

Weekly Exercises 5

Room: 02.09.023

Wednesday, 30.05.2018, 12:15-14:00

Submission deadline: Monday, 28.05.2018, 16:15, Room 02.09.023

Convex conjugate and prox (Due: 28.05) (8+4 Points)

Exercise 1 (4 points). Consider following problems of convex conjugate:

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex. Show that the convex conjugate of the perspective function $g : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$

$$g(x, t) = \begin{cases} tf(x/t), & \text{if } t > 0 \\ \infty, & \text{otherwise} \end{cases}$$

is given by

$$g^*(y, s) = \begin{cases} 0, & \text{if } f^*(y) \leq -s \\ \infty, & \text{otherwise} \end{cases}$$

- Show that the biconjugate of the perspective function g is given by

$$g^{**}(x, t) = \begin{cases} tf(x/t), & \text{if } t > 0 \\ \sigma_{\text{dom}(f^*)}(x), & \text{if } t = 0 \\ \infty, & \text{if } t < 0 \end{cases}$$

where $\sigma_{\text{dom}(f^*)}(x) = \sup_{y \in \text{dom}(f^*)} \langle x, y \rangle$ is the *support function* of $\text{dom}(f^*)$.

Solution. •

$$\begin{aligned} g^*(y, s) &= \sup_{x \in \mathbb{R}^n, t \in \mathbb{R}, t > 0} \langle x, y \rangle + st - tf(x/t) \\ &\stackrel{\xi=x/t}{=} \sup_{\xi \in \mathbb{R}^n, t \in \mathbb{R}, t > 0} t \langle \xi, y \rangle + st - tf(\xi) \\ &= \sup_{t > 0} t \cdot \left(\left[\sup_{\xi} \langle y, \xi \rangle - f(\xi) \right] + s \right) \\ &= \sup_{t > 0} t \cdot (f^*(y) + s) = \begin{cases} 0 & \text{if } f^*(y) + s \leq 0 \\ \infty & \text{otherwise.} \end{cases} \end{aligned}$$

•

$$g^{**}(x, t) = \sup_{f^*(y) \leq -s} \langle y, x \rangle + st$$

For $t < 0$ the supremum is unbounded, since for $s \rightarrow -\infty$ we have $st \rightarrow \infty$.

For $t = 0$ we have

$$g^{**}(x, t) = \sup_{f^*(y) \leq -s} \langle y, x \rangle = \sup_{y \in \text{dom}(f^*)} \langle y, x \rangle = \sigma_{\text{dom}(f^*)}(x)$$

Finally let $t > 0$. The supremum in s is achieved at $\hat{s} = -f^*(y)$ and we have

$$\begin{aligned} g^{**}(x, t) &= \sup_{f^*(y) \leq -s} \langle y, x \rangle + st = \sup_y \langle y, x \rangle - tf^*(y) \\ &= t \sup_y \langle y, x/t \rangle - f^*(y) = tf^{**}(x/t) = tf(x/t). \end{aligned}$$

The last equality holds since f is convex (and is everywhere finite hence continuous and proper) we have $f^{**} = f$.

Exercise 2 (4 points). Compute the proximity operator of the 1, 2-norm, i.e.

$$\text{prox}_{\tau \|X\|_{1,2}},$$

where $X \in \mathbb{R}^{m \times n}$ is a matrix .

Solution. Firstly recall the subdifferential of 1, 2-norm computed in previous sheet:

$$\partial \|X\|_{1,2} = \{P \in \mathbb{R}^{m \times n} : P_i \in \partial \|X_i\|_2\}$$

where P_i and X_i are the i -th row of corresponding matrix and

$$\partial \|X_i\| = \begin{cases} \frac{X_i}{\|X_i\|_2}, & X_i \neq 0 \\ \{P_i \in \mathbb{R}^n : \|P_i\|_2 \leq 1\}, & X_i = 0 \end{cases}$$

Now we use the definition of proximity operator:

$$\text{prox}_{\tau \|X\|_{1,2}}(Y) = \text{argmin}_X \|X\|_{1,2} + \frac{1}{2\tau} \|X - Y\|_2^2$$

which gives us the following by using the optimality condition:

$$0 \in \partial \|X\|_{1,2} + \frac{1}{\tau}(X - Y).$$

Since each row is independently, we can solve it for each row and get:

$$0 \in \partial \|X_i\|_2 + \frac{1}{\tau}(X_i - Y_i).$$

If $\|X_i\| \neq 0$, we have $Y_i = \tau \frac{X_i}{\|X_i\|} + X_i$. If we denote $X_i = te_i$ where $e_i := \frac{X_i}{\|X_i\|}$, previous equation becomes $Y_i = \tau e_i + te_i$. Hence, $\|Y_i\|_2 = \tau + t$, which implies $\|Y_i\| > \tau$.

If $\|X_i\| = 0$, we have $Y_i \in \{P_i \in \mathbb{R}^n : \|P_i\|_2 \leq \tau\}$.

To summary we have:

$$\text{prox}_{\tau\|X\|_{1,2}}(Y) = \left\{ X \in \mathbb{R}^{m \times n} : X_i = \begin{cases} 0, & \text{if } \|Y_i\| \leq \tau \\ (\|Y_i\|_2 - \tau) \frac{Y_i}{\|Y_i\|}, & \text{if } \|Y_i\| > \tau \end{cases} \right\}.$$

Exercise 3 (4 points). Prove that the proximal operator of the nuclear norm is the proximal operator of the ℓ_1 -norm applied to the singular values of the input argument. Formally, let $Y \in \mathbb{R}^{n \times n}$ and let $Y = U\Sigma V^\top$ be the singular value decomposition of Y . Prove that

$$\text{prox}_{\tau\|\cdot\|_{\text{nuc}}}(Y) = U \text{diag}(\{(\sigma_i - \tau)_+\}) V^\top,$$

where $\text{diag}(\{(\sigma_i - \tau)_+\}) := \text{diag}(\{\max\{0, \sigma_i - \tau\}\}) = \text{prox}_{\tau\|\cdot\|_1}(\{\sigma_i\})$ is the shrinkage (or soft thresholding) operator applied to the singular values σ_i of Y .

Solution. Let $Y \in \mathbb{R}^{n \times n}$. We are interested in the solution of

$$\text{argmin}_X \frac{1}{2} \|X - Y\|_F^2 + \tau \|X\|_{\text{nuc}}.$$

Since the above problem is strictly convex there exists a unique solution \hat{X} . The optimality condition of the problem is given as

$$0 \in \hat{X} - Y + \partial\|\cdot\|_{\text{nuc}}(\hat{X}). \quad (1)$$

where $\partial\|\cdot\|_{\text{nuc}}(X)$ is the subdifferential of the nuclear norm at X characterized on exercise sheet 3. Our aim is to show that $\hat{X} := U \text{diag}(\{(\sigma_i - \tau)_+\}) V^\top$ meets the optimality condition. To this end we decompose $V = [V_1 \ V_2]$, $U = [U_1 \ U_2]$ and $\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix}$ so that

$$Y = U_1 \Sigma_1 V_1^\top + U_2 \Sigma_2 V_2^\top,$$

where Σ_1 contains all singular values $\sigma_i > \tau$ and Σ_2 all singular values $\sigma_i \leq \tau$. We may then write \hat{X} as

$$\hat{X} = U \text{diag}(\{(\sigma_i - \tau)_+\}) V^\top = U_1 \underbrace{(\Sigma_1 - \tau I)}_{\sigma_i > 0} V_1^\top + U_2 \underbrace{\text{diag}(\{0\})}_{\sigma_i = 0} V_2^\top.$$

We will now show that \hat{X} meets (1): $Y - \hat{X}$ is given as

$$Y - \hat{X} = \tau(U_1 V_1^\top + U_2 \frac{1}{\tau} \Sigma_2 V_2^\top).$$

By construction $\|\frac{1}{\tau} \Sigma_2\|_{\text{spec}} \leq 1$. And therefore and due to sheet 3

$$Y - \hat{X} \in \tau \partial\|\cdot\|_{\text{nuc}}(\hat{X})$$

Multinomial Logistic Regression(Due:28.05) (16 Points)

Exercise 4 (16 Points). In this exercise you are asked to train a linear model for a multiclass classification task with Logistic regression. The idea is as follows: You are given a set of training samples $\mathcal{I} = \{1, \dots, N\}$ that are represented by their feature vectors $x_i \in \mathbb{R}^d$, for $i \in \mathcal{I}$. Each training sample i is associated with a class label $y_i \in \{1, \dots, C\}$. The aim is to estimate a linear classifier parameterized by $W^* \in \mathbb{R}^{d \times C}$, $b^* \in \mathbb{R}^C$ so that $y_i = \operatorname{argmax}_{1 \leq j \leq C} x_i^\top W_j^* + b_j^*$ for most training samples i . Once you have obtained this “optimal” classifier the hope is, that you are able to classify new unseen and unlabeled samples $x \in \mathbb{R}^d$. In machine learning this is called generalization. For this task you may query your trained model via the classifier rule

$$y = \operatorname{argmax}_{1 \leq j \leq C} x^\top W_j^* + b_j^* \quad (2)$$

and y probably is the true class label of x if your model generalizes well.

In order to estimate the model we solve an optimization problem of the form

$$\min_{W \in \mathbb{R}^{d \times C}, b \in \mathbb{R}^C} \frac{1}{N} \sum_{i=1}^N \ell(W, b, x_i, y_i) + \frac{\lambda_1}{2} \|W\|_2^2 + \frac{\lambda_2}{2} \|b\|_2^2, \quad (3)$$

where

$$\ell(W, b, x_i, y_i) = -\log \left(\frac{\exp(\langle W_{y_i}, x_i \rangle + b_{y_i})}{\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j)} \right) \quad (4)$$

is called the softmax loss. Note that the above problem is smooth and strongly convex and can be solved with gradient descent. In practice however, it may happen, that some features (i.e. components of the vector x_i) do not contain any information about the true class labels, i.e. components that are just noise. In order to filter out the useless features we modify the norm on W . So we have

$$\min_{W \in \mathbb{R}^{d \times C}, b \in \mathbb{R}^C} \frac{1}{N} \sum_{i=1}^N \ell(W, b, x_i, y_i) + \frac{\lambda_1}{2} \|W\|_{1,2} + \frac{\lambda_2}{2} \|b\|_2^2 \quad (5)$$

You are asked to do the following:

- Download the toy data template from the homepage
- Implement a proximal gradient descent algorithm to optimize above objective function (5) (Avoid for-loops)
- Make sure that your objective monotonically decreases. Plot the objective values. Stop your code if the difference of two successive iterates is less than 10^{-12} .

- In order to ensure that your derivative is computed correctly you may first optimize the fully differentiable model (3) with MATLABs *fminunc* with the options '*GradObj*', '*On*' and '*DerivativeCheck*', '*On*'.
- Iteratively compute the test error in percent, i.e. how many test samples are not classified correctly via the rule (2).
- Play around with different parameter settings for λ_1, λ_2 . Can you identify the useless features? Explain why the model generalizes better to unseen test data if you use 1, 2-norm on W^* (answer by comment at the end of your code).
- You may apply your code to the MNIST dataset <http://yann.lecun.com/exdb/mnist/> and see that your are now able to classify handwritten digits.

Solution. We apply the proximal gradient descent scheme to our objective (5). To this end we need compute the partial derivatives $\frac{\partial F(W,b)}{\partial W_{lk}}$ and $\frac{\partial F(W,b)}{\partial b_k}$ of the differentiable part of the objective

$$F(W, b) = \frac{1}{N} \sum_{i=1}^N \ell(W, b, x_i, y_i) + \frac{\lambda_1}{2} \|W\|_2^2 + \frac{\lambda_1}{2} \|b\|_2^2.$$

First we observe, that

$$\frac{\partial F(W, b)}{\partial W_{lk}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(W, b, x_i, y_i)}{\partial W_{lk}} + \lambda_1 W_{lk}$$

and

$$\frac{\partial F(W, b)}{\partial b_k} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(W, b, x_i, y_i)}{\partial b_k} + \lambda_1 b_k.$$

For some class $1 \leq k \leq C$ define

$$h_k(W, b) = \frac{\exp(\langle W_{y_i}, x_i \rangle + b_{y_i})}{\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j)}$$

and

$$\mathbf{1}\{y_i = k\} = \begin{cases} 1 & \text{if } y_i = k \\ 0 & \text{otherwise.} \end{cases}$$

Via the one-dimensional chain rule and the quotient rule the partial derivatives of

the individual loss terms are given as:

$$\begin{aligned}
& \frac{\partial \ell(W, b, x_i, y_i)}{\partial W_{lk}} \\
&= - \frac{1}{h_{y_i}(W, b)} \cdot \frac{\mathbf{1}\{y_i = k\} \cdot \exp(\langle W_{y_i}, x_i \rangle + b_{y_i}) \cdot x_{il} \cdot \left(\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j) \right)}{\left(\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j) \right)^2} \\
&= + \frac{1}{h_{y_i}(W, b)} \cdot \frac{\exp(\langle W_{y_i}, x_i \rangle + b_{y_i}) \cdot \exp(\langle W_k, x_i \rangle + b_k) \cdot x_{il}}{\left(\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j) \right)^2} \\
&= - \frac{1}{h_{y_i}(W, b)} \cdot \mathbf{1}\{y_i = k\} \cdot x_{il} \cdot h_{y_i}(W, b) + \frac{1}{h_{y_i}(W, b)} \cdot h_{y_i}(W, b) \cdot h_k(W, b) \cdot x_{il} \\
&= (h_k(W, b) - \mathbf{1}\{y_i = k\}) \cdot x_{il}.
\end{aligned}$$

Similarly we obtain for the derivative wrt. b_k :

$$\begin{aligned}
& \frac{\partial \ell(W, b, x_i, y_i)}{\partial b_k} \\
&= - \frac{1}{h_{y_i}(W, b)} \cdot \frac{\mathbf{1}\{y_i = k\} \cdot \exp(\langle W_{y_i}, x_i \rangle + b_{y_i}) \cdot \left(\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j) \right)}{\left(\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j) \right)^2} \\
&= + \frac{1}{h_{y_i}(W, b)} \cdot \frac{\exp(\langle W_{y_i}, x_i \rangle + b_{y_i}) \cdot \exp(\langle W_k, x_i \rangle + b_k)}{\left(\sum_{j=1}^C \exp(\langle W_j, x_i \rangle + b_j) \right)^2} \\
&= h_k(W, b) - \mathbf{1}\{y_i = k\}.
\end{aligned}$$