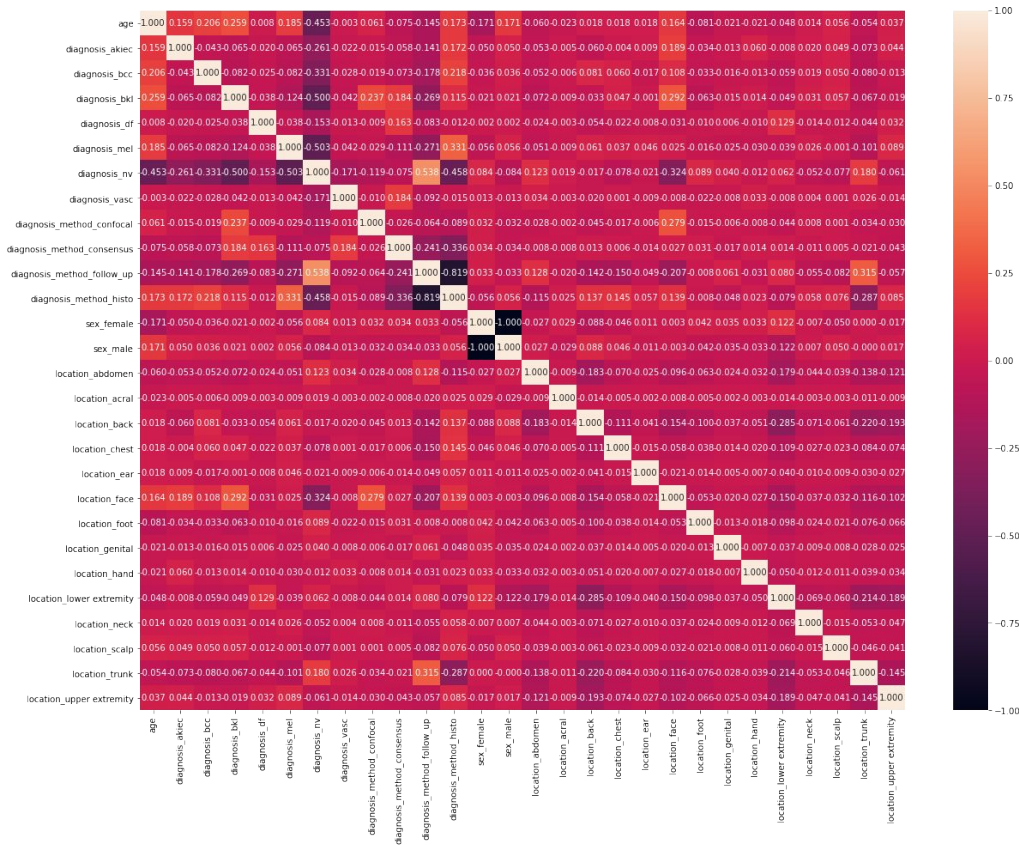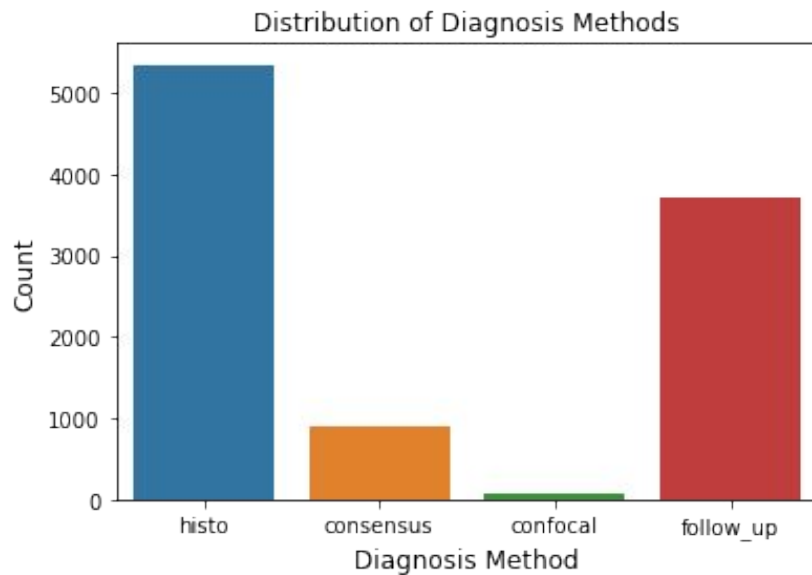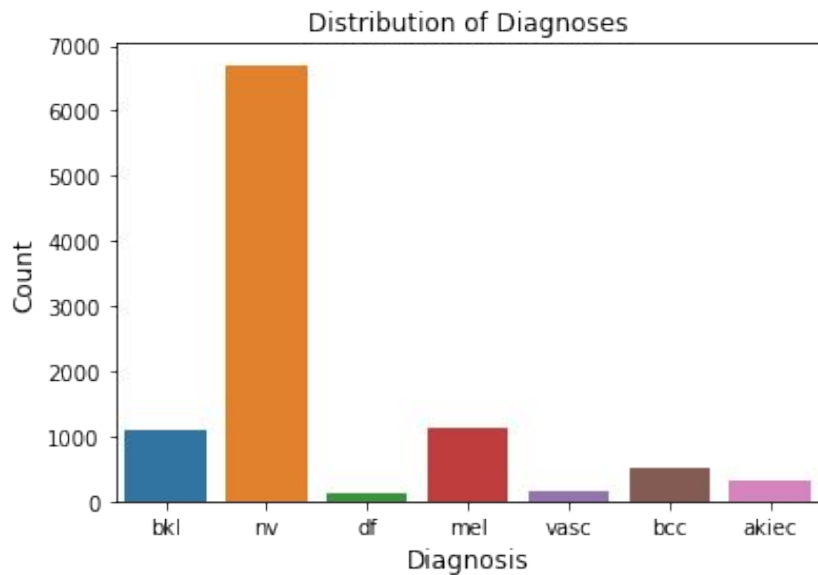    - Male and Female patients

# Data Wrangling

- The metadata was missing some values for different features that needed to be filled
  - All features were combined into one large feature to determine the most common values to fill missing data
- Images were initially all together in two files
  - Not optimal for Tensorflow Keras model
- Images were split based on diagnosis from the metadata and image ID into separate folders labelled by diagnosis
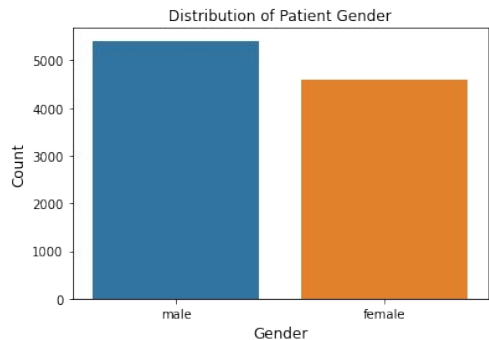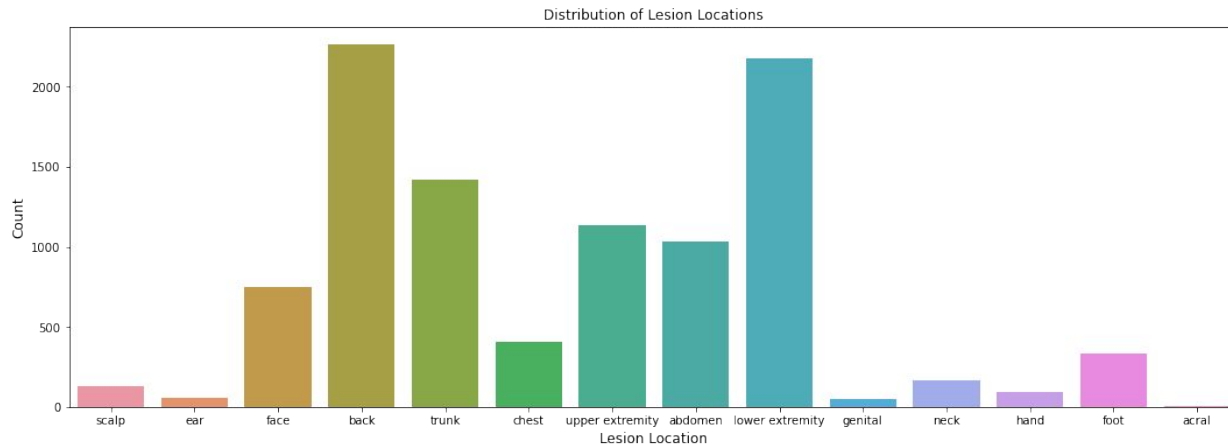
# Exploratory Data Analysis

- The correlation matrix showed some interesting information
  - NV diagnosis had a high negative correlation with age and histopathology and a positive correlation with follow-up
  - The other diagnoses had a range of positive and negative correlations, although not as high as NV
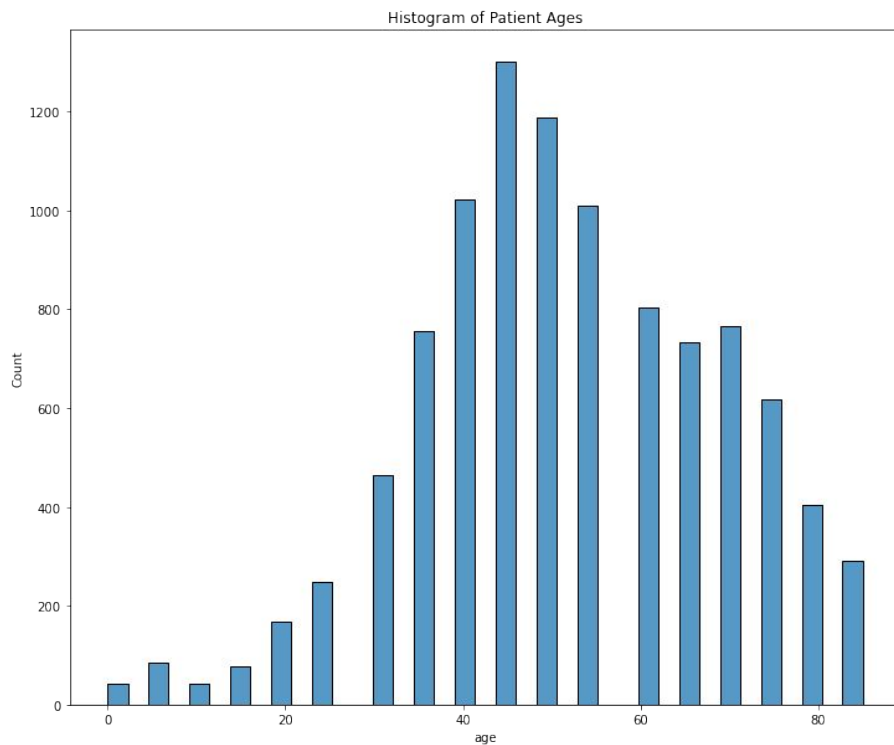
# Exploratory Data Analysis (cont.)
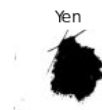
# Exploratory Data Analysis (cont.)



Distribution of Lesion Locations



Distribution of Patient Gender

# Exploratory Data Analysis (cont.)

# Basic Image Analysis

- Another goal of this project was to determine if the images were ready for machine learning and if any interesting information could be gleaned from them as part of the EDA process
- Different processing techniques available through Python libraries were used for this section
- This included:
    - Slicing images by color channels
    - Thresholding
    - Edge detection
    - Exposure adjustment
    - Contour detection

# Basic Image Analysis



Blue Channel Levels



Original / Isodata / Li / Mean / Minimum / Otsu / Triangle / Yen



Adaptive Histogram Equalized



Edges with Canny



Contours

# Basic Image Analysis

- Brought out edge detail as well as underlying surface detail on the skin of the patient.
  - Useful in differentiating between different lesion types
  - Visualizing skin damage not initially visible
- Different lesion types also have different colors and similar shapes
  - Probably not useful to edit the images from their original format

# Machine Learning Setup

- Two types of data were used
  - Images
  - Tabular data
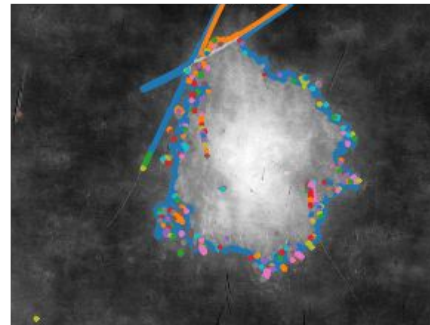- Images were scaled to prepare for machine learning
- Tabular data was optimized using encoding to ensure numerical values
  - Data not relevant to the analysis was dropped
- Machine learning models chosen from typically used models for multi-class classification
  - Images
    - Sequential Keras model with multiple layers
      - Rescaling, 2D Convolutional layers, 2D MaxPooling layers, Dropout layer to prevent overfitting, Flatten layer, and Dense layer
  - Tabular Data
    - Random Forest Classifier
    - Easy Ensemble Classifier
    - KNN Classifier
    - Naive Bayes Classifier
- Image model assessed on accuracy and loss
- Tabular data model assessed on accuracy and classification report metrics

# Machine Learning Metrics

Metrics from best image model:

| Model | Accuracy | Loss |
|---|---|---|
| Keras Sequential | 73% | 0.75 |

Metrics from best tabular data model:

| Model | Accuracy |
|---|---|
| Random Forest Classifier | 73% |

# Keras Sequential Model

- The best image classification model utilized the following features:
  - Data augmentation
    - Random flips
    - Random rotation
    - Random zoom
  - 11 additional layers, including:
    - Rescaling
    - Convolutional 2D
    - MaxPooling 2D
    - Dropout
    - Flatten
    - Dense
  - Adam optimizer
  - Sparse categorical cross entropy loss
  - Accuracy metric
  - 15 Epochs
  - 128 Batch Size

# Random Forest Classifier

- Best tabular data model utilized a generic random forest classifier
  - 200 Estimators
  - Accuracy metric
    - 73%
  - Classification report

# Random Forest Classifier

| Target | Precision | Recall | F1-Score | Support |
|--------|-----------|--------|----------|---------|
| **0** | 0.21 | 0.12 | 0.15 | 69 |
| **1** | 0.29 | 0.24 | 0.26 | 93 |
| **2** | 0.48 | 0.43 | 0.45 | 228 |
| **3** | 0.73 | 0.29 | 0.41 | 28 |
| **4** | 0.36 | 0.24 | 0.29 | 226 |
| **5** | 0.84 | 0.95 | 0.89 | 1338 |
| **6** | 0.46 | 0.29 | 0.35 | 21 |

# Conclusion

# Conclusions

- Machine learning was able to work for both classification using image and text data
- Both models produced very similar accuracy results
- Since the text data had distinct features, we were able to pull out metrics for each classification
  - Based on the precision, recall, and F1 scores, the classes were very difficult to predict from the provided datasets
  - Even after running cross validation, the accuracy didn't increase
  - Our best metrics came from classes which had the highest instances, meaning more data to analyze
- Moving forward, more precise fine tuning may be necessary to adjust for the class imbalances in our modeling
- Experiment with combining the methods in a wide and deep model