

## Homework 3: Due Friday Dec. 1, 11:59 PM

**Instructions:** upload a PDF report using L<sup>A</sup>T<sub>E</sub>X containing your answers to Canvas (remember to include your name and ID number).

$w \in \mathbb{R}^{dx1}$

## Problem 1. Proximal Gradient Descent

Consider solving the following problem

$$\min_{\mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1,$$

where  $X \in \mathbb{R}^{n \times d}$  is the feature matrix (each row is a feature vector),  $\mathbf{y} \in \mathbb{R}^n$  is the label vector,  $\|\mathbf{w}\|_1 := \sum_i |w_i|$  and  $\lambda > 0$  is a constant to balance loss and regularization. This is known as the Lasso regression problem and our goal is to derive the “proximal gradient method” for solving this.

- (10 pt) The gradient descent algorithm cannot be directly applied since the objective function is non-differentiable. Discuss why the objective function is non-differentiable.
- (30 pt) In the class we showed that gradient descent is based on the idea of function approximation. To form an approximation for non-differentiable function, we split the differentiable part and non-differentiable part. Let  $g(\mathbf{w}) = \|X\mathbf{w} - \mathbf{y}\|_2^2$ , as discussed in the gradient descent lecture we approximate  $g(\mathbf{w})$  by

$$g(\mathbf{w}) \approx \hat{g}(\mathbf{w}) := g(\mathbf{w}_t) + \nabla g(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_t) + \frac{\eta}{2} \|\mathbf{w} - \mathbf{w}_t\|^2.$$

strong  
correct

In each iteration of proximal gradient descent, we obtain the next iterate ( $\mathbf{w}_{t+1}$ ) by minimizing the following approximation function:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \left( \hat{g}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \right)$$

Derive the close form solution of  $\mathbf{w}_{t+1}$  given  $\mathbf{w}_t, \nabla g(\mathbf{w}_t), \eta, \lambda$ .

## Problem 2. (60 pt) Implementing LSTM

In this homework, you are asked to implement a sequence model with LSTM to conduct sentiment analysis on SST-2 dataset with positive and negative sentiments (a binary classification problem). You need to finish the pipeline for training and evaluating the model. We provide a skeleton code for data loading and iterations of training data. You are asked to implement the rest of training in Pytorch code. Detailed submission requirements are written in the final section.

You can follow the setup instructions at <https://pytorch.org/get-started/locally/>. A useful tutorial on learning building LSTM at [https://pytorch.org/tutorials/beginner/nlp/sequence\\_models\\_tutorial.html](https://pytorch.org/tutorials/beginner/nlp/sequence_models_tutorial.html). Details of LSTM can be found here: <https://pytorch.org/docs/stable/generated/torch.nn.LSTM.html>.

We use SST-2 for sentiment analysis. SST-2 has only two sentiments, positive and negative. Pytorch/torch-text has provide a useful dataloader to automatically download and load the data into batches. We have written the data loader for you as follow. You can find it in the attached file.

Write a report containing your experiment results. Some requirements for your report:

- You are asked to implement a sequence model with a **bidirectional** LSTM. At most two layers of LSTM can be used. The model accuracy should be around 80%.
- Describe your model structure including vocabulary size and embedding dimension for the embedding layer; input size, hidden size, number of layers for LSTM layer; drouout layer if employed; input dim and out dim for fully connected layer.
- For each of the model, report the  $(\sum_{b=1}^B \sum_{d=1}^{D_b} \frac{\text{loss}(\text{label}_{b,d}, f_b(\text{data}_{b,d}))}{D_b})/B$  for each training epoch, where B is the total number of batches,  $f_b$  is the model after updated by b-th batch and  $D_b$  is the number of data points in b-th batch. An epoch is defined as one iteration of all dataset. Essentially, during a training epoch, you record down the average training loss of that batch after you update the model, and then report the average of all such batch-averaged losses after one iteration of whole dataset. You could plot the results as a figure our simply list down. Please at least report 10 epochs.

- Report the final testing accuracy of trained model.

Also, upload your code in a zip file and show how to run your code in README.

$$W_{t+1} = \underset{W}{\operatorname{argmin}} \underbrace{g(W) + \lambda \|W\|_1}_{\downarrow}$$

$$W_{t+1} = \underset{W}{\operatorname{argmin}} \left( g(W) + \lambda \|W\|_1 \right)$$

$$= \underset{W}{\operatorname{argmin}} \left( \underbrace{g(W_t) + \nabla g(W_t)^T (W - W_t) + \frac{\eta}{2} \|W - W_t\|_2^2}_{\text{Taylor expansion}} + \lambda \|W\|_1 \right)$$

$$g(W_t) + \nabla g(W_t)^T (W - W_t) + \frac{\eta}{2} \|W - W_t\|_2^2$$

$$W_{t+1} = \underset{W}{\operatorname{argmin}} \left( g(W) + \lambda \|W\|_1 \right)$$

$$f(x) = g(x) + h(x)$$

$$g(x) = g(x^k) + \nabla g(x^k)^T (x - x^k) + \frac{1}{2\alpha} \|x - x^k\|_2^2$$

$$h(x) = \lambda \|x\|_1$$

$$\frac{1}{2\alpha} = \frac{\eta}{2}$$

$$g(x) = g(x^k) + \nabla g(x^k)^T (x - x^k) + \frac{\eta}{2} \|x - x^k\|_2^2$$

$$= \frac{\eta}{2} \left( \frac{2}{\eta} g(x^k) + \frac{2}{\eta} \nabla g(x^k)^T (x - x^k) + \|x - x^k\|_2^2 \right) = \frac{\eta}{2} \left( \frac{2}{\eta} g(x^k) + \frac{2}{\eta} \nabla g(x^k)^T (x - x^k) + \|x - x^k\|_2^2 \right)$$

$$= \frac{\eta}{2} \left[ \frac{2}{\eta} g(x^k) + \left( 2 \frac{1}{\eta} \nabla g(x^k)^T (x - x^k) + \|x - x^k\|_2^2 + \left\| \frac{1}{\eta} \nabla g(x^k) \right\|_2^2 \right) - \left\| \frac{1}{\eta} \nabla g(x^k) \right\|_2^2 \right]$$

$$= \frac{\eta}{2} \left[ \frac{2}{\eta} g(x^k) + \left\| x - x^k + \frac{1}{\eta} \nabla g(x^k) \right\|_2^2 - \left\| \frac{1}{\eta} \nabla g(x^k) \right\|_2^2 \right]$$

$$= g(x^k) + \frac{\eta}{2} \left\| x - x^k + \frac{1}{\eta} \nabla g(x^k) \right\|_2^2 - \frac{\eta}{2} \left\| \frac{1}{\eta} \nabla g(x^k) \right\|_2^2$$

$$= g(x^k) + \frac{\eta}{2} \left\| x - \left( x^k - \frac{1}{\eta} \nabla g(x^k) \right) \right\|_2^2 - \frac{\eta}{2} \left\| \frac{1}{\eta} \nabla g(x^k) \right\|_2^2$$

$$= \underbrace{g(x^k)}_{\text{not related to } x} + \frac{\eta}{2} \left\| x - \hat{x} \right\|_2^2 - \underbrace{\frac{\eta}{2} \left\| \frac{1}{\eta} \nabla g(x^k) \right\|_2^2}_{\text{not related to } x}$$

not related to  $x$

not related to  $x$

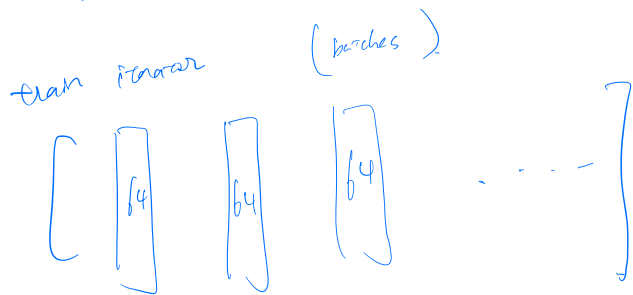
$\therefore$  minimize  $g(x) + h(x)$  is equivalent to

$$\underset{x}{\operatorname{minimize}} \quad \frac{1}{2\alpha} \|x - \hat{x}\|_2^2 + \lambda \|x\|_1$$

$$W = W_t + \frac{1}{\eta} \nabla g(W_t)$$

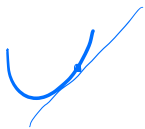
$$W = \left( W_t - \frac{1}{\eta} \nabla g(W_t) \right)$$

# LSTM



convex i.f.f.

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) \quad \forall x, y \in \mathbb{R}^n$$



$$\text{minimize}_x \quad \frac{1}{2\alpha} \|x - \hat{x}\|_2^2 + \lambda \|x\|_1$$

$$x_{k+1} = \arg\min_x \quad \frac{1}{2\alpha} \|x - \hat{x}\|_2^2 + \lambda \|x\|_1$$

proximal mapping

$$\text{prox}_{\alpha f}(x^k) = \arg\min_x \left( f(x) + \frac{1}{2\alpha} \|x - x^k\|_2^2 \right)$$

Soft-thresholding operator

$$S_\lambda(\hat{x}) = \begin{cases} \hat{x} - \lambda \\ 0 \\ \hat{x} + \lambda \end{cases}$$

$$h(x) = \lambda \|x\|_1$$

$$\text{prox}_{\alpha f}(x^k) = \arg\min_x \lambda \|x\|_1 + \frac{1}{2\alpha} \|x - \hat{x}^k\|_2^2$$

$$\Rightarrow \text{prox}_{\alpha f}(x^k) = \arg\min_x \lambda \sum_{i=1}^d |x_i| + \frac{1}{2\alpha} \sum_{i=1}^n (x_i - \hat{x}_i^k)^2$$

$$\Rightarrow \text{element-wise: } \text{prox}_{\alpha f}(x_i^k) =$$

$$\text{prox}_{\lambda f}(x_i^k) = \arg \min_{x_i} \lambda |x_i| + \frac{1}{2\alpha} (x_i - x_i^k)^2$$

$$0 \in \partial \left( \lambda |x_i| + \frac{1}{2\alpha} (x_i - \hat{x}_i)^2 \right)$$

你微分的這坨東西要 可以 = 0

optimality condition

$$0 \in \left\{ \frac{1}{2} (x_i - \hat{x}_i) + \lambda \partial |x_i| \right\}$$

$$x_i^* = \text{prox}_{\lambda f}(x_i^k)$$

$$\text{i.e.} \quad 0 \in \underbrace{\left\{ \frac{1}{\alpha} (x_i - \hat{x}_i) + \lambda \partial |x_i| \right\}}_{\text{gradient}}$$

$x_i$  condition

gradient

$$0 \in \left\{ \begin{array}{ll} x_i > 0 & \frac{1}{\alpha} (x_i - \hat{x}_i) + \lambda \\ x_i = 0 & \begin{array}{l} \frac{1}{\alpha} (x_i + \hat{x}_i) + \lambda \\ \frac{1}{\alpha} (x_i + \hat{x}_i) + 0.9999 \dots \lambda \\ \frac{1}{\alpha} (x_i + \hat{x}_i) + 0.0000 \dots \lambda \\ \frac{1}{\alpha} (x_i + \hat{x}_i) + 0 \end{array} \\ x_i < 0 & \frac{1}{\alpha} (x_i - \hat{x}_i) - \lambda \end{array} \right\}$$

$$\theta = \eta (w_i - w_i') + \lambda v$$

$$\theta = w_i - w_i' + \frac{1}{\eta} \lambda v$$

$$w_i = w_i' - \frac{1}{\eta} \lambda V \quad V = \begin{cases} 1 & \text{if } w_i > 0 \\ [-1, 1] & \text{if } w_i = 0 \\ -1 & \text{if } w_i < 0 \end{cases}$$

$$w_i \begin{cases} w_i' - \frac{1}{\eta} \lambda & , \quad w_i' > \frac{1}{\eta} \lambda \\ 0 & , \quad |w_i| \leq \frac{1}{\eta} \lambda \\ \underbrace{w_i' + \frac{1}{\eta} \lambda}_{} & , \quad w_i' < -\frac{1}{\eta} \lambda \end{cases}$$

$$w_i' + \frac{1}{\eta} \lambda < 0$$