# EP3260: Machine Learning Over Networks Peer-review of CA2 of group 2

Stefanos Antaris[*1], Amaru Cuba Gyllensten[†1,2],
Martin Isaksson[‡1,2,3], Sarit Khirirat[§1], and Klas Segeljakt[¶1,2]

[1]*KTH Royal Institute of Technology*
[2]*RISE AI*
[3]*Ericsson Research*

February, 2019

## Contents

## 1 Computer assignment

### 1.1 General comments

- Using Jupyter notebooks is nice and makes it easy for the reviewer. Excellent!

[*]antaris@kth.se

[†]amaru.cuba.gyllensten@ri.se

[‡]martisak@kth.se

[§]sarit@kth.se

[¶]klasseg@kth.se

## 1.2 Load data and preprocessing

- There is no explanation of where this dataset comes from, or how to download it. Therefore it would be difficult to reproduce this.

- Good that you print out the head of the dataset, and use the name of the columns to create $(X, y)$.

- When you use `print` like that, you loose the HTML formatting of Pandas dataframe.

## 1.3 Train and test sets

- A little bit disappointed by the fact that you talk about dev (validation) set but never construct one. This would have been useful here for hyper-parameter tuning.

## 1.4 Logistic ridge regression with different optimizers

- Using a case structure is not good for readability, the different solvers should have been made into separate functions or classes.

- Code is repeated in the different solvers.

- Interesting use of `resource.getrusage(resource.RUSAGE_SELF)`!

- `x = np.array(X_train.iloc[0:6000,:])` reduces the problem, and makes it faster.

- In the for-loop, the number of iterations the solver runs is increasing in each loop. Is this intentional? This means $i - 1$ iterations are wasted in each turn of the loop.

```
for i in range(50):
    [...]
    gde = solver(x.T,y,w,num_iters=i)
    [...]
```

- The training curves look a bit strange

- Missing validation and test loss in the curves

## 1.5 Tuning $\lambda$

- Consider the use of `hyperopt`.

- There should be a plot of cost for different hyperparameter settings.

- There should be a table comparing memory usage.

- There should also be an explanation for the results of the memory usage and reasoning regarding their accuracy. If there is none, then maybe an analytical comparison would have been better.