

TC1 - Exam 2021-2022 - 2 hours and 30 minutes

Caio Corro

No document, no phone, no computer. You **must** use formal notation and you **must** define all mathematical terms you use.

1 Course questions (11 points)

- (1 point) Why do we want the training problem to be a convex or concave optimization problem?
- (2 points) What important property of convex functions can impact the optimization algorithm? Why?
- (1 point) What does "support vector" refer to in Support Vector Machines?
- (2 point) Why are slack variables introduced in the constrained formulation of SVM ? (you can use a picture to illustrate your explanation - for this question you are allowed to not be very formal, but the answer must be clear)
- (2 points) Why is it interesting to optimize the dual formulation of a SVM via the coordinate ascent algorithm? (at least 2 reasons)
- (2 points) Let $\mathbf{w} \in \mathbb{R}^k$ be a vector of scores and $\hat{\mathbf{y}} \in E(k)$ be a one-hot vector representing the gold output. The negative log-likelihood loss for the softmax prediction function can be written as:

$$l(\mathbf{w}; \hat{\mathbf{y}}) = -\log(\hat{\mathbf{y}}^\top \text{softmax}(\mathbf{w}))$$

Compute the gradient $\nabla_{\mathbf{w}} l(\mathbf{w}; \hat{\mathbf{y}})$. How can it be interpreted?

- (1 point) What could go wrong if we use the negative log-likelihood loss function when using sparsemax as a prediction function?

2 Convex analysis (24 points)

- (2 points) Let $f : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ be the negative logarithm function, i.e. $f(x) = -\log x$. Compute (explicitly) the Fenchel conjugate and the biconjugate of f . Is this result expected?
- (1 point) If we consider the function $h(x) = \log x$ instead, would we expect to have $h^{**}(x) = h(x)$? Why? (no need to compute anything for this question)
- (2 points) Let $f : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ be a function defined as $f(x) = ax + b$ where $a \in \mathbb{R}$ and $b \in \mathbb{R}$. Compute (explicitly) the Fenchel conjugate and biconjugate of f .
- (3 points) Let $f : X \rightarrow \mathbb{R}$ be a convex function where $X \subset \mathbb{R}^n$ is a convex subset of \mathbb{R}^n . Prove that its extended real-valued extension $\tilde{f} : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex.

$$\tilde{f}(x) = \begin{cases} f(x) & \text{if } x \in X \\ \infty & \text{otherwise} \end{cases}$$

- (2 points) Why is it useful to rely on extended real-valued extension of functions? Why is it important that $\tilde{f}(x)$ is convex? (you can give an example)
- (2 points) Let $f(x) = |x|$. Prove that g defined as follows is a subgradient of f at x :

$$g = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \\ 0 & \text{otherwise} \end{cases}$$

7. (2 points) Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function defined as follows:

$$f(\mathbf{x}) = \sum_i f_i(x_i)$$

where f_1, \dots, f_n are convex functions defined as $f_i : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$. Prove that the conjugate of f can be written as:

$$f^*(\mathbf{y}) = \sum_i f_i^*(y_i)$$

8. (1 points) In which situation is the result of the previous function useful in Machine Learning? (one example, from the course)

9. (2 points) Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function defined as follows:

$$f(\mathbf{x}) = \alpha h(\mathbf{x}) \quad \text{with} \quad \alpha > 0$$

Prove that $f^*(\mathbf{y}) = \alpha h^*(\mathbf{y}/\alpha)$.

10. (1 points) In which situation is the result of the previous function useful in Machine Learning? (one example, from the course)

11. (2 points) Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. Let $\mathbf{y} \in \text{dom } f^*$. Prove that $\mathbf{g} = \arg \max_{\mathbf{x} \in \text{dom } f} \mathbf{y}^\top \mathbf{x} - f(\mathbf{x})$ is a subgradient of $f^*(\mathbf{y})$.

12. (1 points) In which situation is the result of the previous function useful in Machine Learning? (one example, from the course)

13. (2 points) Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be the function defined as follows:

$$f(w) = \max(0, 1 - w)$$

Compute the Fenchel conjugate of f .

14. (1 points) In which situation is the result of the previous function useful in Machine Learning? (one example, from the course)

3 Fenchel-Young losses (24 points)

Note: for questions 3 and 4, you are expected to describe both the regularized prediction function \hat{y}_Ω and the associated Fenchel-Young loss L_Ω .

1. (4 points) Let $\Omega : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex regularizer with $\text{dom } \Omega \subseteq \Delta^k$. Prove that \hat{y}_Ω is invariant to constants, that is:

$$\forall w \in \mathbb{R}^k, c \in \mathbb{R} : \quad \hat{y}_\Omega(\mathbf{w}) = \hat{y}_\Omega(\mathbf{w} + c)$$

2. Let $\Omega : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex regularizer. Prove that the associated Fenchel-Young loss is always non-negative, that is:

$$\forall \mathbf{w} \in \mathbb{R}^k, \mathbf{y} \in E(k) : \quad L_\Omega(\mathbf{w}; \mathbf{y}) \geq 0$$

3. (4 points) (8 points) We assume the following scoring and prediction function for regression:

$$\begin{aligned} \hat{w}(\mathbf{x}) &= \mathbf{a}^\top \mathbf{x} & \hat{w} : \mathbb{R}^n &\rightarrow \mathbb{R} \\ \hat{y}(w) &= w & \hat{y} : \mathbb{R} &\rightarrow \mathbb{R} \end{aligned}$$

where $\mathbf{a} \in \mathbb{R}^n$ are the parameters of the model. The standard training loss is the squared loss defined as $l(w; \hat{y}) = \frac{1}{2} \|\hat{y} - w\|_2^2$ where \hat{y} is the gold output and w is the output of the weighting function. Prove formally that we can express this prediction and loss function in the Fenchel-Young loss framework.

4. (8 points) We assume the following scoring and prediction for probabilistic binary classification:

$$\begin{aligned} \hat{w}(\mathbf{x}) &= \mathbf{a}^\top \mathbf{x} & \hat{w} : \mathbb{R}^n &\rightarrow \mathbb{R} \\ \hat{y}(w) &= \sigma(w) & \hat{y} : \mathbb{R} &\rightarrow]0, 1[\end{aligned}$$

where $\sigma(w) = \frac{\exp(w)}{1 + \exp(w)}$ is the sigmoid function and $\mathbf{a} \in \mathbb{R}^n$ are the parameters of the model. In other words, the function $\hat{y}(w)$ returns the parameter μ of a Bernoulli distribution defined as follows:

$$p(y; \mu) = \mu^y (1 - \mu)^{(1-y)}$$

The entropy of a Bernoulli distribution is defined as follows:

$$H[p(\cdot; \mu)] = -y \log y - (1 - y) \log(1 - y)$$

Build formally the regularized prediction function and the Fenchel-Young loss associated with this probabilistic binary classifier and derive their closed form expressions (as you would implement them).

Cheat sheet

4 KKT conditions

Assume we have a maximization problem defined as follows:

$$\begin{aligned} \max_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g^{(i)}(\mathbf{x}) \geq 0 && \forall 1 \leq i \leq m \\ & h^{(i)}(\mathbf{x}) = 0 && \forall 1 \leq i \leq n \end{aligned}$$

where $\mathbf{x} \in \mathbb{R}^k$, $g^{(i)} \geq 0$ is a set of m inequality and $h^{(i)}(\mathbf{x}) = 0$ is a set of n equality. An optimal solution $\mathbf{x}^* \in \mathbb{R}^k$ of the mathematical program satisfies the following constraints:

$$\begin{aligned} (\text{stationarity}) \quad & \forall i : \frac{\partial}{\partial x_i} f(\mathbf{x}^*) + \sum_j \mu_j \frac{\partial}{\partial x_i} g^{(j)}(\mathbf{x}^*) - \sum_j \lambda_j \frac{\partial}{\partial x_i} h^{(j)}(\mathbf{x}^*) = 0 \\ (\text{primal feasibility}) \quad & \forall i : g^{(i)}(\mathbf{x}^*) \geq 0 \\ & \forall i : h^{(i)}(\mathbf{x}^*) = 0 \\ (\text{dual feasibility}) \quad & \forall i : \mu_i \geq 0 \\ (\text{complementary slackness}) \quad & \sum_i \mu_i g^{(i)}(\mathbf{x}^*) \geq 0 \end{aligned}$$

where $\mu \in \mathbb{R}_+^m$ and $\lambda \in \mathbb{R}^n$ are dual variables associated with primal inequalities and equalities, respectively.

5 Fenchel conjugate

Let $f : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be a function. The the Fenchel conjugate of f is the function $f^* : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ defined as:

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom } f} \mathbf{x}^\top \mathbf{y} - f(\mathbf{x})$$

6 Subgradient

Let $f : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be a function. A vector $\mathbf{g} \in \mathbb{R}^k$ is a subgradient of f at $\mathbf{x} \in \mathbb{R}^k$ if and only if:

$$\forall \mathbf{x}' \in \mathbb{R}^k : f(\mathbf{x}') \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{x}' - \mathbf{x})$$

7 Regularized prediction functions and Fenchel-Young losses

Let $\Omega : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex regularizer. The prediction function associated with Ω is:

$$\hat{\mathbf{y}}(\mathbf{w}) = \arg \max_{\boldsymbol{\mu} \in \text{dom } \Omega} \mathbf{w}^\top \boldsymbol{\mu} - \Omega(\boldsymbol{\mu})$$

Then Fenchel-Young loss associated with Ω is:

$$L_\Omega(\mathbf{w}; \mathbf{y}) = \Omega^*(\mathbf{w}) + \Omega(\mathbf{y}) - \mathbf{y}^\top \mathbf{w}$$

8 Inverse function

Let σ be the sigmoid function, i.e. $\sigma(w) = \frac{\exp(w)}{1+\exp(w)}$. The inverse of the sigmoid is: $\sigma^{-1}(\mu) = \log \frac{\mu}{1-\mu}$.