

# ECE 50024: Homework 3

Manish Kumar Krishne Gowda, 0033682812 (Spring 2023)

## Exercise 1:

a) To PT  $x^T A x = \text{tr}[A x x^T]$

w.k.t  $\text{tr}(ABC) = \text{tr}(BCA)$

$$\Rightarrow \text{tr}(x^T A x) = \text{tr}(A x x^T)$$

$$= x^T A x = \text{tr}(A x x^T). \rightarrow \textcircled{1}$$

$x^T A x$  is a scalar &  $\text{tr}(\text{scalar}) = \text{Scalar}$   
( $1 \times 1$  matrix)

b)  $\hat{\Sigma}_{ML} = \underset{\Sigma}{\arg \max} p(D | \Sigma).$

likelihood of D is  $p(D | \Sigma) = \prod_{n=1}^N \left\{ \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu) \right] \right\}$

$$= \frac{1}{(2\pi)^{Nd/2} |\Sigma|^{N/2}} \exp \left\{ -\frac{1}{2} \sum_{n=1}^N \underbrace{(x_n - \mu)^T}_{\mathbb{x}^T} \underbrace{\Sigma^{-1}}_{\mathbb{A}} \underbrace{(x_n - \mu)}_{\mathbb{x}} \right\}$$

$$= \frac{1}{(2\pi)^{Nd/2} |\Sigma|^{N/2}} \exp \left\{ -\frac{1}{2} \text{tr} \left( \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T \right) \right\}$$

from \textcircled{1}

9 we have  $A = \Sigma^{-1} \tilde{\Sigma}$ ,  $\lambda_1, \lambda_2, \dots, \lambda_d$  are eigenvalues of  $A$ .

w.k.t  $\text{tr}[A] = \sum_{i=1}^d \lambda_i$

$$|A| = \prod_{i=1}^d \lambda_i \quad (\text{i.e. det of a matrix is the product of its eigenvalues})$$

thus substituting for  $\Sigma^{-1} = \frac{A}{\tilde{\Sigma}}$

$$\Rightarrow p(D | \Sigma) \propto \left( \frac{1}{(2\pi)^{Nd/2} |\Sigma|^{N/2}} \exp \left[ -\frac{1}{2} \text{tr} \left( \frac{A}{\tilde{\Sigma}} \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T \right) \right] \right)$$

$$= \frac{1}{(2\pi)^{Nd/2} |\sum|^N/2} \left( \prod_{i=1}^d \lambda_i \right)^{N/2} \exp \left\{ -\frac{N}{2} \sum_{i=1}^d \lambda_i \right\}.$$

d). We have

$$P(D|\Sigma) = \frac{1}{(2\pi)^{Nd/2} |\sum|^N/2} \left( \prod_{i=1}^d \lambda_i \right)^{N/2} \exp \left\{ -\frac{N}{2} \sum_{i=1}^d \lambda_i \right\}.$$

$$\lambda_i^{\max} = \underset{\lambda_i}{\operatorname{argmax}} P(D|\Sigma)$$

$$\lambda_1^{\max} = \underset{\lambda_1}{\operatorname{argmax}} \log P(D|\Sigma),$$

$$= \underset{\lambda_1}{\operatorname{argmax}} \left\{ -\log(2\pi)^{Nd/2} - \log |\sum|^{N/2} + \frac{N}{2} \sum_{i=1}^d (\log \lambda_i) - \frac{N}{2} \sum_{i=1}^d \lambda_i \right\}$$

$$\frac{\partial \lambda_1^{\max}}{\partial \lambda_1} = 0 \quad = \underset{\lambda_1}{\operatorname{argmax}} \left\{ \frac{N}{2} \sum_{i=1}^d (\log \lambda_i) - \frac{N}{2} \sum_{i=1}^d \lambda_i \right\}$$

$\Rightarrow$

$$\frac{\partial \lambda_1^{\max}}{\partial \lambda_1} = 0 \Rightarrow \frac{N}{2} \frac{1}{\lambda_1} - \frac{N}{2} = 0$$

$$\Rightarrow \underline{\underline{\lambda_1 = 1}}$$

$$\text{Hence } \lambda_2 = \lambda_3 = \dots = \lambda_d = 1$$

as

e) as  $\lambda_1 = \lambda_2 = \dots = \lambda_d = 1$  for  $A \in A$  is diagonalizable, then  $A = I_{d \times d}$

$$\Rightarrow I = \Sigma^{-1} \tilde{\Sigma}$$

$$\Rightarrow \tilde{\Sigma} = (\Sigma \Sigma)^{-1} \Sigma$$

$$\Rightarrow \tilde{\Sigma} = \hat{\Sigma} = \hat{\Sigma}_{ML}$$

thus  $\hat{\Sigma}_{ML} = \tilde{\Sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T$

f) alternative way would be to maximise  
over the likelihood w.r.t  $\Sigma$ ; here we take  
the gradient of  $-\log(P(D|\Sigma))$  w.r.t  $\Sigma$  & set  
it to zero.

g)  $E\left(\hat{\Sigma}_{ML}\right) = \Sigma$  if we desire  $\hat{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu})(x_n - \hat{\mu})^T$

## Exercise 2

a)

2a). we have

$$P_{Y|X}(C_i|x) \geq \sum_{C_0}^C P_{Y|X}(C_0|x)$$

$$= \frac{P_{XY}(x|C_i) P_Y(C_i)}{P_X(x)} \geq \sum_{C_0}^C \frac{P_{XY}(x|C_0) P_Y(C_0)}{P_X(x)}$$

$$\Rightarrow \left( \frac{1}{(\frac{\partial x}{\partial \mu_1})^{1/2} |\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right\} \pi_1 \right) \geq$$

$$\left( \frac{1}{(\frac{\partial x}{\partial \mu_0})^{1/2} |\Sigma_0|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \right\} \pi_0 \right).$$

taking log on b.s. (as log is a monotonically increasing function)

$$-\frac{1}{2} \log |\Sigma_1| + \log \pi_1 - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \geq -\frac{1}{2} |\Sigma_0| + \log \pi_0$$

$$-\frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0)$$

### (b) Estimate $\mu_1, \mu_0, \Sigma_1, \Sigma_0, \pi_1$ and $\pi_0$ in Python. Report:

(i) The first 2 entries of the vector  $\mu_1$  and the first 2 entries of the vector  $\mu_0$ .

$$\mu_0 = [[0.48249575], [0.4864399]]$$

$$\mu_1 = [[0.44080734], [0.43871359]]$$

(ii) The first  $2 \times 2$  entries of the matrix  $\Sigma_1$  and the first  $2 \times 2$  entries of the matrix  $\Sigma_0$ .

$$\Sigma_0 = [0.06447725 0.03691294] \\ [0.03691294 0.06622764]]$$

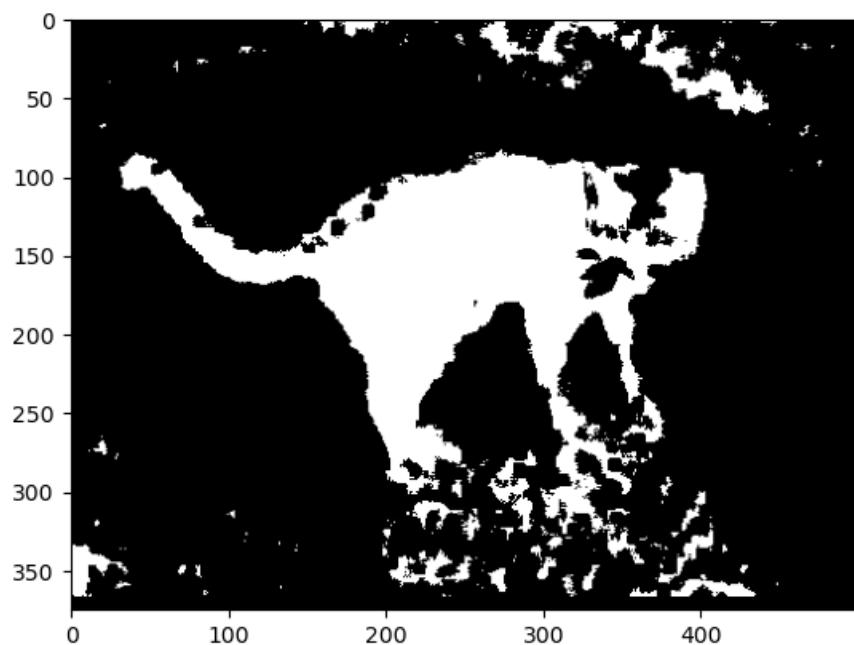
$$\Sigma_1 = [[0.04305652 0.03533616], [0.03533616 0.042466]]$$

(iii) The values of  $\pi_1$  and  $\pi_0$ .

$$\pi_0 = 0.828650711064863$$

$$\pi_1 = 0.171349288935137$$

c)



d)

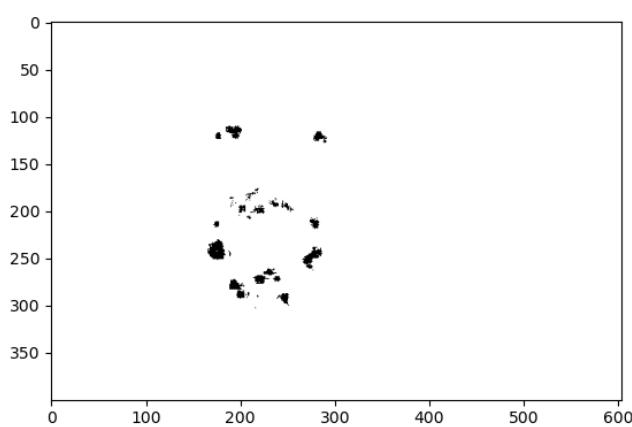
MAE = 51.200803343436455

e)

Internet image I used :



Predicted mask:

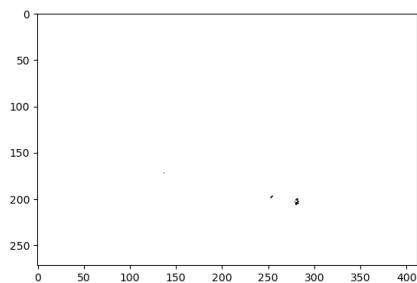


As we see the model performed very bad as we did not have training data corresponding to the new image. We used the old mean, and covariance matrix which led to sub-optimal performance. We can say the model overfit to the original image. Note that this is not a one of case that the prediction is bad. I used multiple images and this was the “best” output. Following is another example of the bad performance

Input:



Predicted mask:



CODE:

```
1 import numpy as np
2 import scipy.stats as stats
3 import matplotlib.pyplot as plt
4
5 train_cat = np.matrix(np.loadtxt('data/train_cat.txt', delimiter = ','))
6 train_grass = np.matrix(np.loadtxt('data/train_grass.txt', delimiter = ','))
7 print(train_cat.shape)
8 print(train_grass.shape)
9
10
11 mu_cat = np.mean(train_cat, axis=1)
12 sigma_cat = np.cov(train_cat,bias=True)
13 mu_grass = np.mean(train_grass, axis=1)
14 sigma_grass = np.cov(train_grass,bias=True)
15 pi0 = train_grass.size/(train_grass.size+train_cat.size)
16 pi1 = train_cat.size/(train_grass.size+train_cat.size)
17
18 print(mu_cat[:2])
19 print(mu_grass[:2])
20 print(sigma_cat[:2,:2])
21 print(sigma_grass[:2,:2])
22 print(pi0)
23 print(pi1)
24
25 Y = plt.imread('data/cat_grass.jpg') / 255
26 print(Y.shape)
27 truth = plt.imread('data/truth.png')
28 Y = np.asarray(Y)
29 truth = np.asarray(truth)
30 M,N = Y.shape
```

```

31 result = np.empty((M-8,N-8), dtype=float)
32 MAE = 0
33
34 def dec_rule(x, mu, sigma, prior):
35     param1 = -0.5*(x-mu).transpose() @ np.linalg.inv(sigma) @ (x-mu)
36     #print(param1.shape)
37     param2 = np.log(prior)
38     param3 = -0.5*np.log(np.linalg.det(sigma))
39     return (param1 + param2 + param3)
40 for i in range(M-8):
41     for j in range(N-8):
42         block = Y[i:i+8, j:j+8] #
43         x = block.flatten()
44         x = x[:, None]
45         value0 = dec_rule(x,mu_grass,sigma_grass,pi0)
46         #print(value0)
47         value1 = dec_rule(x,mu_cat,sigma_cat,pi1)
48         #print(value1)
49         result[i,j] = 0 if value0>value1 else 255 #can set the else to 1,
50 and later mul matrix by 255
51         #print(result[i,j])
52     MAE = MAE + abs(result[i,j] - truth[i,j])
53
54 MAE = MAE/result.size
55 print(MAE)
56 print(Y.shape)
57 print(result)
58 plt.imshow(result, cmap=plt.cm.gray)
59 #plt.imshow(result, cmap='Greys', interpolation='nearest')
    plt.show()

```

3.

a)

a) we have

$$P_{Y|X}(C_1|x) \geq_{C_0} P_{Y|X}(C_0|x) \text{ as posteriors}$$

$$\Rightarrow \frac{P_{X|Y}(x|C_1) P_Y(C_1)}{P_X(x)} \geq_{C_0} \frac{P_{X|Y}(x|C_0) P_Y(C_0)}{P_X(x)}$$

$$\Rightarrow \frac{P_{X|Y}(x|C_1)}{P_{X|Y}(x|C_0)} \geq_{C_0} \frac{P_Y(C_0)}{P_Y(C_1)} = \frac{\pi_0}{\pi_1} = c \rightarrow (\star)$$

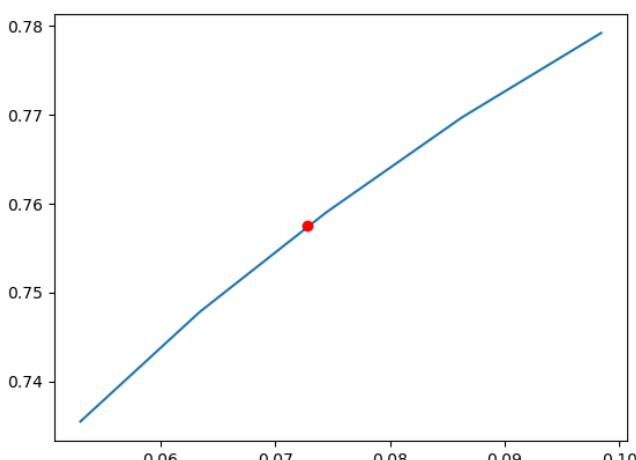
From prev. exercise we have

$$-\frac{1}{2} \log |\Sigma_1| + \log \pi_1 - \frac{1}{2} (x - \mu_0)^T \sum_1^{-1} (x - \mu_1) \geq_{C_0} -\frac{1}{2} \log |\Sigma_0| + \log \pi_0 - \frac{1}{2} (x - \mu_0)^T \sum_0^{-1} (x - \mu_0)$$

$$\Rightarrow -\frac{1}{2} \log |\Sigma_1| + \frac{1}{2} \log |\Sigma_0| - \frac{1}{2} (x - \mu_0)^T \sum_1^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_0)^T \sum_0^{-1} (x - \mu_0) \geq_{C_0} \underbrace{\cdot \log \pi_0 - \log \pi_1}_{\log \left( \frac{\pi_0}{\pi_1} \right)}$$

which is the log form of  $(\star)$

b &amp; c)



```

1 import numpy as np
2 import scipy.stats as stats
3 import matplotlib.pyplot as plt
4 import pickle
5
6 train_cat = np.matrix(np.loadtxt('data/train_cat.txt', delimiter = ','))
7 train_grass = np.matrix(np.loadtxt('data/train_grass.txt', delimiter = ',')) 
8 print(train_cat.shape)
9 print(train_grass.shape)
10
11
12 mu_cat = np.mean(train_cat, axis=1)
13 sigma_cat = np.cov(train_cat,bias=True)
14 mu_grass = np.mean(train_grass, axis=1)
15 sigma_grass = np.cov(train_grass,bias=True)
16 pi0 = train_grass.size/(train_grass.size+train_cat.size)
17 pi1 = train_cat.size/(train_grass.size+train_cat.size)
18
19 Y = plt.imread('data/cat_grass.jpg') / 255
20 print(Y.shape)
21 truth = plt.imread('data/truth.png')
22 truth = np.asarray(truth)
23 print(truth.shape)
24 grnd_truth_negatives = truth[np.where(truth==0)].shape[0]
25 grnd_truth_positives = truth[np.where(truth!=0)].shape[0]
26 print(grnd_truth_positives)
27 print(grnd_truth_negatives)
28 Y = np.asarray(Y)
29 M,N = Y.shape
30 result = np.empty((M-8,N-8), dtype=float)
31
32
33 def dec_rule(x):
34     param1_cat = -0.5*(x-mu_cat).transpose() @ np.linalg.inv(sigma_cat) @ (x-
35 mu_cat)
36     param1_grass = -0.5*(x-mu_grass).transpose() @ np.linalg.inv(sigma_grass) @
37 (x-mu_grass)
38     #print(param1.shape)
39     param2_cat = -0.5*np.log(np.linalg.det(sigma_cat))
40     param2_grass = -0.5*np.log(np.linalg.det(sigma_grass))
41     return (param1_cat + param2_cat - param1_grass - param2_grass)
42
43 tou = np.geomspace(start=0.01, stop=1000, num=5)
44 prob_detection = []
45 prob_false_alarm = []
46 for t in range(0,len(tou)):
47     print(t)
48     true_positives = 0
49     false_positives = 0
50     for i in range(M-8):
51         for j in range(N-8):
52             block = Y[i:i+8, j:j+8] #
53             x = block.flatten()
54             x = x[:, None]
55             value = dec_rule(x)
56             #result = 1 if value > np.log(tou[t]) else 0
57             if(value > np.log(tou[t])):
58                 if(truth[i,j] == 1):
59                     true_positives += 1
60                 else:
61                     false_positives += 1
62     prob_detection.append(true_positives/grnd_truth_positives)

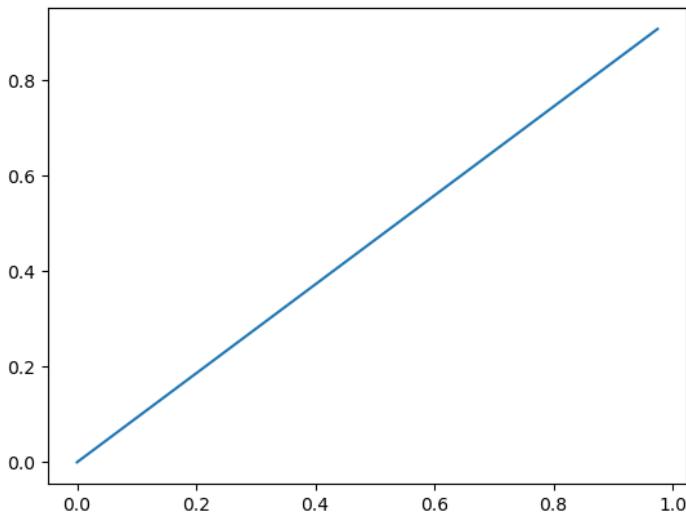
```

```

63     prob_false_alarm.append(false_positives/grnd_truth_negatives)
64
65 tou_bay = pi0/pi1 #bayesian tou
66 true_positives_bay = 0
67 false_positives_bay = 0
68 for i in range(M-8):
69     for j in range(N-8):
70         block = Y[i:i+8, j:j+8] #
71         x = block.flatten()
72         x = x[:, None]
73         value = dec_rule(x)
74         if(value > np.log(tou_bay)):
75             if(truth[i,j] == 1):
76                 true_positives_bay += 1
77             else:
78                 false_positives_bay += 1
79 prob_detection_bay = true_positives_bay/grnd_truth_positives
80 prob_false_alarm_bay = false_positives_bay/grnd_truth_negatives
81 prob_detection = np.asarray(prob_detection)
82 prob_false_alarm = np.asarray(prob_false_alarm)
83 plt.plot(prob_false_alarm,prob_detection,markersize=12)
plt.plot(prob_false_alarm_bay,prob_detection_bay,'ro')
plt.show()

```

d)



```

1 import numpy as np
2 import scipy.stats as stats
3 import matplotlib.pyplot as plt
4 import cvxpy as cp
5
6 train_cat = np.matrix(np.loadtxt('data/train_cat.txt', delimiter = ','))
7 train_grass = np.matrix(np.loadtxt('data/train_grass.txt', delimiter = ',')) 
8 print(train_cat.shape)
9 print(train_grass.shape)
10
11
12 mu_cat = np.mean(train_cat, axis=1)
13 sigma_cat = np.cov(train_cat,bias=True)

```

```

14 mu_grass = np.mean(train_grass, axis=1)
15 sigma_grass = np.cov(train_grass,bias=True)
16
17 Y = plt.imread('data/cat_grass.jpg') / 255
18 print(Y.shape)
19 truth = plt.imread('data/truth.png')
20 truth = np.asarray(truth)
21 print(truth.shape)
22 grnd_truth_negatives = truth[np.where(truth==0)].shape[0]
23 grnd_truth_positives = truth[np.where(truth!=0)].shape[0]
24 print(grnd_truth_positives)
25 print(grnd_truth_negatives)
26 Y = np.asarray(Y)
27 M,N = Y.shape
28 result = np.empty((M-8,N-8), dtype=float)
29
30 train_cat = train_cat.transpose()
31 train_grass = train_grass.transpose()
32 N = train_cat.shape[0] + train_grass.shape[0]
33 print(N)
34 d = train_cat.shape[1] + 1
35 print(d)
36 x = np.vstack((train_cat,train_grass))
37 X = np.column_stack((x, np.ones(N))) #consider the basis function as 1 + x1 +
38 x2+...+x64
39 y = np.vstack((np.ones((train_cat.shape[0],1)),-
40 1*np.ones((train_grass.shape[0],1))))
41 print(X.shape)
42
43 theta = cp.Variable(d)
44 y_cvx = cp.reshape(y, (y.shape[0],))
45 cost = cp.Minimize(cp.sum_squares(X@theta-y_cvx))
46 prob = cp.Problem(cost)
47 prob.solve()
48 theta_cap_cvx = theta.value
49 print(f"theta_cap by cvx method : {theta_cap_cvx}")
50 theta_cap_cvx = theta_cap_cvx[:, None]
51 theta_cap_cvx = theta_cap_cvx[:-1]
52 print(theta_cap_cvx.shape)
53
54 tou = np.geomspace(start=-1, stop=1, num=20)
55 prob_detection = []
56 prob_false_alarm = []
57
58 M,N = Y.shape
59 for t in range(0,len(tou)):
60     print(t)
61     true_positives = 0
62     false_positives = 0
63     for i in range(M-8):
64         for j in range(N-8):
65             block = Y[i:i+8, j:j+8] #
66             x = block.flatten()
67             x = x[:, None]
68             #print(x.shape)
69             value = theta_cap_cvx.transpose() @ x
70             #result = 1 if value > np.log(tou[t]) else 0
71             if(value > tou[t]):
72                 if(truth[i,j] == 1):
73                     true_positives += 1
74                 else:
75                     false_positives += 1

```

```
76     prob_detection.append(true_positives/grnd_truth_positives)
77     prob_false_alarm.append(false_positives/grnd_truth_negatives)
78
79 prob_detection = np.asarray(prob_detection)
80 prob_false_alarm = np.asarray(prob_false_alarm)
81 plt.plot(prob_false_alarm,prob_detection,markersize=12)
82 plt.show()
```

---

# Learning to Reweight Examples for Robust Deep Learning

---

Manish Kumar Krishne Gowda, Purdue University, ECE Department, On Campus MS Student<sup>1</sup>

## Abstract

Deep Neural Networks easily overfit to training set biases and label noises. Regularizers and example reweighting algorithms are popular solutions to these problems. But these algorithms require careful tuning of additional hyperparameters. Traditionally, validation is performed at the end of training, which can be prohibitively expensive if the example weights are treated as some hyperparameters to optimize. To circumvent this, the paper proposes to perform validation at every training iteration to dynamically determine the example weights of the current batch. This approach significantly increases the robustness to training set biases.

## 1. Introduction

The paper proposes a learning algorithm that learns to assign weights to training examples based on their gradient directions. They perform gradient descent on the mini-batch example weights which are initialized to zero. Thereby they minimize the loss on clean unbiased validation set.

The model in the paper is derived from a meta-learning objective towards an online approximation that can fit into any regular supervised training. Instead of minimizing the expected loss for the training set each input example is weighted equally and this reweighted input is used to model the dataset. There are multiple examples of automatic differentiation techniques which will be the correct starting point as it is needed to compute the gradient of the validation loss. The paper itself for example references a couple of other papers which implements it using popular deep learning frameworks such as TensorFlow.

### 1.1. Unfamiliar Concepts

As of now I have limited knowledge of the concept of deep neural networks and its implementation in software. This paper further narrows down the deep learning problem to batch learning which needs to be practically implemented. The paper itself was read in haste and hence a thorough reading to understand the algorithm mentioned and its translation to a practical dataset needs to be realized.

## 2. Issues with Implementation

The core of the paper is the Algorithm 1 Learning to Reweight Examples using Automatic Differentiation in page 4. I am trying to understand the different terminologies in the paper like Deep neural networks, hyperparameters, metaa-learning, MNIST and CIFAR benchmarks etc etc. With these I should be able to replicate one of the two experiments mentioned in the paper i.e. MNIST data imbalance experiments or CIFAR noisy label experiments. Both are computationally expensive so, trying on my local system seems infeasible.

With deep neural network to be in agenda of discussion in class next week i hope to get a few questions cleared and get some hints on the algorithm implementation during office hours.

Some information on the MNIST data imbalance experiment implementation can be found at [shorturl.at/iprM7](http://shorturl.at/iprM7). The results are well replicated and provide some confidence for implementation.

## References

Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.

## A. Do not have an appendix here

**Do not put content after the references.** Put anything that you might normally include after the references in a separate supplementary file.

We recommend that you build supplementary material in a separate document. If you must create one PDF and cut it up, please be careful to use a tool that doesn't alter the margins, and that doesn't aggressively rewrite the PDF file. pdftk usually works fine.

**Please do not use Apple's preview to cut off supplementary material.** In previous years it has altered margins, and created headaches at the camera-ready stage.