

EE 127 Final: Solutions

1. (8 points, Topic: Projections.) We consider a line in \mathbf{R}^n , described as

$$\mathcal{L} := \{x_0 + tu : t \in \mathbf{R}\},$$

where x_0, u are given n -vectors. Without loss of generality we assume that $\|u\|_2 = 1$.

- Find a closed-form expression for the projection $z(x)$ of an arbitrary point $x \in \mathbf{R}^n$ on the line \mathcal{L} .
- Show that the squared minimum distance D^2 from $x = 0$ to the line \mathcal{L} is $D^2 = x_0^T x_0 - (u^T x_0)^2$. Why is the latter expression non-negative?
- On what condition on a scalar $r \geq 0$ does a sphere \mathcal{S}_r , with center 0 and radius r , intersect the line? Make sure to express your condition clearly in mathematical form, involving x_0 and u .
- Assume that the sphere \mathcal{S}_r does intersect the line. Express the intersection as a segment, in the form

$$\mathcal{L} \cap \mathcal{S}_r = \{x_0 + tu : t \in [t_-, t_+]\},$$

where $t_- \leq t_+$ are scalars, which you will determine.

Solution:

- We solve the problem of minimizing the squared Euclidean distance:

$$\min_t \|x_0 + tu - x\|_2^2.$$

The objective function is a quadratic. Exploiting the fact that $\|u\|_2 = 1$, we obtain

$$\begin{aligned} \|x_0 + tu - x\|_2^2 &= t^2 \|u\|_2^2 - 2tu^T(x - x_0) + \|x - x_0\|_2^2 \\ &= (t - u^T(x - x_0))^2 + \|x - x_0\|_2^2 - (u^T(x - x_0))^2. \end{aligned}$$

The smallest value is attained with $t = t(x)$, where

$$t(x) = u^T(x - x_0).$$

The corresponding projection is

$$z(x) = x_0 + t(x)u = x_0 + u^T(x - x_0) \cdot u.$$

(b) We have

$$D^2 = \min_t \|x_0 + tu\|_2^2 = \|z(0)\|_2^2 = \|x_0 - (u^T x_0)u\|_2^2 = \|x_0\|_2^2 - (u^T x_0)^2,$$

as claimed. We have $\|x_0\|_2^2 \geq (u^T x_0)^2$, due to the Cauchy-Schwartz inequality:

$$|u^T x_0| \leq \|u\|_2 \cdot \|x_0\|_2 = \|x_0\|_2.$$

(c) The condition is $r \geq D$, that is:

$$r^2 \geq \|x_0\|_2^2 - (u^T x_0)^2.$$

(d) The intersection condition is that

$$\|x_0 + tu\|_2^2 \leq r^2.$$

Expanding the square, and using $\|u\|_2 = 1$, we obtain

$$(t + u^T x_0)^2 \leq r^2 - x_0^T x_0 + (u^T x_0)^2 = r^2 - D^2.$$

This is the same as $t_- \leq t \leq t_+$, with

$$t_{\pm} = -u^T x_0 \pm \sqrt{r^2 - D^2}.$$

2. (12 points, Topics: machine learning, SOCP.) We consider a binary classification problem where the number of data points in the negative class is far greater than that of the positive class. To avoid too large a computational burden due to the number of negatively-labelled points, we will only rely on a very simplified information: all of those points belong to a given (hyper-) sphere:

$$\mathcal{S}_- := \{\hat{x}_- + \rho u : \|u\|_2 \leq 1\}.$$

Here, $\hat{x}_- \in \mathbf{R}^n$ is the center of the hyper-sphere, and $\rho \geq 0$ a measure of its size. We assume that all the positively-labelled points, denoted $x_i^+ \in \mathbf{R}^n$, $i = 1, \dots, m_+$, are outside the sphere.

We consider linear separation, and will parametrize a candidate separating hyperplane as

$$\mathcal{H}(w, b) = \{x : w^T x + b = 0\},$$

where $w \in \mathbf{R}^n$, $b \in \mathbf{R}$ contain the parameters of the classifier. We will impose the following requirements on $\mathcal{H}(w, b)$:

- (i) All the negatively-labelled points should be in the half-space defined by $w^T x + b \leq 0$. In other words: there are no errors on the negative class.
- (ii) The number of positively-labelled points that are in the same half-space is minimal. In other words: there are as few errors as possible on the positive class.

In this exercise, we develop a convex programming approach to this problem.

- (a) Show that condition (i) is implied by the sufficient condition:

$$w^T \hat{x}_- + b + \rho \|w\|_2 \leq 0.$$

Hint: you may use the Cauchy-Schwartz inequality, more precisely, the fact that, for any $w \in \mathbf{R}^n$:

$$\|w\|_2 = \max_u u^T w : \|u\|_2 \leq 1.$$

- (b) Express the number of errors made on the positively-labelled class, in terms of (w, b) and the “error” function

$$E(z) = \begin{cases} 1 & \text{if } z \leq 0 \\ 0 & \text{otherwise.} \end{cases}$$

- (c) Explain why the following problem:

$$\min_{w, b} \sum_{i=1}^{m_+} \max(0, 1 - (w^T x_i^+ + b)) : w^T \hat{x}_- + b + \rho \|w\|_2 \leq 0,$$

is a sensible heuristic to address both requirements (i) and (ii). Make sure to explain how you arrive at the formulation above.

- (d) Formulate the above as an SOCP. Make sure to define the variables, the objective function, and the constraints, precisely.

Solution:

- (a) Condition (i) is equivalent to: for every $u \in \mathbf{R}^n$ with $\|u\|_2 \leq 1$, we have

$$w^T(\hat{x}^- + \rho u) + b \leq 0.$$

The above condition can be written as:

$$0 \geq w^T \hat{x}^- + b + \rho \cdot \max_{u: \|u\|_2 \leq 1} w^T u.$$

Using the hint we get the desired result.

- (b) The number of errors on the positive set is

$$\sum_{i=1}^{m_+} E(w^T x_i^+ + b),$$

where E is the error function.

- (c) We have

$$E(z) \leq \max(0, 1 - z).$$

This means that we can get an upper bound on the number of errors on the positive class, as

$$\sum_{i=1}^{m_+} E(w^T x_i^+ + b) \leq \sum_{i=1}^{m_+} \max(0, 1 - (w^T x_i^+ + b)).$$

The proposed formulation makes use of the above bound.

- (d) We can write the problem as an SOCP:

$$\min_{w, b, s} \sum_{i=1}^{m_+} s_i \quad : \quad \begin{aligned} s_i &\geq 0, \quad s_i \geq 1 - (w^T x_i^+ + b), \quad i = 1, \dots, m_+, \\ w^T \hat{x}^- + b + \rho \|w\|_2 &\leq 0. \end{aligned}$$

The variables are (w, b) and $s \in \mathbf{R}^{m_+}$.

3. (20 points, Topic: SVD.) We are given two data sets encoded in matrices $X = [x_1, \dots, x_m]$ and $Y = [y_1, \dots, y_m]$, where $x_i, y_i \in \mathbf{R}^n$, $i = 1, \dots, m$. We would like to analyze the correlations between these two multi-dimensional data sets. We will use the sample covariance matrices associated with X, Y :

$$\Sigma_{xx} := \frac{1}{m} \sum_{i=1}^m (x_i - \hat{x})(x_i - \hat{x})^T, \quad \Sigma_{yy} := \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y})(y_i - \hat{y})^T,$$

where $\hat{x} = (1/m)(x_1 + \dots + x_m) \in \mathbf{R}^n$ is the sample average of x , and \hat{y} is defined similarly. We will also use the so-called cross-covariance matrix, which is a $n \times n$ matrix defined as

$$\Sigma_{xy} := \frac{1}{m} \sum_{i=1}^m (x_i - \hat{x})(y_i - \hat{y})^T.$$

Note that Σ_{xy} is not symmetric in general. We assume that Σ_{xy} is non-zero, and that Σ_{xx}, Σ_{yy} are both positive-definite.

- (a) Define the *sample correlation* between two non-zero vectors $\alpha \in \mathbf{R}^m$ and $\beta \in \mathbf{R}^m$, as the number

$$\text{Corr}(\alpha, \beta) := \frac{1}{m\sigma_\alpha\sigma_\beta} \sum_{i=1}^m (\alpha_i - \hat{\alpha})(\beta_i - \hat{\beta}),$$

where $\hat{\alpha}, \hat{\beta}$ are the sample averages:

$$\hat{\alpha} = \frac{1}{m}(\alpha_1 + \dots + \alpha_m), \quad \hat{\beta} = \frac{1}{m}(\beta_1 + \dots + \beta_m),$$

and $\sigma_\alpha^2, \sigma_\beta^2$ are the sample variances:

$$\sigma_\alpha^2 := \frac{1}{m} \sum_{i=1}^m (\alpha_i - \hat{\alpha})^2, \quad \sigma_\beta^2 := \frac{1}{m} \sum_{i=1}^m (\beta_i - \hat{\beta})^2.$$

Explain precisely why we always have $|\text{Corr}(\alpha, \beta)| \leq 1$.

- (b) On what condition on α, β do we have $|\text{Corr}(\alpha, \beta)| = 1$?
(c) On to the multi-dimensional case: we seek linear combinations of the data points:

$$\alpha = u^T X, \quad \beta = v^T Y,$$

where the vector variables $u, v \in \mathbf{R}^n$, $u \neq 0$, $v \neq 0$, are to be chosen such that the sample correlation between the (row) vectors α, β is maximal. Express the problem as a maximization problem, of the form:

$$\max_{u \neq 0, v \neq 0} \frac{u^T \Sigma_{xy} v}{\sqrt{u^T \Sigma_{xx} u} \sqrt{v^T \Sigma_{yy} v}}.$$

(d) Show that the problem can be equivalently expressed as

$$\max_{u,v} u^T \Sigma_{xy} v : u^T \Sigma_{xx} u = 1, \quad v^T \Sigma_{yy} v = 1. \quad (1)$$

- (e) Explain how to reduce the above problem to the case when both matrices Σ_{xx}, Σ_{yy} are the identity. Make sure to describe precisely the relations between u, v and any transformed variables you are using. *Hint:* for a positive-definite matrix S , we can define the *matrix square-root* of S , denoted $S^{1/2}$, as the symmetric matrix with the same system of eigenvectors, and eigenvalues set at the square-root of the eigenvalues of S ; by construction, that matrix satisfies $S^{1/2} \cdot S^{1/2} = S$.
- (f) Explain how to solve problem (1) via the SVD of an appropriate matrix, which you will determine. Make sure to clearly state the optimal value of the problem, and optimal points u, v .
- (g) (5-point bonus question) Consider the problem with inequality constraints instead of equalities:

$$\phi := \max_{u,v} u^T \Sigma_{xy} v : u^T \Sigma_{xx} u \leq 1, \quad v^T \Sigma_{yy} v \leq 1.$$

Explain why the above problem has the same optimal value and optimal points as problem (1). *Hint:* argue first that $\phi > 0$, and then prove that both inequalities are equalities at optimum, using a scaling argument.

Solution:

- (a) Define the centered vectors $\bar{\alpha} = \alpha - \hat{\alpha} \mathbf{1}$, $\bar{\beta} = \beta - \hat{\beta} \mathbf{1}$, with $\mathbf{1}$ the vector of ones in \mathbf{R}^m . We have

$$\text{Corr}(\alpha, \beta) = \frac{\bar{\alpha}^T \bar{\beta}}{\|\bar{\alpha}\|_2 \|\bar{\beta}\|_2},$$

which implies the desired inequality $|\text{Corr}(\alpha, \beta)| \leq 1$, as a consequence of the Cauchy-Schwartz inequality.

- (b) Equality holds when $\bar{\alpha}, \bar{\beta}$ are collinear. This means that there is an affine relationship between α, β .
- (c) Define the centers of data points:

$$\hat{x} = \frac{1}{m}(x_1 + \dots + x_m), \quad \hat{y} = \frac{1}{m}(y_1 + \dots + y_m).$$

We have, with $\alpha_i = u^T x_i$, $\beta_i = v^T y_i$, $i = 1, \dots, m$:

$$\begin{aligned} \sigma_\alpha^2 &= \frac{1}{m} \sum_{i=1}^m (u^T (x_i - \hat{x}))^2 \\ &= \frac{1}{m} \sum_{i=1}^m u^T (x_i - \hat{x})(x_i - \hat{x})^T u \\ &= u^T \Sigma_{xx} u, \end{aligned}$$

and similarly,

$$\sigma_\beta^2 = v^T \Sigma_{yy} v.$$

Finally,

$$\begin{aligned} \sigma_\alpha \sigma_\beta \text{Corr}(\alpha, \beta) &= \frac{1}{m} \sum_{i=1}^m u^T (x_i - \hat{x})(x_i - \hat{x})^T v \\ &= u^T \Sigma_{xy} v. \end{aligned}$$

Our problem then can be written as:

$$\max_{u \neq 0, v \neq 0} \frac{u^T \Sigma_{xy} v}{\sqrt{u^T \Sigma_{xx} u} \sqrt{v^T \Sigma_{yy} v}}.$$

(d) This comes from the change of variables

$$u \rightarrow \frac{u}{\sqrt{u^T \Sigma_{xx} u}}, \quad v \rightarrow \frac{v}{\sqrt{v^T \Sigma_{yy} v}}.$$

(e) Let $\bar{u} := \Sigma_{xx}^{-1/2} u$, $\bar{v} := \Sigma_{yy}^{-1/2} v$. We can express our problem in an equivalent form:

$$\max_{\bar{u}, \bar{v}} \bar{u}^T \bar{\Sigma}_{xy} \bar{v} : \bar{u}^T \bar{u} = 1, \quad \bar{v}^T \bar{v} = 1,$$

where

$$\bar{\Sigma}_{xy} := \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2}.$$

(f) Thanks to the Cauchy-Schwartz inequality, the above problem reduces to

$$\max_{\bar{u}, \bar{v}} \|\bar{\Sigma}_{xy} \bar{v}\|_2 : \bar{v}^T \bar{v} = 1.$$

Using the Cauchy-Schwartz inequality again, we obtain that the optimal value of the problem is the largest singular value of $\bar{\Sigma}_{xy}$, and an optimal \bar{v} is any singular vector associated with the largest singular value of $\bar{\Sigma}_{xy}$.

An optimal \bar{u} is related to such an optimal \bar{v} by

$$\bar{u} = \frac{\bar{\Sigma}_{xy} \bar{v}}{\|\bar{\Sigma}_{xy} \bar{v}\|_2}.$$

Once \bar{u}, \bar{v} are found, we set

$$u = \Sigma_{xx}^{-1/2} \bar{u}, \quad v = \Sigma_{yy}^{-1/2} \bar{v}.$$

(g) Consider the problem with inequality constraints instead of equalities:

$$\phi := \max_{u,v} u^T \Sigma_{xy} v : u^T \Sigma_{xx} u \leq 1, v^T \Sigma_{yy} v \leq 1.$$

Since $(u, v) = (0, 0)$ is feasible, we have $\phi \geq 0$. Assume that $\phi = 0$, meaning that, for every u, v such that $u^T \Sigma_{xx} u \leq 1, v^T \Sigma_{yy} v \leq 1$, we have

$$u^T \Sigma_{xy} v \leq 0.$$

Since the constraints are unchanged upon replacing u by $-u$, we get: for every u, v such that $u^T \Sigma_{xx} u \leq 1, v^T \Sigma_{yy} v \leq 1$, we have

$$u^T \Sigma_{xy} v = 0.$$

This in turn implies that for *every* (u, v) we have $u^T \Sigma_{xy} v = 0$. Applying this to u, v equal to the i -th and j -th unit vectors in \mathbf{R}^n , respectively, we obtain that $\Sigma_{xy}(i, j) = 0$ for every i, j . This is a contradiction since we assumed $\Sigma_{xy} \neq 0$.

We conclude that $\phi > 0$. Now assume that say $\sigma^2 := u^T \Sigma_{xx} u < 1$ at optimum. We can replace u by u/σ , and *increase* the corresponding objective value, due to $\phi > 0$. This leads to a contradiction again. Thus, at optimum, we must have $u^T \Sigma_{xx} u = 1$. A similar result holds for the other constraint.