

Name: YongKyung Oh  
Andrew ID: yongkyuo

# Machine Learning for Text Mining

## Homework 4

### 1. Statement of Assurance

1. Did you receive any help whatsoever from anyone in solving this assignment? Yes / **No**.

If you answered 'yes', give full details? (e.g. "Jane explained to me what is asked in Question 3.4").

2. Did you give any help whatsoever to anyone in solving this assignment? Yes / **No**.

If you answered 'yes', give full details? (e.g. "I pointed Joe to section 2.3 to help him with Question 2").

3. Did you find or come across code that implements any part of this assignment? Yes / **No**.

If you answered 'yes', give full details? (e.g. book & page, URL & location within the page, etc)

### 2. Writeup (40 pts)

(1) [10 pts] Gradient

$$f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{\lambda}{n} \sum_{i=1}^n \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)^2 = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{\lambda}{n} L$$
$$\nabla f(\mathbf{w}) = \mathbf{w} + \frac{\lambda}{n} \nabla L$$

Here  $L$  is non differentiable function, so we should conduct sub-gradient method. Assume that  $g(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i = \sum_j^d (w_i x_{ij})$  where  $\mathbf{x}_i \in \mathbb{R}^d$

$$L = \sum_{i=1}^n (1 - y_i * g(\mathbf{x}_i))^2 \text{ where } \forall i \in I$$
$$\nabla L = \sum_{i=1}^n 2(1 - y_i * g(\mathbf{x}_i)) \left( -y_i * \frac{\partial g}{\partial \mathbf{w}} \right)$$

$g(\mathbf{x}_i)$  is dot product and it can be re-written as follow.

$$g(\mathbf{x}_i) = \sum_j^d (w_i x_{ij}) = \text{sum}(\mathbf{w} \otimes \mathbf{x}) = u = \text{sum}(\mathbf{v})$$
$$\frac{\partial \mathbf{v}}{\partial \mathbf{w}} = \frac{\partial (\mathbf{w} \otimes \mathbf{x})}{\partial \mathbf{w}} = \text{diag}(\mathbf{x})$$
$$\frac{\partial u}{\partial \mathbf{v}} = \frac{\partial (\text{sum}(\mathbf{v}))}{\partial \mathbf{v}} = \vec{1}^T$$

$$\frac{\partial u}{\partial \mathbf{w}} = \frac{\partial u}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial \mathbf{w}} = \vec{1}^T * \text{diag}(\mathbf{x}) = \mathbf{x}^T$$

Therefore, our original equation can be changed as follow

$$\begin{aligned} L &= \sum_{i=1}^n (1 - y_i \mathbf{w}^T \mathbf{x}_i)^2 \text{ where } \forall i \in I \\ \nabla L &= \sum_{i=1}^n 2(1 - y_i \mathbf{w}^T \mathbf{x}_i)(-y_i \mathbf{x}_i^T) \\ &= \sum_{i=1}^n 2\mathbf{x}_i^T (y_i \mathbf{w}^T \mathbf{x}_i y_i - y_i) \\ &= \sum_{i=1}^n 2\mathbf{x}_i^T (\mathbf{w}^T \mathbf{x}_i - y_i) \\ &= 2\mathbf{X}_I^T (\mathbf{X}_I \mathbf{w} - \mathbf{y}_I) \end{aligned}$$

Since  $y_i \in \{-1, 1\}$ , so  $y_i^2 = 1$ .

$$\nabla f(\mathbf{w}) = \mathbf{w} + \frac{\lambda}{n} \nabla L = \mathbf{w} + \frac{2\lambda}{n} \mathbf{X}_I^T (\mathbf{X}_I \mathbf{w} - \mathbf{y}_I)$$

**(2) [10 pts] Hessian**

$$\begin{aligned} \nabla f(\mathbf{w}) &= \mathbf{w} + \frac{\lambda}{n} \nabla L = \mathbf{w} + \frac{2\lambda}{n} \mathbf{X}_I^T (\mathbf{X}_I \mathbf{w} - \mathbf{y}_I) \\ \nabla^2 f(\mathbf{w}) &= \frac{\partial \mathbf{w}}{\partial \mathbf{w}} + \frac{\partial}{\partial \mathbf{w}} \left\{ \frac{2\lambda}{n} \mathbf{X}_I^T (\mathbf{X}_I \mathbf{w} - \mathbf{y}_I) \right\} \\ &= I_d + \frac{2\lambda}{n} \mathbf{X}_I^T (\mathbf{X}_I) \\ &= I_d + \frac{2\lambda}{n} \mathbf{X}^T D \mathbf{X} \end{aligned}$$

Where  $I_d$  is the identity matrix of dimension  $d$  and  $D \in \mathbb{R}^{n \times n}$  is a diagonal matrix with the following elements.

$$D_{ii} = \begin{cases} 1, & \text{if } i \in I \\ 0, & \text{otherwise} \end{cases}$$

(3) [10 pts] Optimality

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{\lambda}{n} \sum_{i=1}^n \xi^2 \\ \text{s. t.} \quad & y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i \\ & \xi_i \geq 0 \\ & \xi_i \geq 1 - y_i \mathbf{w}^T \mathbf{x}_i \end{aligned}$$

When  $1 - y_i \mathbf{w}^T \mathbf{x}_i > 0$ ,  $\xi_i = 1 - y_i \mathbf{w}^T \mathbf{x}_i$ , because of the objective function is min.

$$\begin{cases} 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0 \rightarrow \min \xi^2 = (1 - y_i \mathbf{w}^T \mathbf{x}_i)^2 \\ 1 - y_i \mathbf{w}^T \mathbf{x}_i = 0 \rightarrow \min \xi^2 = 0 \\ 1 - y_i \mathbf{w}^T \mathbf{x}_i < 0 \rightarrow \min \xi^2 = 0 \geq 0 \end{cases}$$
$$\therefore \xi_i = \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

Substitute to the original equation

$$f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{\lambda}{n} \sum_{i=1}^n \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)^2$$

(4) [10 pts] Algorithm Pseudo Code

(Outline the pseudo code of the optimization update procedure for mini-batch stochastic gradient method and Newton method)

**Mini-batch stochastic gradient method (Pegasos Algorithm)<sup>1</sup>**

Input: a list of feature vectors  $X$ ; a list of output  $y$ ; regularization parameter  $\lambda$

Output: weight vector  $w$

$w = (0, 0, \dots, 0)$

**for**  $t$  in max\_iteration

    select a position  $i$  randomly (mini-batch set with batchsize  $n$ )

$\eta = \frac{1}{\lambda t}$  step\_size (learning rate with decreasing factor)

    score =  $y_i(\mathbf{w}^T \mathbf{x}_i)$

**if** score < 1

$\mathbf{w} = (1 - \eta\lambda)\mathbf{w} + (\eta y_i)\mathbf{x}_i$

**else**

$\mathbf{w} = (1 - \eta\lambda)\mathbf{w}$

the end result is  $\mathbf{w}$

---

<sup>1</sup> [https://svn.spraakdata.gu.se/repos/richard/pub/ml2015\\_web/a2\\_clarification.pdf](https://svn.spraakdata.gu.se/repos/richard/pub/ml2015_web/a2_clarification.pdf)

## Newton method<sup>2</sup>

Input: a list of feature vectors  $X$ ; a list of output  $y$ ; regularization parameter  $\lambda$

Output: weight vector  $w$

$w = (0, 0, \dots, 0)$

**for**  $t$  in max\_iteration

**if**  $f(w_t)$  is sufficiently small **then**

$w = w_t$

**return**  $w^*$

**end**

$$w_{t+1} = w_t - \frac{f(w_t)}{f'(w_t)}$$

**if**  $|w_{t+1} - w_t|$  is sufficiently small **then**

$w = w_{t+1}$

**return**  $w^*$

**end**

the end result is  $w$

---

<sup>2</sup> <https://www.math.usm.edu/lambers/mat460/fall09/lecture10.pdf>

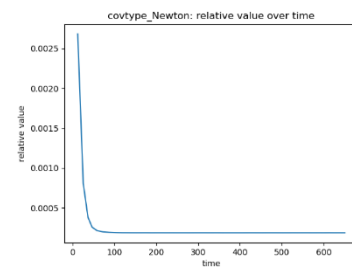
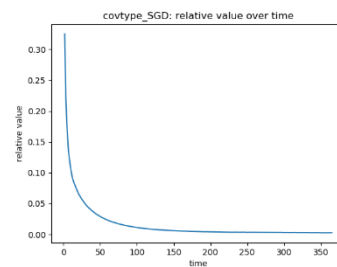
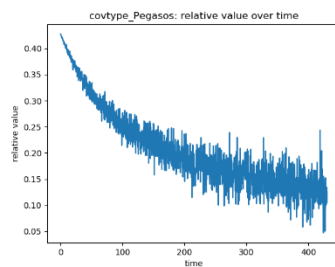
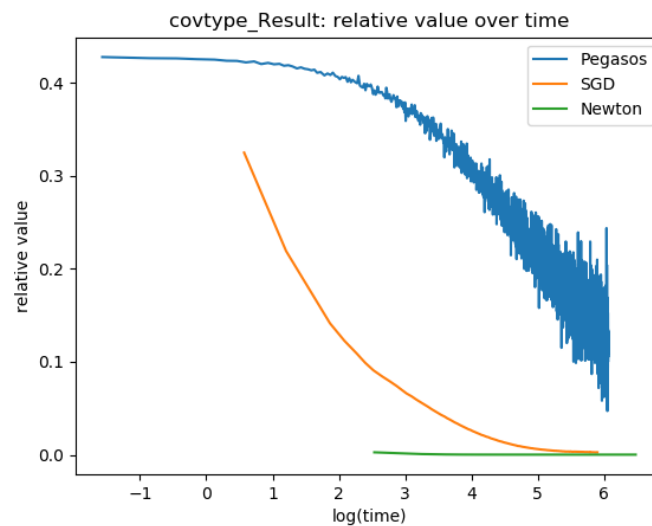
### 3. Experiments (20 pts)

Plot the figures for **both** of two datasets and **both** the approaches

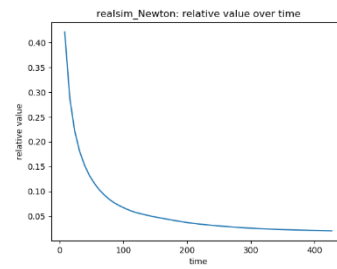
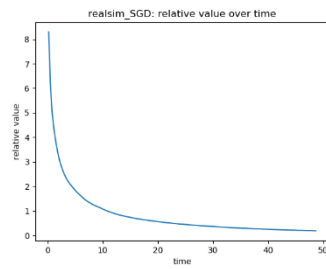
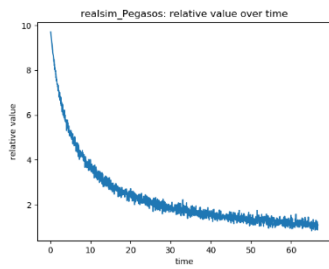
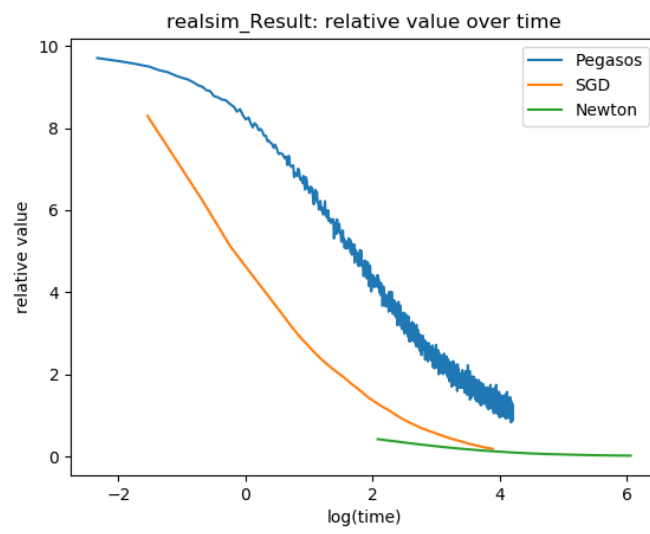
I conducted both mini-batch SGD (Pegasos) and SGD to compare the result. Three different approaches are plotted with log(time) axis.

(1) [5 pts] Relative function value difference versus training time

<contype>

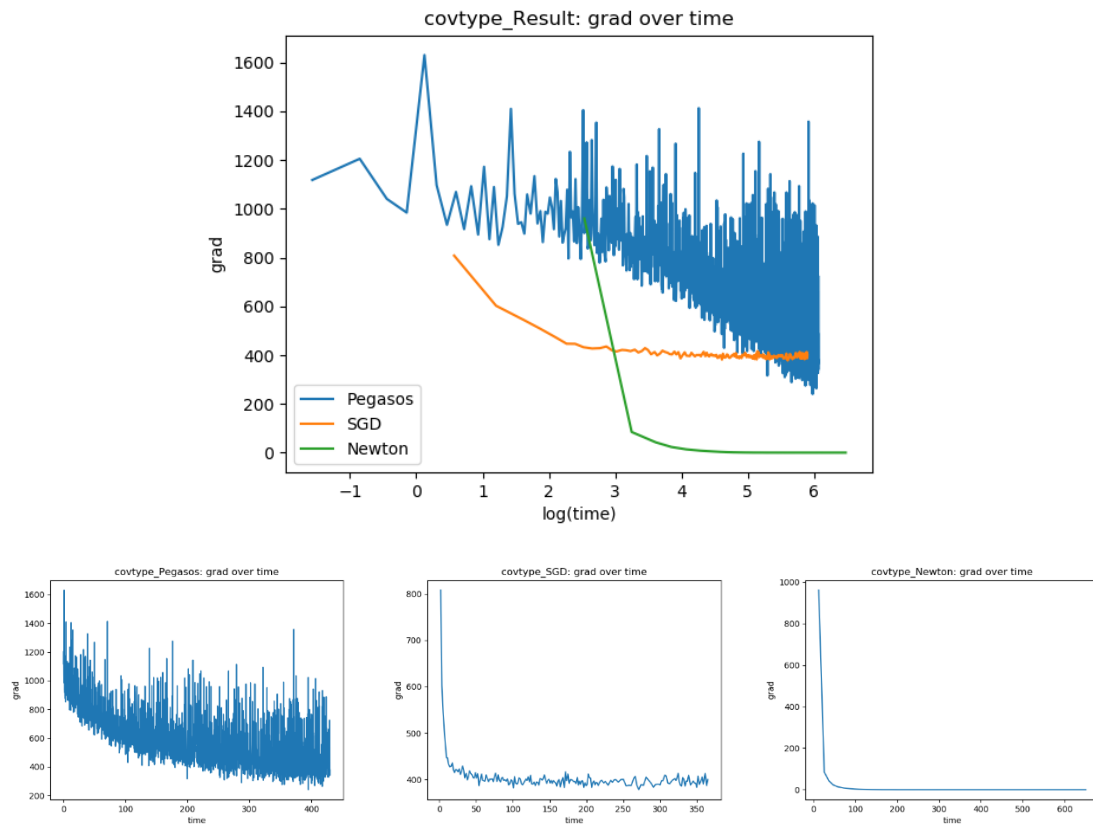


<realsim>

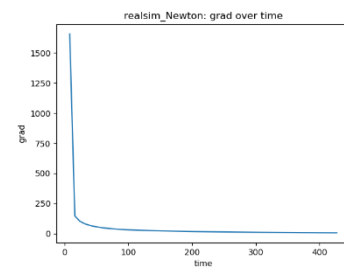
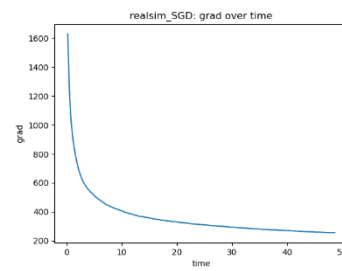
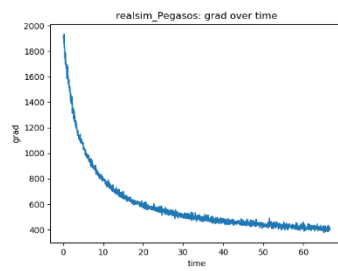
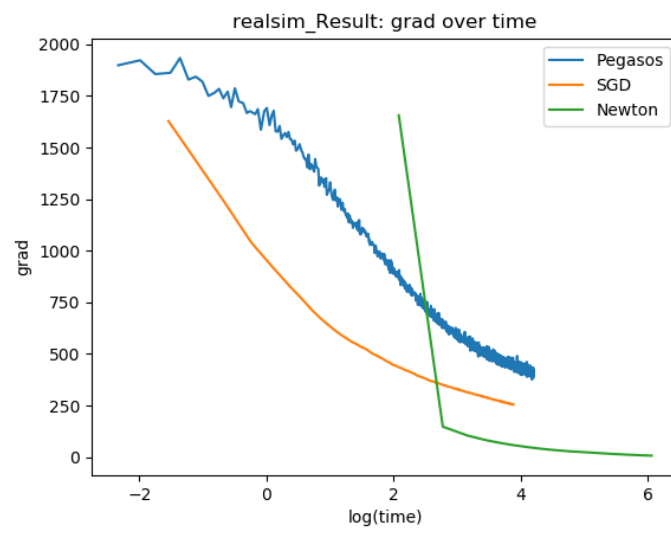


(2) [5 pts] Gradient norm versus training time

<contype>

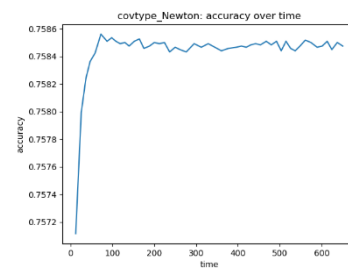
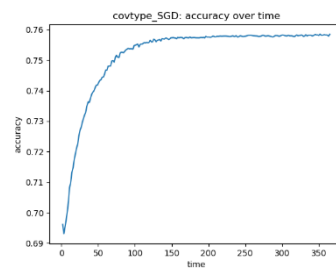
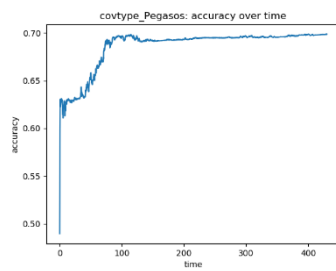
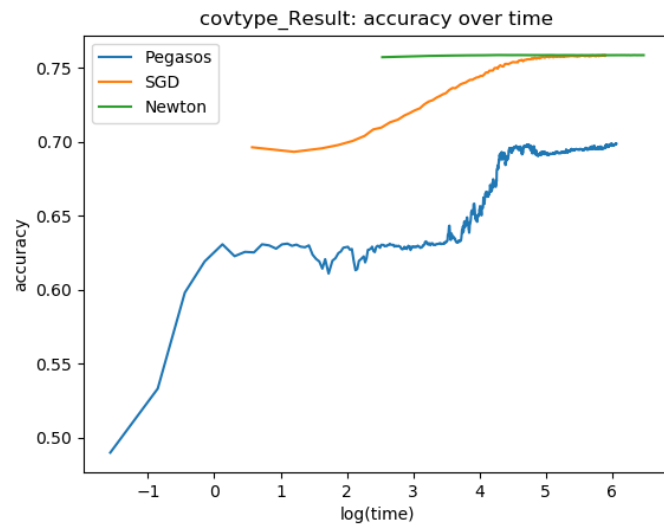


<realsim>

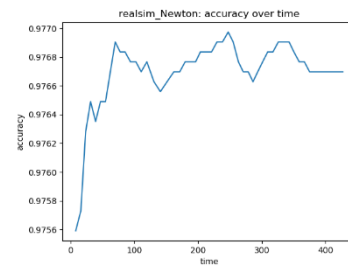
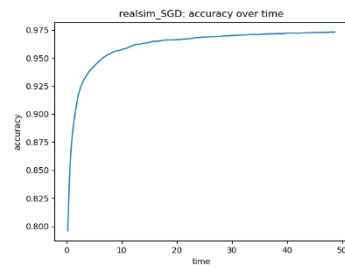
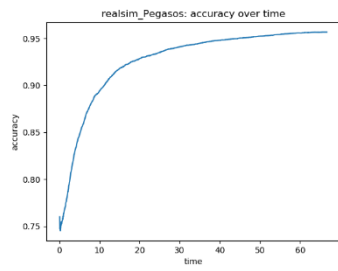
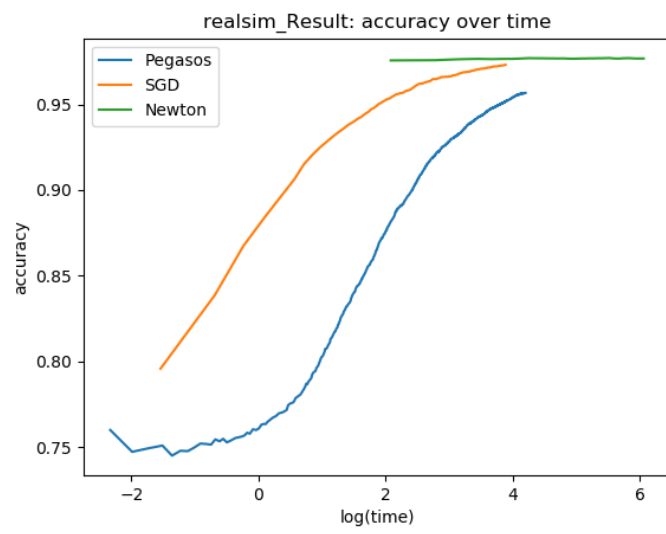




**(3) [5 pts] Test set accuracies versus training time**  
**<contype>**



<realsim>



(4) [5 pts] Discuss the difference between mini-batch SGD and Newton method in terms of the three types of figures

- All methods work well and almost reach the target values. Mini-batch SGD requires more iteration to converge. Appropriate stopping rule is needed.
- We use the conjugate gradient method in the Newton's method. Therefore, Newton's method is converging faster than mini batch gradient decent or SGD.
- In the mini-batch gradient, there are fluctuation in the gradient value because of the variation of each batches. Annealing the learning rate help to converge comparatively stable.
- Newton's method is a second order optimization technique and it can sometimes outperform than stochastic gradient decent. In some cases, it requires more computation cost. So, appropriate method should be considered.