



# Find the nuclei in divergent images to advance medical discovery

## Capstone Project Proposal

15.05.2018

---

Jimit Jaishwal

---

## Domain Background

Imagine speeding up research for almost every disease, from lung cancer and heart disease to rare disorders. We've all seen people suffer from diseases like cancer, heart disease, chronic obstructive pulmonary disease, Alzheimer's, and diabetes. Many have seen their loved ones pass away. Think how many lives would be transformed if cures came faster. By automating nucleus detection, you could help unlock cures faster—from rare disorders to the common cold. Please see this video you can better understand the problem, this is a part of Data Science Bowl 2018 in Kaggle - [View This Video](#)

## Problem Statement

Identifying the cells' nuclei is the starting point for most analysis because human body's 30 trillion cells contain a nucleus full of DNA, the genetic code that programs each cell. Identifying nuclei allows researchers to determine each cell in a sample, and by measuring how cells react to various treatments, the researcher can understand the underlying biological processes at work.

I've used this dataset - [Find the nuclei in divergent images to advance medical discovery](#) and thus created a computer model that can identify a range of nuclei across varied conditions. By observing patterns, asking questions, and building a model.

## Datasets and Inputs

Our dataset contains a large number of nuclei images. The images were acquired under a variety of conditions and vary in the cell type, and imaging modality (brightfield vs. fluorescence). The dataset is designed to challenge an algorithm's ability to generalize across these variations.

Each image is represented by an associated `ImageId`. Files belonging to an image are contained in a folder with this `ImageId`. Within this folder are two subfolders:

- `Images` this folder contains the image file.
- `masks` this folder contains the segmented masks of each nucleus. This folder is only included in the training set. Each mask contains one nucleus.

## File descriptions

- `/train/*` - training set images (images and annotated masks) [Download](#)
- `/test/*` - stage 1 test set images (images only, you are predicting the masks) [Download](#)

---

The most important file is train training set images and annotated masks, we have 670 training examples each image is 128 \* 128 pixels and depth is 3

## Solution Statement

I am using deep learning algorithm and make with Tensorflow/Keras and trained with training data. In this project i use Convolutional neural network will be implemented Tensorflow/Keras library and optimized to minimize multi-class logarithmic loss as defined in Evolution and Metrics section, In many CNN archistrucre I use U-net Convolutional neural network. The u-net is convolutional network architecture for fast and precise segmentation of images. Prediction will be made on the test dataset and will be evaluated.

## Benchmark Model

The model with the Private Leaderboard score(Intersection Over Union) of 0.614 will be used as a benchmark model. Attempt will be made so that score(Intersection Over Union) obtained will be among the top 50% of the Private Leaderboard submissions.

## Evaluation Metrics

Submissions are evaluated using the Multiclass logarithmic loss. Multiclass logarithmic loss is used to assess predictive models in high-cardinality classification problems. Each image has been labeled with one mask image file. For each image, you must submit a set of predicted mask image (one for every image). The formula is then,

$$\text{log-loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij})$$

where N is the number of images in the test set, M is the number of image class labels as different type Image Mask file, log is the natural logarithm, y is the true mask image and p is predicted mask image.

---

## Project Design

The general sequence of steps are as follows:

- **Data Visualization** - Visual various type images and mask file of nuclei, we also use visualization to show predicted mask file of give microscopic images.
- **Data Preprocessing** - we scale and normalize our all features. The values should be in the range of 0 to 1
- **Model Selection** - In our project dataset is images so i choose deep learning U-net CNN model, because that was good image segmentation.

We have already training and testing dataset. When we train our model then we use our test dataset and predict mask files and save results.

## References

- <https://www.kaggle.com/c/data-science-bowl-2018>
- <https://arxiv.org/pdf/1505.04597.pdf>
- <https://www.datasciencecentral.com/profiles/blogs/polymorphic-malware-detection-using-sequence-classification>
- <https://storage.googleapis.com/kaggle-media/competitions/dsb-2018/dsb.jpg>
- <https://www.youtube.com/watch?v=eHwkfhmJexs&feature=youtu.be>