

# Cancer Detection in Histopathological Slides

Nitin Mahajan  
Milestone Report 2

## OVERVIEW

Cancer is a deadly disease. Researchers and clinicians are trying to find the methods to detect it at early stages. Early diagnosis will play an important role in planning the treatment plan and improvement of the patient's survival rate. Cancer can be benign (localized) or metastatic (spread to distant organs). One of the most important early diagnosis is detection in lymph nodes to find out whether the cancer has metastasized or not. The method to do this is H & E staining of histological slides of lymph nodes taken from biopsies.

## GOAL

Currently pathologists manually examine the slides in the laboratory and decide if the patient has metastatic cancer or not. Reading the slides and making a report based on human judgement which can be inconsistent and vary from person to person and from day to day. Therefore, developing a computation model to read the slides would provide and can automate the process to give unbiased results.

## DATA SOURCE

The data for this project are downloaded from Kaggle website  
<https://www.kaggle.com/c/histopathologic-cancer-detection/data>

## WHAT WE ARE LOOKIING AT (IMAGES)

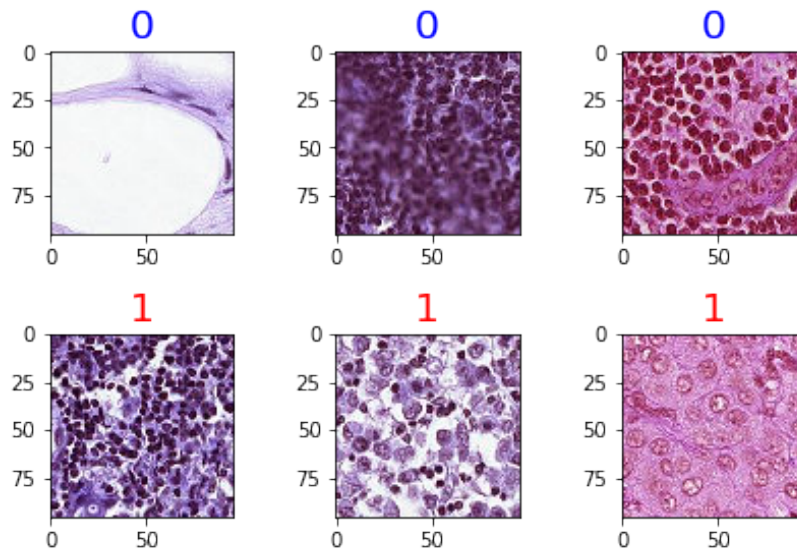
The histopathological images are glass slide microscope images of lymph nodes that are stained with hematoxylin and eosin (H&E). This staining method is one of the most widely used in medical diagnosis and it produces blue, violet and red colors. Dark blue hematoxylin binds to negatively charged substances such as nucleic acids and pink eosin to positively charged substances like amino-acid side chains (most proteins). Typically, nuclei are stained blue, whereas cytoplasm and extracellular parts in various shades of pink.

Lymph nodes are small glands that filter the fluid in the lymphatic system and they are the first place a breast cancer is likely to spread. Histological assessment of lymph node metastases is part of determining the stage of breast cancer. The diagnostic procedure for pathologists is tedious and time-consuming as a large area of tissue has to be examined and small metastases can be easily missed.

## EXPLORATORY DATA ANALYSIS

Following data are provided:-

1. Sample\_submission.csv - a sample submission file
2. Train\_labels.csv – A CSV file with labels of 0 or 1 (0 for cancer not detected and 1 for cancer detected) for corresponding images in training dataset.
3. Train – A directory with 220,025 images. TRAINING DATA SET
4. Test - A directory with 57,458 images. TEST DATA SET
5. The training data set have 130908 and 89117 images with '0' and '1' label, suggesting the data is imbalanced.
6. Below are the representative images for normal and cancer



## DATA WRANGLING & SAMPLING OF IMAGES

Since the dataset is very large and neural networks take very long to train on all the datasets, I decided to sample 10,000 images in each class (a total of 20,000 images) to make the dataset balanced and smaller yet containing enough images to train my models. Once sampling is completed, move the images to separate folders to be consistent in training different models' multiple times. As the dataset has already splitted the training and the test data set, there is no need to split data in train and test. However, I will split the training data in validation and training data sets in a ratio of 20/80. data

## MACHINE & RESOURCES FOR BUILDING CONVOLUTIONAL NEURAL NETWORK MODELS

We need GPUs for sure to handle the image data and building the CNN models. As motioned earlier, I don't have GPUs and so I subsampled the data set. I will be building the CNN models on my computer which has following configuration.

Laptop: Apple MacBook Pro

Processor: 2.6GHz 6-Core Intel Core i7

Memory: 16GB 2400 MHz DDR4

Modules and libraries: Python and Keras with TensorFlow in backend will be the main language and library

Reading Material: Deep Learning with Python (Francois Chollet), Udemy Course – AZ Machine Learning and other online resources.

## MODEL DESCRIPTIONS

### Model #1

**Model:** "sequential\_1"

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 95, 95, 32)	416
conv2d_2 (Conv2D)	(None, 94, 94, 32)	4128
conv2d_3 (Conv2D)	(None, 93, 93, 32)	4128
max_pooling2d_1 (MaxPooling2D)	(None, 46, 46, 32)	0
conv2d_4 (Conv2D)	(None, 45, 45, 32)	4128
conv2d_5 (Conv2D)	(None, 44, 44, 32)	4128
max_pooling2d_2 (MaxPooling2D)	(None, 22, 22, 32)	0
conv2d_6 (Conv2D)	(None, 21, 21, 64)	8256
conv2d_7 (Conv2D)	(None, 20, 20, 64)	16448
max_pooling2d_3 (MaxPooling2D)	(None, 10, 10, 64)	0
flatten_1 (Flatten)	(None, 6400)	0
dense_1 (Dense)	(None, 64)	409664
dropout_1 (Dropout)	(None, 64)	0
dense_2 (Dense)	(None, 2)	130
Total params: 451,426		
Trainable params: 451,426		
Non-trainable params: 0		

## Model #2

Model: "sequential\_2"

Layer (type)	Output Shape	Param #
=====		
conv2d_8 (Conv2D)	(None, 95, 95, 32)	416
conv2d_9 (Conv2D)	(None, 94, 94, 32)	4128
conv2d_10 (Conv2D)	(None, 93, 93, 32)	4128
max_pooling2d_4 (MaxPooling2D)	(None, 46, 46, 32)	0
conv2d_11 (Conv2D)	(None, 45, 45, 32)	4128
conv2d_12 (Conv2D)	(None, 44, 44, 32)	4128
conv2d_13 (Conv2D)	(None, 43, 43, 32)	4128
max_pooling2d_5 (MaxPooling2D)	(None, 21, 21, 32)	0
conv2d_14 (Conv2D)	(None, 20, 20, 64)	8256
conv2d_15 (Conv2D)	(None, 19, 19, 64)	16448
conv2d_16 (Conv2D)	(None, 18, 18, 64)	16448
max_pooling2d_6 (MaxPooling2D)	(None, 9, 9, 64)	0
conv2d_17 (Conv2D)	(None, 8, 8, 128)	32896
conv2d_18 (Conv2D)	(None, 7, 7, 128)	65664
conv2d_19 (Conv2D)	(None, 6, 6, 128)	65664
max_pooling2d_7 (MaxPooling2D)	(None, 3, 3, 128)	0
flatten_2 (Flatten)	(None, 1152)	0
dense_3 (Dense)	(None, 64)	73792
dropout_2 (Dropout)	(None, 64)	0
dense_4 (Dense)	(None, 2)	130
=====		
Total params: 300,354		
Trainable params: 300,354		
Non-trainable params: 0		

## Model # 3

Model: "model\_1"

Layer (type) to	Output Shape	Param #	Connected
=====			
input_2 (InputLayer)	(None, 96, 96, 3)	0	
=====			
mobilenet_1.00_224 (Model)	multiple	3228864	input_2[0] [0]
=====			
global_max_pooling2d_1 (GlobalM	(None, 1024)	0	mobilenet_1.00_224[1][0]
=====			
global_average_pooling2d_1 (Glo	(None, 1024)	0	mobilenet_1.00_224[1][0]
=====			
flatten_3 (Flatten)	(None, 9216)	0	mobilenet_1.00_224[1][0]
=====			
concatenate_1 (Concatenate)	(None, 11264)	0	global_max_pooling2d_1[0][0]
=====			
global_average_pooling2d_1[0][0]			global_max_pooling2d_1[0][0]
=====			
dropout_3 (Dropout)	(None, 11264)	0	concatenate_1[0][0]
=====			
3_ (Dense)	(None, 2)	22530	dropout_3[0][0]
=====			
=====			
Total params: 3,251,394			
Trainable params: 3,229,506			
Non-trainable params: 21,888			
=====			
=====			

## Model performances:

	val_loss	val_acc	roc_auc_scores
Model1	0.740433	0.839625	0.922942
Model2	0.590117	0.845375	0.920965
Model8 (MobileNet)	0.134585	0.908625	0.968650

## Confusion Matrix

		Model1		Model2		Model8	
True Label	Cancer	1586	414	1654	346	1835	165
	Normal	227	1773	284	1716	192	1808
		Cancer	Normal	Cancer	Normal	Cancer	Normal
		Predicted Label		Predicted Label		Predicted Label	

	Precision		Recall		f1-score	
	Cancer	Normal	Cancer	Normal	Cancer	Normal
Model1	0.87	0.81	0.79	0.89	0.83	0.85
Model2	0.85	0.83	0.83	0.86	0.84	0.84
Model8	0.91	0.92	0.9	0.9	0.91	0.91

## Conclusion

Above results clearly demonstrates Model 8(MobileNet) performed the best.

## FUTURE DIRECTIONS

- Try some other transfer learning models like ResNet50, NASNetMobile etc,
- Trained the model used only 10,000 images in each class. However, with high power computers/GPU one can increase or use all set for the analysis.