

## Machine Learning Homework Sheet 06

### Constrained Optimization and SVM

## 1 Constrained Optimization

**Problem 1:** Solve the following constrained optimization problem using the recipe described in the lecture (slide 17).

$$\begin{aligned} &\text{minimize} && f_0(\boldsymbol{\theta}) = \theta_1 - \sqrt{3}\theta_2 \\ &\text{subject to} && f_1(\boldsymbol{\theta}) = \theta_1^2 + \theta_2^2 - 4 \leq 0 \end{aligned}$$

Write down the **Lagrangian**

$$L(\boldsymbol{\theta}, \alpha) = \theta_1 - \sqrt{3}\theta_2 + \alpha (\theta_1^2 + \theta_2^2 - 4) .$$

Calculate the derivative of  $L(\boldsymbol{\theta}, \alpha)$  w.r.t.  $\boldsymbol{\theta}$  and set it to zero,

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \alpha) &= \begin{pmatrix} 1 + 2\alpha\theta_1 \\ -\sqrt{3} + 2\alpha\theta_2 \end{pmatrix} = 0 \\ \Rightarrow \theta_1 &= -\frac{1}{2\alpha}, \quad \theta_2 = \frac{\sqrt{3}}{2\alpha} . \end{aligned}$$

Substituting this back into  $L(\boldsymbol{\theta}, \alpha)$  gives the **dual function**

$$g(\alpha) = -\frac{1}{2\alpha} - \frac{3}{2\alpha} + \frac{1}{4\alpha} + \frac{3}{4\alpha} - 4\alpha = -\frac{1}{\alpha} - 4\alpha .$$

We must now **maximize the dual function**  $g(\alpha)$  subject to  $\alpha \geq 0$ .

Since  $g(\alpha)$  is concave for  $\alpha \geq 0$ , we set its derivative to zero and solve for  $\alpha$ .

$$\begin{aligned} \frac{dg}{d\alpha} &= \frac{1}{\alpha^2} - 4 = 0 \\ \Rightarrow \alpha^2 &= \frac{1}{4} \\ \Rightarrow \alpha_{1,2} &= \pm \frac{1}{2} \\ \Rightarrow \alpha^* &= \frac{1}{2} \end{aligned}$$

since we require  $\alpha \geq 0$ . The optimization problem is convex and Slater's condition holds, therefore the minimal value of  $f_0(\boldsymbol{\theta})$  is

$$\mathbf{p}^* = g\left(\frac{1}{2}\right) = -4 .$$

We calculate the minimizer  $\boldsymbol{\theta}^*$  by substituting  $\alpha = 1/2$  into  $\theta_1 = -\frac{1}{2\alpha}$  and  $\theta_2 = \frac{\sqrt{3}}{2\alpha}$  and get

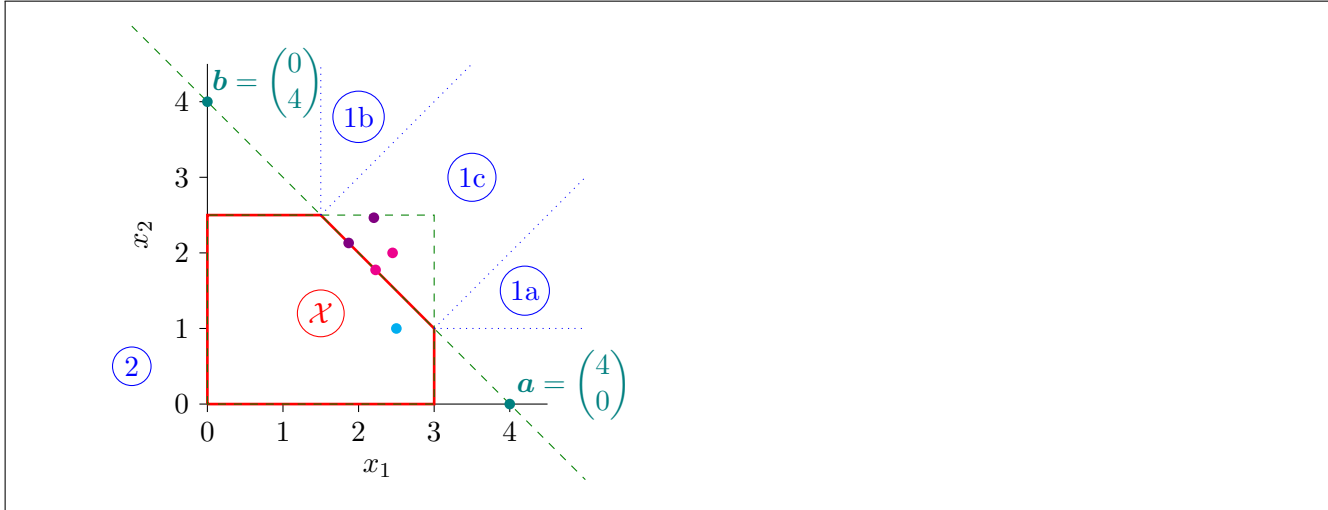
$$\theta_1^* = -1, \quad \theta_2^* = \sqrt{3} .$$

## 2 Projected Gradient Descent

**Problem 2:** Given is the following (convex) domain defined by a set of linear constraints

$$\mathcal{X} \subset \mathbb{R}^2 = \{\mathbf{x} \in \mathbb{R}^2 : (x_1 + x_2 \leq 4) \wedge (0 \leq x_1 \leq 3) \wedge (0 \leq x_2 \leq 2.5)\}.$$

a) Visualize the set  $\mathcal{X}$ .



b) Derive a closed form for the projection  $\pi_{\mathcal{X}}(\mathbf{p}) = \arg \min_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x} - \mathbf{p}\|_2^2$ . That is, given an arbitrary point  $\mathbf{p} \in \mathbb{R}^2$ , what is its projection on  $\mathcal{X}$ ?

*Hint:* For one part of  $\mathbb{R}^2$  you might want to use the line projection  $\pi_{\text{line}}(\mathbf{p}) = \mathbf{a} + \frac{(\mathbf{p}-\mathbf{a})^T(\mathbf{b}-\mathbf{a})}{\|\mathbf{b}-\mathbf{a}\|_2^2}(\mathbf{b}-\mathbf{a})$ , where  $\mathbf{a} \in \mathbb{R}^2$  and  $\mathbf{b} \in \mathbb{R}^2$  specify the line.

$$\pi_{\mathcal{X}}(\mathbf{p}) = \begin{cases} \mathbf{p} & \text{if } \mathbf{p} \in \mathcal{X} \\ (3, 1)^T & \text{if } \mathbf{p} \in \textcircled{1a} = \{\mathbf{p} | p_2 \geq 1 \wedge -p_1 + p_2 \leq -2\} \\ (1.5, 2.5)^T & \text{if } \mathbf{p} \in \textcircled{1b} = \{\mathbf{p} | p_1 \geq 1.5 \wedge -p_1 + p_2 \geq 1\} \\ \pi_{\text{line}}(\mathbf{p}) & \text{if } \mathbf{p} \in \textcircled{1c} = \{\mathbf{p} | -2 < -p_1 + p_2 < 1 \wedge p_1 + p_2 > 4\} \\ \pi_{\text{box}}(\mathbf{p}) & \text{if } \mathbf{p} \in \textcircled{2} = \{\mathbf{p} | \mathbf{p} \notin \mathcal{X} \wedge (p_1 < 1.5 \vee p_2 < 1)\} \end{cases}$$

$$\begin{aligned} \pi_{\text{line}}(\mathbf{p}) &= \mathbf{a} + \frac{(\mathbf{p}-\mathbf{a})^T(\mathbf{b}-\mathbf{a})}{\|\mathbf{b}-\mathbf{a}\|_2^2}(\mathbf{b}-\mathbf{a}) = \begin{pmatrix} 4 \\ 0 \end{pmatrix} + \frac{(p_1-4, p_2) \begin{pmatrix} -4 \\ 4 \end{pmatrix}}{32} \begin{pmatrix} -4 \\ 4 \end{pmatrix} \\ &= \begin{pmatrix} 4 \\ 0 \end{pmatrix} + \frac{16-4p_1+4p_2}{32} \begin{pmatrix} -4 \\ 4 \end{pmatrix} = \begin{pmatrix} 4 \\ 0 \end{pmatrix} + \left(2 - \frac{1}{2}p_1 + \frac{1}{2}p_2\right) \begin{pmatrix} -1 \\ 1 \end{pmatrix} \\ &= \begin{pmatrix} 2 + \frac{1}{2}p_1 - \frac{1}{2}p_2 \\ 2 - \frac{1}{2}p_1 + \frac{1}{2}p_2 \end{pmatrix} \\ \pi_{\text{box}}(\mathbf{p}) &= \begin{pmatrix} \max(0, \min(3, p_1)) \\ \max(0, \min(2.5, p_2)) \end{pmatrix} \end{aligned}$$

Upload a single PDF file with your solution to Moodle by 2.12.2018, 23:59pm CET. We recommend to typeset your solution (using L<sup>A</sup>T<sub>E</sub>X or Word), but handwritten solutions are also accepted.

If your handwritten solution is illegible, it won't be graded and you waive your right to dispute that.

c) Given is the following constrained optimization problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & (x_1 - 2)^2 + (2x_2 - 7)^2, \\ \text{subject to } & \mathbf{x} \in \mathcal{X}. \end{aligned}$$

Perform two steps of projected gradient descent starting from the point  $\mathbf{x}^{(0)} = (2.5, 1)^T$ . Use a constant learning rate/step size of  $\tau = 0.05$ .

Projected gradient descent is a two-step algorithm consisting of regular gradient descent and projection:

$$\begin{aligned} 1. \quad & \mathbf{p}^{t+1} = \mathbf{x}^t - \tau \nabla_{\mathbf{x}} f(\mathbf{x}^t) \\ 2. \quad & \mathbf{x}^{t+1} = \pi_{\mathcal{X}}(\mathbf{p}^{t+1}) \end{aligned}$$

The gradient of this function is

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} 2x_1 - 4 \\ 8x_2 - 28 \end{pmatrix}.$$

Using this gradient and the projection from above we can perform two steps of the algorithm:

$$\begin{aligned} 1.1 \quad & \mathbf{p}^{(1)} = \mathbf{x}^{(0)} - \tau \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) = \begin{pmatrix} 2.5 \\ 1 \end{pmatrix} - 0.05 \begin{pmatrix} 1 \\ -20 \end{pmatrix} = \begin{pmatrix} 2.45 \\ 2 \end{pmatrix} \\ 1.2 \quad & \mathbf{x}^{(1)} = \pi_{\mathcal{X}}(\mathbf{p}^{(1)}) = \pi_{\text{line}}\left(\begin{pmatrix} 2.45 \\ 2 \end{pmatrix}\right) = \begin{pmatrix} 2.225 \\ 1.775 \end{pmatrix} \\ 2.1 \quad & \mathbf{p}^{(2)} = \mathbf{x}^{(1)} - \tau \nabla_{\mathbf{x}} f(\mathbf{x}^{(1)}) = \begin{pmatrix} 2.225 \\ 1.775 \end{pmatrix} - 0.05 \begin{pmatrix} 0.45 \\ -13.8 \end{pmatrix} = \begin{pmatrix} 2.2025 \\ 2.465 \end{pmatrix} \\ 2.2 \quad & \mathbf{x}^{(2)} = \pi_{\mathcal{X}}(\mathbf{p}^{(2)}) = \pi_{\text{line}}\left(\begin{pmatrix} 2.2025 \\ 2.465 \end{pmatrix}\right) = \begin{pmatrix} 1.86875 \\ 2.13125 \end{pmatrix} \end{aligned}$$

This process is also illustrated in the figure above, with points showing the **starting point (cyan)**, **step 1 (magenta)** and **step 2 (violet)**.

### 3 SVM

**Problem 3:** Explain the similarities and differences between the SVM and perceptron algorithms.

Both algorithms are looking for a hyperplane that separates the two classes.

The difference is that SVM also tries to maximize the margin, while the perceptron only cares about separation.

**Problem 4:** Show that the duality gap is zero for SVM.

The objective is convex and the constraints are affine in  $\mathbf{w}$ , thus the Slater's condition is satisfied, so the duality gap is zero.

**Problem 5:** Recall that the dual function for SVM (slide 41) can be written as

$$g(\boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{1}_N$$

- (a) Show how the matrix  $\mathbf{Q}$  can be computed. (*Hint: You might want to use Hadamard product, denoted as  $\odot$* ).
- (b) Prove that the matrix  $\mathbf{Q}$  is negative (semi-)definite.
- (c) Explain what the negative (semi-)definiteness means for our optimization problem. Why is this property important?

(a)  $\mathbf{Q} = (-\mathbf{y}\mathbf{y}^T \odot \mathbf{X}\mathbf{X}^T)$

(b) We will first show that  $\mathbf{M} = (\mathbf{y}\mathbf{y}^T \odot \mathbf{X}\mathbf{X}^T) \in \mathbb{R}^{N \times N}$  is a positive (semi-)definite (PSD) matrix.

For  $\mathbf{M}$  to be PSD, for every  $\mathbf{a} \in \mathbb{R}^N$  it must hold that  $\mathbf{a}^T \mathbf{M} \mathbf{a} \geq 0$ .

$$\begin{aligned} \mathbf{a}^T \mathbf{M} \mathbf{a} &= \sum_{i=1}^N \sum_{j=1}^N a_i y_i y_j \mathbf{x}_i^T \mathbf{x}_j a_j \\ &= \sum_{i=1}^N \sum_{j=1}^N (a_i y_i \mathbf{x}_i)^T (a_j y_j \mathbf{x}_j) \\ &= (\mathbf{a} \odot \mathbf{y})^T \mathbf{X} \mathbf{X}^T (\mathbf{a} \odot \mathbf{y}) \\ &= (\mathbf{X}^T (\mathbf{a} \odot \mathbf{y}))^T (\mathbf{X}^T (\mathbf{a} \odot \mathbf{y})) \\ &= (\mathbf{X}^T (\mathbf{a} \odot \mathbf{y}))^2 \\ &\geq 0 \end{aligned}$$

As  $\mathbf{M}$  is PSD, it follows that  $\mathbf{Q} = -\mathbf{M}$  is negative (semi-)definite.

- (c) This means that our maximization problem is **concave**, thus every local maximum is a global maximum.

**Problem 6:** Download the notebook `homework_06_notebook.ipynb` from Piazza. Fill in the missing code and run the notebook. Convert the evaluated notebook to pdf and add it to the printout of your homework (see printing instructions inside the notebook).