

Instructions: All your solutions should be prepared in \LaTeX and the PDF and .tex should be submitted to canvas. For each question, the best and correct answers will be selected as sample solutions for the entire class to enjoy. If you prefer that we do not use your solutions, please indicate this clearly on the first page of your assignment.

The programming parts can be written in the programming language of your choice and the code should be submitted alongside your solutions.

1. Entropy maximization. In this problem, we will consider the *entropy maximization problem*. Let us consider a probability distribution $\mathbf{x} \in \mathbb{R}^n$ over a finite set of size n . We have $\mathbf{x} \geq 0$ and $\sum_{i=1}^n x_i = 1$. The entropy of \mathbf{x} is defined by:

$$H(\mathbf{x}) = \sum_{i=1}^n x_i \log \frac{1}{x_i}$$

We are interested in maximizing entropy, or equivalently, solving the following problem:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \sum_{i=1}^n x_i \log x_i \\ \text{s.t.} \quad & \sum_{i=1}^n x_i = 1 \\ & \mathbf{x} \geq 0 \end{aligned}$$

- a. Prove Jensen's inequality: let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a strictly convex function, let $\mathbf{x}_1, \dots, \mathbf{x}_m$ be m vectors in \mathbb{R}^n , and let $\lambda_1, \dots, \lambda_m$ be such that $\sum_{i=1}^m \lambda_i = 1$ and $\lambda_i \geq 0$, $1 \leq i \leq m$, then:

$$f\left(\sum_{i=1}^m \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^m \lambda_i f(\mathbf{x}_i)$$

and prove that the inequality is an equality if and only if $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_m$.

- b. Using Jensen's inequality, what is the optimal solution to the entropy maximization problem above? Specify both the distribution \mathbf{x} of maximum entropy and the value of its entropy.
- c. We now add the constraint $A\mathbf{x} \leq b$ to the entropy maximization problem, where $A \in \mathbb{R}^{m \times n}$

and $\mathbf{b} \in \mathbb{R}^m$. The problem now becomes:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \sum_{i=1}^n x_i \log x_i \\ \text{s.t.} \quad & \sum_{i=1}^n x_i = 1 \\ & \mathbf{x} \geq 0 \\ & A\mathbf{x} \leq \mathbf{b} \end{aligned}$$

Show that the dual of this problem can be written in the following form:

$$\begin{aligned} \max_{\nu \in \mathbb{R}^m} \quad & -\mathbf{b}^\top \nu - \log \left(\sum_{i=1}^n e^{-\mathbf{a}_i^\top \nu} \right) \\ \text{s.t.} \quad & \nu \geq 0 \end{aligned}$$

where \mathbf{a}_i is the i th column of A . Assuming that strong duality holds for this problem, re-derive the result of part b. by considering a pair of primal/dual optimal solutions.

2. Minimum volume ellipsoid. An ellipsoid in \mathbb{R}^d is the image of the unit ball by a linear invertible map, i.e a set \mathcal{E} defined by:

$$\mathcal{E} = \{A\mathbf{x} : \mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|_2 \leq 1\}$$

for some invertible linear map $A : \mathbb{R}^d \mapsto \mathbb{R}^d$. In this case, we define the volume of the ellipsoid to be $|\det A|$. An equivalent parametrization of the ellipsoid is:

$$\mathcal{E} = \{\mathbf{y} \in \mathbb{R}^d : \mathbf{y}^\top W \mathbf{y} \leq 1\}$$

with $W = (A^{-1})^\top A^{-1}$. Note that W is symmetric positive definite and that under this parametrization, the volume of the ellipsoid is $(\det W)^{-1/2}$.

Let us denote by \mathbf{S}_d^{++} the set of symmetric positive definite matrices of size $d \times d$. Given n points $\mathbf{x}_1, \dots, \mathbf{x}_n$ in \mathbb{R}^d , the *minimum volume ellipsoid* problem consists in finding the ellipsoid of minimum volume containing all points $\mathbf{x}_1, \dots, \mathbf{x}_n$, that is:

$$\begin{aligned} \min_{W \in \mathbf{S}_d^{++}} \quad & (\det W)^{-1/2} \\ \text{s.t.} \quad & \mathbf{x}_i^\top W \mathbf{x}_i \leq 1, \quad 1 \leq i \leq n \end{aligned}$$

a. Show that \mathbf{S}_d^{++} is convex.

b. Let us define $d : \mathbf{S}_d^{++} \rightarrow \mathbb{R}$ by $d(W) = (\det W)^{-1/2}$. Is d convex over \mathbf{S}_d^{++} ?

Using the fact that \log is increasing over $\mathbb{R}^+ \setminus \{0\}$, we consider the following problem which is equivalent to the minimum volume ellipsoid problem:

$$\begin{aligned} \min_{W \in \mathbf{S}_d^{++}} \quad & \log \det(W^{-1}) \\ \text{s.t.} \quad & \mathbf{x}_i^\top W \mathbf{x}_i \leq 1, \quad 1 \leq i \leq n \end{aligned} \tag{1}$$

- c. Show that the function f defined by $f(W) = \log \det(W^{-1})$ is convex and differentiable over \mathbf{S}_d^{++} and that $\nabla f(W) = -W^{-1}$.
- d. Show that the dual of problem (1) is:

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^n} \log \det \left(\sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}_i^\top \right) - \sum_{i=1}^n \lambda_i + d \\ \text{s.t. } \lambda \geq 0, \sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}_i^\top \in \mathbf{S}_d^{++} \end{aligned}$$

- e. Show that the dual can be further simplified to:

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^n} \log \det \left(\sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}_i^\top \right) + d \log d \\ \text{s.t. } \lambda \geq 0, \sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}_i^\top \in \mathbf{S}_d^{++}, \sum_{i=1}^n \lambda_i = 1 \end{aligned} \tag{2}$$

3. Support vector machines. In this problem, we will reuse the dataset from Problem Set 3 on forged banknotes. The dataset is available at <http://thibaut.horel.org/banknotes.data>. This time, we will use support vector machines to construct the classifier. As we saw in class, after modifying the dataset (see the paragraph on Data Initialization in Lecture 2, Section 3.1), this amounts to finding $\mathbf{w} \in \mathbb{R}^d$ such that $\mathbf{w}^\top \mathbf{x}_i \geq 0$, for $1 \leq i \leq n$, where $\mathbf{x}_1, \dots, \mathbf{x}_n$ are the (modified) data points.

As seen in class, the optimization problem for support vector machines now takes the following form:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t. } \mathbf{w}^\top \mathbf{x}_i \geq 0, 1 \leq i \leq n \end{aligned}$$

In cases where the dataset is not linearly separable, it is not possible to find \mathbf{w} satisfying the constraints of the above problem. In particular, we might have $\mathbf{w}^\top \mathbf{x}_i < 0$ for some i . If this is the case, there exists $\xi_i \geq 0$ such that $\mathbf{w}^\top \mathbf{x}_i + \xi_i \geq 0$. The number ξ_i quantifies the “misclassification” of data point i . Since we want to discourage these misclassifications, we incorporate them into the objective function and consider the following optimizing problem instead:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \xi \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{i=1}^n \xi_i \\ \text{s.t. } \mathbf{w}^\top \mathbf{x}_i + \xi_i \geq 0, 1 \leq i \leq n \\ \xi_i \geq 0, 1 \leq i \leq n \end{aligned} \tag{3}$$

where λ is a parameter that we can choose depending on how much we want to penalize misclassified data points.

- a. Reuse your implementation of the perceptron algorithm and run it on the banknote dataset. Which behavior do you observe? Can you explain why?

- b. Use a convex solver to solve the convex program (3). Note that the objective function is quadratic, so you can use a function specific to quadratic problems. In CVXOPT, this is the `cvxopt.solvers.qp` function. Solve the problem for different values of λ and plot the classification accuracy (fraction of correctly classified data points, defined in Problem Set 3, problem 4) as a function of λ . How do you explain the shape of this plot? How does the best accuracy (for the best value of λ) compare to the accuracy obtained using boosting in Problem Set 3?
- c. Show that the program (3) is equivalent to the following problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{i=1}^n \max(0, -\mathbf{w}^\top \mathbf{x}_i) \quad (4)$$

- d. **[Optional, for bonus credits]** The advantage of problem (4) is that it is unconstrained. So we can use subgradient descent to solve it. Run the subgradient descent algorithm to solve Problem (4) for the best value of λ found in part b.
- e. **[Optional, for bonus credits]** Note that (4) also has a “separable” objective function as seen in Stochastic Gradient Descent (section 7). Implement the Stochastic Gradient descent algorithm and use it to solve (4). Compare the number of iterations required to reach the same accuracy with gradient descent (part d.) and stochastic gradient descent (part e.).