# Finding Generalizations in Multi-class Image Classifiers
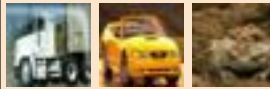
Claudia Richoux - CMSC35200 Final Project

## Introduction

Seeing how neural nets do things, and what they learn to look for while doing that, is both important for interpretability, and also really cool. People who study thought often study "different" brains in order to better understand the general case- for example, hippocampal lesions for studying memory, or Freud's case studies, or Descartes and his dogs. I thought applying this to neural nets might yield something interesting, so I mistrained a bunch of CNNs and then did a bunch of statistics on their performance.

## Methodology

I made a small CNN to train on the CIFAR-10 image dataset. Then, for each choice of two CIFAR-10 classes, I repeatedly randomly selected half the members of each class to mislabel as the second, and trained a net on them. This is so that the net would learn nothing about either of the confused classes (or at least, the information it DID pick up about them would cancel out in further computations!), and no other classes would be affected beyond the direct effect of not learning these two classes. I repeated this process 50 times per selection of classes to confuse, to ensure I could draw accurate assumptions about the loss/accuracy of models trained in this way- fortunately, these statistics tended to have small $\sigma$ for each experiment, so averaging was more or less valid :)
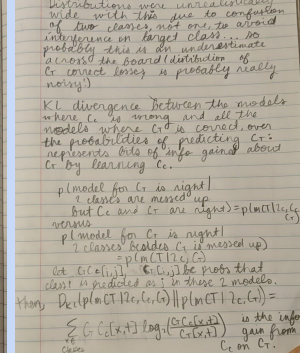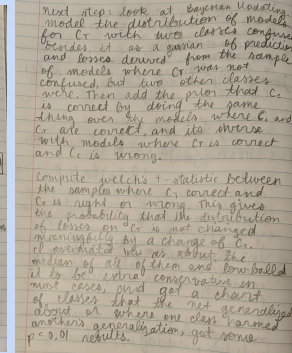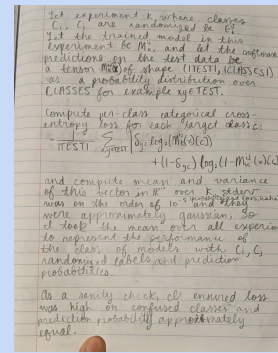


truck, car, frog

```
Layer (type)                Output Shape        Param #
conv2d_6 (Conv2D)           (None, 30, 30, 32)  896
max_pooling2d_4 (MaxPooling2 (None, 15, 15, 32)  0
conv2d_7 (Conv2D)           (None, 13, 13, 64)  18496
max_pooling2d_5 (MaxPooling2 (None, 6, 6, 64)    0
conv2d_8 (Conv2D)           (None, 4, 4, 64)    36928
flatten_2 (Flatten)         (None, 1024)        0
dense_4 (Dense)             (None, 64)          65600
dense_5 (Dense)             (None, 10)          650
Total params: 122,570
Trainable params: 122,570
Non-trainable params: 0
```

model architecture

## Data Analysis (couldn't typeset this, so I took pics)



## Results :)

Welch's t-test shows that it generalized about cars and trucks (and also cars and frogs- look at headlights vs. frog eyes :P). Bird-dog significance was a statistical error with estimation of \nu. Learning about both animals and man-made things tended to confuse it- they look pretty different!
KL-entropy shows that classes whose accuracies correlate can gain up to 0.025 bits of information on classifying T by not confusing C, or lose almost the same amount by classes that clash. On the entropy plot, classes 2-7 are the animals- you see patches of green and purple because the net tends to correlate animals and vehicles separately :)

| target | confused | stdevs |
|---|---|---|
| airplane | dog | −2.78, p < 0.01 |
| automobile | frog | 2.116, p < 0.05 |
| automobile | truck | 2.541, p < 0.02 |
| bird | dog | 2.129, p < 0.05 |
| bird | ship | −1.99, p < 0.05 |
| deer | airplane | −2.21, p < 0.05 |
| deer | truck | −2.91, p < 0.01 |
| frog | airplane | −2.23, p < 0.05 |
| truck | deer | −2.14, p < 0.05 |

## Future Ideas?

- try this with other types of models (text generation? sentiment analysis?)
- deeper models, more classes
- identify weird combos, identify important weights for those generalizations, target Lucid @ those…
- lots of other ways to mistrain...



bits of entropy for discrimination depending on losing a class