

# On the existence of solutions to adversarial training in multiclass classification

Nicolás García Trillo<sup>1</sup>, Matt Jacobs<sup>2</sup> and Jakwang Kim<sup>3</sup>

<sup>1</sup>Department of Statistics, University of Wisconsin-Madison, Madison, WI, USA

<sup>2</sup>Department of Mathematics, UC Santa Barbara, Santa Barbara, CA, USA

<sup>3</sup>Department of Mathematics, University of British Columbia, Vancouver, British Columbia, Canada

**Corresponding author:** Nicolás García Trillo; Email: [garciatrillo@wisc.edu](mailto:garciatrillo@wisc.edu)

**Received:** 22 February 2024; **Revised:** 03 November 2024; **Accepted:** 05 November 2024

**Keywords:** existence of solutions for minimax problems; nonparametric robustness; general topics in artificial intelligence; problem-solving in the context of artificial intelligence

**2020 Mathematics Subject Classification:** 49J35 (Primary); 62G35, 68T20 (Secondary)

## Abstract

Adversarial training is a min-max optimization problem that is designed to construct robust classifiers against adversarial perturbations of data. We study three models of adversarial training in the multiclass agnostic-classifier setting. We prove the existence of Borel measurable robust classifiers in each model and provide a unified perspective of the adversarial training problem, expanding the connections with optimal transport initiated by the authors in their previous work [21]. In addition, we develop new connections between adversarial training in the multiclass setting and total variation regularization. As a corollary of our results, we provide an alternative proof of the existence of Borel measurable solutions to the agnostic adversarial training problem in the binary classification setting.

## 1. Introduction

Modern machine learning models, in particular those generated with deep learning, perform remarkably well, in many cases much better than humans, at classifying data in a variety of challenging application fields like image recognition, medical image reconstruction, and natural language processing. However, the robustness of these learning models to data perturbations is a completely different story. For example, in image recognition, it has been widely documented (e.g., [26]) that certain structured but human-imperceptible modifications of images at the pixel level can fool an otherwise well-performing image classification model. These small data perturbations, known as *adversarial attacks*, when deployed at scale, can make a model’s prediction accuracy drop substantially and in many cases collapse altogether. As such, they are a significant obstacle to the deployment of machine learning systems in security-critical applications, e.g. [8]. To defend against these attacks, many researchers have investigated the problem of adversarial training, i.e., training methods that produce models that are robust to attacks. In adversarial training, one typically pits the adversary against the learner during the training step, forcing the learner to select a model that is robust against attacks. Nonetheless, despite the attention that has been devoted to understanding these problems, both theoretically and algorithmically, there are still several important mathematical questions surrounding them that have not been well understood.

A fundamental difficulty in adversarial training, in contrast to standard training of learning models, is the fact that the adversary has the power to alter the underlying data distribution. In particular, model training becomes an implicit optimisation problem over a space of measures. As a result, one may be forced to leave the prototypical setting of equivalence classes of functions defined over a single fixed

measure space. In general, measurability issues become more delicate for adversarial training problems at the moment of providing a rigorous mathematical formulation for the problem. Due to these difficulties, there are several subtle variations of the adversarial training model in the literature, and it has not been clear whether these models are fully equivalent. More worryingly, for some models, even the *existence* of optimal robust classifiers is unknown, essentially due to convexity and compactness issues.

Let us emphasise that these issues arise even in what can be regarded as the simplest possible setting of the agnostic learner, i.e., where the space of classifiers is taken to be the set of all possible Borel measurable weak (probabilistic) classifiers. While this setting is trivial in the absence of an adversary (there the optimal choice for the learner is always the Bayes classifier), the structure of the problem is much more subtle in the adversarial setting (in other words the analogue of the Bayes classifier is not fully understood). With an adversary, the training process can be viewed as a two-player min-max game (learner versus adversary) [4, 12, 32, 35], and as a result, the optimal strategies for the two players are far from obvious. By relaxing the problem to the agnostic setting, one at least is working over a convex space, but again measurability issues pose a problem for certain formulations of adversarial training.

In light of the above considerations, the purpose of this paper is twofold. On one hand, we provide rigorous justification for the existence of *Borel* measurable robust classifiers in the multiclass classification setting for three different models of adversarial training. Notably, our analysis includes a widely used model for which the existence of Borel classifiers was not previously known. On the other hand, we develop a series of connections between the three mathematical models of adversarial training discussed throughout the paper exploiting ideas from optimal transportation and total variation minimisation. By developing these connections, we hope to present a unified formulation of adversarial training and highlight the prospective advantages of using tools in computational optimal transport for solving these problems in the agnostic-classifier setting (and perhaps beyond the agnostic setting too). We also highlight, in concrete terms, the connection between adversarial training and the direct regularisation of learning models. To achieve all the aforementioned goals, we expand and take advantage of our previous work [21] as well as of the work [15] exploring the connection between adversarial training and perimeter minimisation in the binary classification setting.

### 1.1 Organisation of the paper

The rest of the paper is organised as follows. In Section 2, we introduce three different models for adversarial training in the multiclass classification setting that we will refer to as the open-ball model, the closed-ball model and the distributional-perturbing model. In Section 2.1, we state our main mathematical results, and in Section 2.2, we discuss related literature and some of the implications of our results. In Section 3, we lay down the main mathematical tools for analysing the distributional-perturbing model. Part of these tools come directly from our previous work [21], while others are newly developed. In Section 4, we prove our main results: first, we prove the existence of solutions for the distributional-perturbing model (Section 4.1); then, we prove that solutions to the distributional-perturbing model are solutions to the closed-ball model (Section 4.2); finally, we relate the closed-ball model to the open-ball model in Section 4.3. Lastly, in Section 5, we wrap up the paper and discuss future research directions.

## 2. Set-up and main results

The setting of our problem will be a feature space  $(\mathcal{X}, d)$  (a Polish space with metric  $d$ ) and a label space  $\mathcal{Y} := \{1, \dots, K\}$ , which will represent a set of  $K$  labels for a given classification problem of interest. For each  $x \in \mathcal{X}$ , we will use  $B_\varepsilon(x)$  ( $\overline{B}_\varepsilon(x)$ , respectively) to denote an open (closed) ball with radius  $\varepsilon$  centred at  $x$ . We denote by  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$  the set of input-to-output pairs and endow it with a Borel probability measure  $\mu \in \mathcal{P}(\mathcal{Z})$ , representing a ground-truth data distribution. For convenience, we will often describe the measure  $\mu$  in terms of its class probabilities  $\mu = (\mu_1, \dots, \mu_K)$ , where each  $\mu_i$  is the positive Borel measure (not necessarily a probability measure) over  $\mathcal{X}$  defined according to:

$$\mu_i(A) = \mu(A \times \{i\}),$$

for  $A \in \mathfrak{B}(\mathcal{X})$ , i.e.,  $A$  is a Borel measurable subset of  $\mathcal{X}$ . Notice that the measures  $\mu_i$  are, up to normalisation factors, the conditional distributions of inputs/features given the different output labels.

Typically, a (multiclass) classification rule in the above setting is simply a Borel measurable map  $f: \mathcal{X} \rightarrow \mathcal{Y}$ . In this paper, however, it will be convenient to expand this notion slightly and interpret general classification rules as Borel measurable maps from  $\mathcal{X}$  into  $\Delta_{\mathcal{Y}} := \{(u_i)_{i \in \mathcal{Y}} : 0 \leq u_i \leq 1, \sum_{i \in \mathcal{Y}} u_i \leq 1\}$ , the set of (up to normalisation constants) probability distributions over  $\mathcal{Y}$  (see Remark 2.2); oftentimes these functions are known as *soft-classifiers*. For future reference, we denote by  $\mathcal{F}$  the set

$$\mathcal{F} := \{f: \mathcal{X} \rightarrow \Delta_{\mathcal{Y}} : f \text{ is Borel measurable}\}. \quad (1)$$

Given  $f \in \mathcal{F}$  and  $x \in \mathcal{X}$ , the vector  $f(x) = (f_1(x), \dots, f_K(x))$  will be interpreted as the vector of probabilities over the label set  $\mathcal{Y}$  that the classifier  $f$  assigns to the input data point  $x$ . In practice, from one such  $f$ , one can induce actual (hard) class assignments to the different inputs  $x$  by selecting the coordinate in  $f(x)$  with largest entry. The extended notion of classifier considered in this paper is actually routinely used in practice as it fares well with the use of standard optimisation techniques (in particular,  $\mathcal{F}$  is natural as it can be viewed as a convex relaxation of the space of maps from  $\mathcal{X}$  to  $\mathcal{Y}$ ).

The goal in the standard (unrobust) classification problem is to find a classifier  $f \in \mathcal{F}$  that gives accurate class assignments to inputs under the assumption that data points are distributed according to the ground-truth distribution  $\mu$ . This aim can be mathematically modelled as an optimisation problem of the form:

$$\inf_{f \in \mathcal{F}} R(f, \mu), \quad (2)$$

where  $R(f, \mu)$  is the risk of a classifier  $f$  relative to the data distribution  $\mu$ :

$$R(f, \mu) := \mathbb{E}_{(X, Y) \sim \mu} [\ell(f(X), Y)].$$

Here, a loss function is defined as  $\ell: \Delta_{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}$ . For general and reasonable  $\ell$ , one can observe that solutions to the risk minimisation Problem (2) are the standard multiclass Bayes classifiers from statistical learning theory (e.g., see [13, 40]). These classifiers are characterised by the condition  $f_{Bayes, i}^*(x) = \operatorname{argmax}_{i \in \mathcal{Y}} \mathbb{P}(Y = i | X = x)$  with an arbitrary tie-breaking rule, and such Bayes classifier is written as  $f_{Bayes}^*(x) = (\mathbb{1}_{A_1^*}(x), \dots, \mathbb{1}_{A_K^*}(x))$ , where  $A_1^*, \dots, A_K^*$  form a measurable partition of  $\mathcal{X}$ . In other words, there always exist hard classifiers that solve the risk minimisation Problem (2).

By definition, a solution to (2) classifies *clean* data optimally; by clean data here we mean data distributed according to the original distribution  $\mu$ . However, one should not expect the standard Bayes classifier to perform equally well when inputs have been adversarially contaminated, and the goal in adversarial training is precisely to create classifiers that are less susceptible to data corruption. One possible way to enforce this type of robustness is to replace the objective function in (2) with one that incorporates the actions of a well defined adversary and then search for the classifier that minimises the new notion of (adversarial) risk. This adversarial risk can be defined in multiple ways, but two general ways stand out in the literature and will be the emphasis of our discussion; we will refer to these two alternatives as the *data-perturbing adversarial model* and the *distribution-perturbing adversarial model*. More precisely, the data-perturbing adversarial model is generally defined as

$$\inf_{f \in \mathcal{F}} \mathbb{E}_{(X, Y) \sim \mu} \left[ \sup_{\tilde{X} \in A_{\varepsilon}(X)} \ell(f(\tilde{X}), Y) \right]$$

where  $A_{\varepsilon}(\cdot)$  is a measurable set parameterised by and *adversarial budget*  $\varepsilon$  (usually, closed/open ball with radius  $\varepsilon$ ), and the distribution-perturbing adversarial model is defined as:

$$\inf_{f \in \mathcal{F}} \sup_{\tilde{\mu}} \mathbb{E}_{(\tilde{X}, \tilde{Y}) \sim \tilde{\mu}} [\ell(f(\tilde{X}), \tilde{Y}) - C_{\varepsilon}(\mu, \tilde{\mu})]$$

where  $C_\varepsilon(\mu, \tilde{\mu})$  is some cost or distance over  $\mathcal{P}(\mathcal{X})$ . Although the interpretation of these models is straightforward, their theoretical underpinnings are not completely understood. This is true even for the simple linear loss function

$$\ell(u, i) = 1 - u_i, \quad (u, i) \in \Delta_{\mathcal{Y}} \times \mathcal{Y},$$

which, in lieu of the fact that  $\ell(e_j, i)$  is equal to 1 if  $i \neq j$  and 0 if  $i = j$  ( $e_j$  is the extremal point of  $\Delta_{\mathcal{Y}}$  with entry one in its  $j$ -th coordinate), will be referred to as the 0-1 loss. We will focus on this important case from now on. With this choice,  $R(f, \mu)$  can be written explicitly as:

$$R(f, \mu) = \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} (1 - f_i(x)) d\mu_i(x).$$

As it turns out, some results in the literature have established some connections between the two types of models (see [35] for more details), and we will develop further connections shortly. For the *data-perturbing adversarial model*, we will consider the following two versions:

$$R_{open}^*(\varepsilon) := \inf_{f \in \mathcal{F}} R_{open}^*(f; \varepsilon) := \inf_{f \in \mathcal{F}} \left\{ \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \sup_{\tilde{x} \in B_\varepsilon(x)} \{1 - f_i(\tilde{x})\} d\mu_i(x) \right\}, \quad (3)$$

and

$$\inf_{f \in \mathcal{F}} \left\{ \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \sup_{\tilde{x} \in \bar{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\} d\mu_i(x) \right\}. \quad (4)$$

Recall  $B_\varepsilon(x)$  ( $\bar{B}_\varepsilon(x)$ , respectively) denotes an open (a closed) ball. In both versions, the adversary can substitute any given input  $x$  with a  $\tilde{x}$  that belongs to a small ball of radius  $\varepsilon$  around the original  $x$ . The adversary has more power to perturb an input  $x$  by  $\tilde{x}$  with a larger  $\varepsilon$  since more powerful deviation will be possible with larger  $\varepsilon$ . Due to measurability issues that we will discuss next, at this stage, Problem (4) is introduced informally.

In both Problems (3) and (4), the learner's goal is to minimise the worst-loss that the adversary may induce by carrying out one of their feasible actions. Although at the heuristic level the difference between the two models is subtle (in (3) the adversary optimises over open balls and in (4) over closed balls), at the mathematical level these two models can be quite different. For starters, the Problem (4) is not well formulated, as it follows from a classical result in [30], which discusses that, in general, the function  $x \mapsto \sup_{\tilde{x} \in \bar{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\}$  may not be Borel measurable when only the Borel measurability of the function  $f_i$  has been assumed. For this reason, the integral with respect to  $\mu_i$  in (4) (which is a Borel positive measure, i.e., it is only defined over the Borel  $\sigma$ -algebra) may not be defined for all  $f \in \mathcal{F}$ . In Subsection 2.1, we will provide a rigorous formulation of (4), which we will call the closed-ball Model (9) because the adversarial attack lies in the closed ball. This reformulation will require the use of an extension of the Borel  $\sigma$ -algebra, known as the universal  $\sigma$ -algebra, as well as an extension of the measures  $\mu_i$  to this enlarged  $\sigma$ -algebra. Problem (3), on the other hand, is already well formulated, as no measurability issues arise when taking the sup over open balls. At a high level, this is a consequence of the fact that arbitrary unions of open balls are open sets and thus Borel measurable; see, for example, [15, Remark 2.3]. Regardless of which of the two models one adopts, and putting aside for a moment the measurability issues mentioned above, it is unclear whether it is possible to find minimisers for any of the Problems (3) and (4) within the family  $\mathcal{F}$ .

Moving on to the *distributional-perturbing adversarial model*, notice that it is defined as a minimax problem that can be described as follows: after the adversary selects a new data distribution  $\tilde{\mu} \in \mathcal{P}(\mathcal{X})$  and, by paying some cost  $C(\mu, \tilde{\mu})$ , attempts to make the risk  $R(\cdot, \tilde{\mu})$  be as large as possible, the learner has chosen a classifier  $f \in \mathcal{F}$ . The reverse interpretation is also true due to the strong duality of (5) in Ref. [21]. Theorem 2.5, furthermore, implicitly implies the strong duality as a corollary since it proves a stronger result, the existence of saddle points (Nash equilibria).

Precisely, the distributional-perturbing adversarial model can be rewritten as:

$$R_{DRO}^*(\varepsilon) := \inf_{f \in \mathcal{F}} \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \{R(f, \tilde{\mu}) - C(\mu, \tilde{\mu})\}, \quad (5)$$

where we assume the cost  $C: \mathcal{P}(\mathcal{Z}) \times \mathcal{P}(\mathcal{Z}) \rightarrow [0, \infty]$  takes the form:

$$C(\mu, \tilde{\mu}) := \inf_{\pi \in \Gamma(\mu, \tilde{\mu})} \int c_{\mathcal{Z}}(z, \tilde{z}) d\pi(z, \tilde{z}),$$

for some Borel measurable cost function  $c_{\mathcal{Z}}: \mathcal{Z} \times \mathcal{Z} \rightarrow [0, \infty]$ . Here and in the remainder of the paper, we use  $\Gamma(\cdot, \cdot)$  to represent the set of couplings between two positive measures over the same space; for example,  $\Gamma(\mu, \tilde{\mu})$  denotes the set of positive measures over  $\mathcal{Z} \times \mathcal{Z}$  whose first and second marginals are  $\mu$  and  $\tilde{\mu}$ , respectively. Note that Problem (5) is an instance of the *distributionally robust optimisation* (DRO) problem. Problem (5) is well defined given that all its terms are written as integrals of Borel measurable integrands against Borel measures.

In the remainder, we will assume that the cost  $c_{\mathcal{Z}}: \mathcal{Z} \times \mathcal{Z} \rightarrow [0, \infty]$  has the form

$$c_{\mathcal{Z}}(z, \tilde{z}) := \begin{cases} c(x, \tilde{x}) & \text{if } y = \tilde{y}, \\ \infty & \text{otherwise} \end{cases} \quad (6)$$

for a lower semi-continuous function  $C: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$ . Note that when  $c_{\mathcal{Z}}$  has the above structure, we can rewrite  $C(\mu, \tilde{\mu})$  as:

$$C(\mu, \tilde{\mu}) = \sum_{i=1}^K C(\mu_i, \tilde{\mu}_i),$$

where on the right-hand side, we slightly abuse notation and use  $C(\mu_i, \tilde{\mu}_i)$  to represent

$$C(\mu_i, \tilde{\mu}_i) := \inf_{\pi \in \Gamma(\mu_i, \tilde{\mu}_i)} \int c(x, \tilde{x}) d\pi(x, \tilde{x}).$$

Although  $R_{DRO}^*(\varepsilon)$  seems not to depend on  $\varepsilon$  explicitly unlike Problems (3) and (4), readers can understand that a cost function  $c$  relies implicitly on the choice of adversarial budget  $\varepsilon$ . The most popular example, in both theory and practice, of a cost function  $c$ , that we discuss in detail throughout this paper is the  $0\text{-}\infty$  cost function:

$$c(x, \tilde{x}) = c_{\varepsilon}(x, \tilde{x}) := \begin{cases} \infty & \text{if } d(x, \tilde{x}) > \varepsilon \\ 0 & \text{if } d(x, \tilde{x}) \leq \varepsilon. \end{cases} \quad (7)$$

Similar to (3) and (4), a larger  $\varepsilon$  means more powerful adversarial attacks, as the adversary may select arbitrary  $\tilde{\mu}$  at no cost.

**Remark 2.1.** Throughout the paper, we use the convention that  $C(\mu_i, \tilde{\mu}_i) = \infty$  whenever the set of couplings  $\Gamma(\mu_i, \tilde{\mu}_i)$  is empty. This is the case when  $\mu_i$  and  $\tilde{\mu}_i$  have different total masses.

**Remark 2.2.** Given the structure of the  $0\text{-}1$  loss function considered here, in all the adversarial models introduced above, we may replace the set  $\mathcal{F}$  with the set of those  $f \in \mathcal{F}$  for which  $\sum_i f_i = 1$ . Indeed, given  $f \in \mathcal{F}$ , we can always consider  $\tilde{f} \in \mathcal{F}$  defined according to  $\tilde{f}_{i_0} := f_{i_0} + (1 - \sum_{i \in \mathcal{Y}} f_i)$  and  $\tilde{f}_i = f_i$  for  $i \neq i_0$  to obtain a value of risk that is no greater than the one of the original  $f$ .

## 2.1 Main results

Our first main theorem discusses the existence of (Borel) solutions for Problem (5) under the assumptions on the cost  $c: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$  stated below.

**Assumption 2.3.** We assume that the cost  $c: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$  is a lower semi-continuous and symmetric function satisfying  $c(x, x) = 0$  for all  $x \in \mathcal{X}$ . We also assume that the following compactness

property holds: if  $\{x_n\}_{n \in \mathbb{N}}$  is a bounded sequence in  $(\mathcal{X}, d)$  and  $\{x'_n\}_{n \in \mathbb{N}}$  is a sequence satisfying  $\sup_{n \in \mathbb{N}} c(x_n, x'_n) < \infty$ , then  $\{(x_n, x'_n)\}_{n \in \mathbb{N}}$  is precompact in  $\mathcal{X} \times \mathcal{X}$  (endowed with the product topology).

**Remark 2.4.** Notice that Assumption 2.3 implicitly requires bounded subsets of  $\mathcal{X}$  to be precompact. In particular, when the cost has the form as in (7), Assumption 2.3 implies that every bounded subset in  $\mathcal{X}$  should be precompact. We want to emphasise that this assumption is not too strong and is satisfied in many natural settings, e.g., Euclidean spaces or smooth manifolds of finite dimension endowed with its geodesic distance.

A simple cost function  $c$ , which does not satisfy Assumption 2.4, is

$$c'(x, \tilde{x}) = \begin{cases} \infty & \text{if } d(x, \tilde{x}) \geq \varepsilon \\ 0 & \text{if } d(x, \tilde{x}) < \varepsilon. \end{cases}$$

This is almost the same as (7), but  $c'$  is no longer lower semi-continuous. One can easily see that replacing  $c_\varepsilon$  by  $c'$  changes the closed-ball model (4) to the open-ball model (3). It is reasonable to expect that  $c'$  plays almost the same as  $c_\varepsilon$ : we will see that indeed this is true in Theorem 2.8.

**Theorem 2.5.** Suppose that  $c: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$  satisfies Assumption 2.3. Then there exists a (Borel) solution  $f^*$  of the DRO model (5). Furthermore, there exists  $\tilde{\mu}^* \in \mathcal{P}(\mathcal{X})$  such that  $(f^*, \tilde{\mu}^*)$  is a saddle point for (5). In other words, the following holds: for any  $\tilde{\mu} \in \mathcal{P}(\mathcal{X})$  and any  $f \in \mathcal{F}$ , we have

$$R(f^*, \tilde{\mu}) - C(\mu, \tilde{\mu}) \leq R(f^*, \tilde{\mu}^*) - C(\mu, \tilde{\mu}^*) \leq R(f, \tilde{\mu}^*) - C(\mu, \tilde{\mu}^*). \quad (8)$$

When the cost function  $c$  is regular enough or when  $\mu$  is an empirical measure, we can reduce the problem of finding a solution  $f^*$  of (5) to the problem of solving the dual of a *generalised barycenter problem* or the dual of a *multimarginal optimal transport* problem. These connections were first put forward in our earlier work [21] and will be discussed again in Section 3, concretely in Proposition 3.10. Unfortunately, when the cost is only lower semi-continuous (e.g., for  $c$  as in (7)) and when  $\mu$  is an arbitrary Borel probability measure, we cannot directly use the content of Proposition 3.10 to guarantee the existence of (Borel) solutions  $f^*$ . To overcome this issue, we approximate  $c$  with a sequence of continuous costs  $c_n$  such that the previous theory applies. We then show that both the optimality and the Borel measurability of the optimal classifier are preserved in the limit. At a high level, we can thus reduce finding solutions for the DRO Problem (5) to that of an multimarginal optimal transport (MOT) or a generalised barycenter (or sequences thereof).

Next, we revisit Problem (4) and reformulate it in a rigorous way. Let us first introduce the *universal*  $\sigma$ -algebra of the space  $\mathcal{X}$ .

**Definition 1.** [33, Definition 2.2]. Let  $\mathcal{B}(\mathcal{X})$  be the Borel  $\sigma$ -algebra over  $\mathcal{X}$ , and let  $\mathcal{M}(\mathcal{X})$  be the set of all signed  $\sigma$ -finite Borel measures over  $\mathcal{X}$ . For each  $\nu \in \mathcal{M}(\mathcal{X})$ , let  $\mathcal{L}_\nu(\mathcal{X})$  be the completion of  $\mathcal{B}(\mathcal{X})$  with respect to  $\nu$ . The universal  $\sigma$ -algebra of  $\mathcal{X}$  is defined as:

$$\mathcal{U}(\mathcal{X}) := \bigcap_{\nu \in \mathcal{M}(\mathcal{X})} \mathcal{L}_\nu(\mathcal{X}).$$

We will use  $\overline{\mathcal{P}}(\mathcal{X})$  to denote the set of probability measures  $\gamma$  over  $\mathcal{X}$  for which  $\gamma_i$  is a universal positive measure (i.e., it is defined over  $\mathcal{U}(\mathcal{X})$ ) for all  $i \in \mathcal{Y}$ . For a given probability measure  $\mu \in \mathcal{P}(\mathcal{X})$ , we will denote by  $\overline{\mu}$  its universal extension, which will be interpreted as:

$$\overline{\mu}(A \times \{i\}) := \overline{\mu}_i(A), \quad \forall A \in \mathcal{U}(\mathcal{X}),$$

where  $\overline{\mu}_i$  is the extension of  $\mu_i$  to  $\mathcal{U}(\mathcal{X})$ .

Having introduced the above notions, we can reformulate Problem (4) as:

$$R_{closed}^*(\varepsilon) := \inf_{f \in \mathcal{F}} R_{closed}^*(f; \varepsilon) := \inf_{f \in \mathcal{F}} \left\{ \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \sup_{\tilde{x} \in \tilde{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\} d\overline{\mu}_i(x) \right\}. \quad (9)$$

Although the difference with (4) is subtle (in [4] we use  $\mu_i$ 's whereas in [9] we use  $\overline{\mu_i}$ 's), Problem (9) is rigorously well defined. Indeed, combining [35, Lemma 4.2] with [6, Corollary 7.42.1], which was originally proved in Ref. [30], it follows that for any Borel measurable  $f_i$ , the function  $x \mapsto \sup_{\tilde{x} \in \overline{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\}$  is universally measurable and thus the integrals on the right-hand side of (9) are well defined.

**Remark 2.6.** Similar to the universal extension of Borel measure, one might want to extend the solution space  $\mathcal{F}$  to  $\overline{\mathcal{F}}$ , which is the set of all  $f = (f_1, \dots, f_K)$  for which each  $0 \leq f_i \leq 1$  is universally measurable and  $\sum f_i \leq 1$ . However, unless  $\mathcal{X} = \mathbb{R}^d$ , extending the solution space requires hierarchical extension of measures, i.e. universal extension of universal extension: see [33]. Note that on  $\mathbb{R}^d$ , every universally measurable set is Lebesgue measurable.

Our second main result relates solutions of (5) with solutions of (9).

**Theorem 2.7.** *There exists a Borel solution of (5) for the cost function  $c$  defined in (7) that is also a solution of (9). In particular, there exists a (Borel) solution for (9).*

Finally, we connect Problem (9) with Problem (3).

**Theorem 2.8.** *For all but at most countably many  $\varepsilon \geq 0$ , we have  $R_{\text{open}}^*(\varepsilon) = R_{\text{closed}}^*(\varepsilon)$ . Moreover, for those  $\varepsilon \geq 0$  for which this equality holds, every solution  $f^*$  of (9) is also a solution of (3).*

**Remark 2.9.** Theorem 2.8 is optimal in general. One cannot expect the optimal adversarial risks of (3) and (9) to agree for all values of  $\varepsilon$ . To illustrate this, consider the simple setting of a two-class problem (i.e.,  $K = 2$ ) where  $\mu_1 = \frac{1}{2}\delta_{x_1}$  and  $\mu_2 = \frac{1}{2}\delta_{x_2}$ , where  $\delta_x$  denotes the usual Dirac delta measure at  $x$ . Let  $\epsilon_0 = \frac{1}{2}d(x_1, x_2)$ . It is straightforward to check that  $R_{\text{open}}^*(\epsilon_0) = 0$  whereas  $R_{\text{closed}}^*(\epsilon_0) = 1/2$ . Naturally, if we had selected any other value for  $\varepsilon > 0$  different from  $\epsilon_0$ , we would have obtained  $R_{\text{open}}^*(\varepsilon) = R_{\text{closed}}^*(\varepsilon)$ .

From Theorems 2.5, 2.7 and 2.8, we may conclude that it is essentially sufficient to solve Problem (5) to find a solution for all other formulations of the adversarial training problem discussed in this paper. Our results thus unify all notions of adversarial robustness into the single DRO problem, (5). The advantage of (5) over the other formulations of the adversarial training problem is that it can be closely related to a generalised barycenter problem or an MOT problem, as has been discussed in detail in our previous work [21] (see also section 3 below). In turn, either of those problems can be solved using computational optimal transport tools. Recently, in [36], the authors propose fast algorithms for solving (5). From a practical perspective, hence, it is thus easier to work with the DRO formulation than with the other formulations of adversarial training.

## 2.2 Discussion and literature review

The existence of measurable ‘‘robust’’ solutions to optimisation problems has been a topic of interest not only in the context of adversarial training of classification problem [2, 3, 19, 20, 35] but also in the general DRO literature, e.g., [9–11, 28]. Most previous studies of robust classifiers use the *universal  $\sigma$ -algebra* not only to formulate optimisation problems rigorously, but also as a feasible search space for robust classifiers. The proofs of these existence results rely on the pointwise topology of a sequence of universally measurable sets, the weak topology on the space of probability measures, and lower semi-continuity properties of  $R_{\text{closed}}^*(\cdot; \varepsilon)$ . The (universal) measurability of a minimiser is then guaranteed immediately by the definition of the universal  $\sigma$ -algebra. We want to emphasise that all the works [2, 3, 19, 20, 35] prove their results in the binary ( $K = 2$ ) classification setting where  $\mathcal{X}$  is a subset of Euclidean space.

In contrast to the closed-ball model formulation, the objective in (5) is well defined for all Borel probability measures  $\tilde{\mu}$  and all  $f \in \mathcal{F}$ , as has been discussed in previous sections. The papers [2, 3, 35] can only relate, in the binary setting, problems (5) and (9) when problem (5) is appropriately extended to the universal  $\sigma$ -algebra, yet it is not clear that such extension is necessary. Recently, the authors of [19, 20] remove this technical redundancy and prove the existence of Borel robust classifiers (and consistency)

for the binary classification over Euclidean space. However, it is not obvious that the proof techniques in those papers can be extended to the multiclass setting. The proof technique that we implement in this paper allows us to consider all cases at once, while at the same time allows us to reveal the connections between many of the different models for adversarial training that exist in the literature.

For concreteness, we summarise some related results in the literature in the following theorem. Let  $\bar{\mathcal{F}}$  denote the set of all  $f = (f_1, \dots, f_K)$  for which each  $0 \leq f_i \leq 1$  is universally measurable and  $\sum f_i \leq 1$ .

**Theorem 2.10.** [2, 3, 19, 20, 35]. *Suppose  $K = 2$  and  $\bar{\mu} \in \bar{\mathcal{P}}(\mathcal{X})$ . Then, for any  $f \in \bar{\mathcal{F}}$ , we have  $\sup_{\tilde{x} \in \bar{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\} \in \mathcal{U}(\mathcal{X})$  and*

$$\sum_{i=1}^2 \int_{\mathcal{X}} \sup_{\tilde{x} \in \bar{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\} d\bar{\mu}_i(x) = \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{X})} \{R(f, \tilde{\mu}) - C(\bar{\mu}, \tilde{\mu})\},$$

where  $C$  is defined in terms of the cost function  $c$  from (7).

Assume further that  $(\mathcal{X}, d) = (\mathbb{R}^d, \|\cdot\|)$ . Then, it holds that  $\sup_{\tilde{x} \in \bar{B}_\varepsilon(\cdot)} \{1 - f_i(\tilde{x})\}$  is universally measurable for any  $f \in \bar{\mathcal{F}}$ . Also, there exists a universally measurable minimiser ([2, 3, 35]) and a Borel measurable minimiser ([19, 20]). Finally, the values of (9) and (5) for the binary setting coincide.

In this paper, we use the universal  $\sigma$ -algebra only for the rigorous description of the objective function in (9). Instead, we will focus on Borel measurable classifiers in the adversarial training in the multiclass setting. First, based on some of our previous results in [21], we will prove the existence of Borel measurable robust classifiers of (5) for general lower semi-continuous  $c$  satisfying Assumption 6 only in the multiclass setting. This extension is not immediate from all previous results since most techniques there are tailored for the binary setting. Our proof, however, is much simpler even in this general setting. Employing the combination of duality, an approximation argument, and well known theorems of optimal transport literature, we will be able to achieve our goal. Then, back to the closed-ball model, we prove the existence of Borel robust classifiers of (9) by proving that solutions to the DRO model with cost  $c_\varepsilon$  are also solutions to the closed-ball model problem. Furthermore, based on a simple but important observation connecting (9) and (3), Theorems 2.7 and 2.8 imply that (5), (9) and (3) are indeed equivalent (at least for almost all values of adversarial budget) in the sense that if  $f^*$  is a Borel minimiser of one of them, then it is a minimiser of others automatically, which is a new and stronger result in the adversarial training literature.

When we specialise our results to the binary classification setting (i.e.,  $K = 2$ ), we obtain the following improvement upon the results from [7, 20, 34].

**Corollary 2.11.** *Let  $K = 2$  and let  $f^* \in \mathcal{F}$  be any solution to the problem (9). Then, for Lebesgue a.e.  $t \in [0, 1]$ , the pair  $(\mathbb{1}_{\{f_1^* \geq t\}}, \mathbb{1}_{\{f_1^* \geq t\}})$  is also a solution to (9).*

*In particular, there exist solutions to the problem*

$$\min_{A \in \mathcal{B}(\mathcal{X})} \int_{\mathcal{X}} \sup_{\tilde{x} \in \bar{B}_\varepsilon(x)} \mathbb{1}_{A^c}(\tilde{x}) d\bar{\mu}_1(x) + \int_{\mathcal{X}} \sup_{\tilde{x} \in \bar{B}_\varepsilon(x)} \mathbb{1}_A(\tilde{x}) d\bar{\mu}_2(x).$$

Notice that Corollary 2.11 implies, for the binary setting, the existence of robust hard classifiers for the adversarial training problem, a property shared with the standard risk minimisation Problem (2) that we discussed at the beginning of Section 2. Analogous results on the equivalence of the hard-classification and soft-classification problems in adversarial training under the binary setting have been obtained in [15, 24, 34, 35]. Unfortunately, when the number of classes is such that  $K > 2$ , the hard-classification and soft-classification problems in adversarial training may not be equivalent, as has been discussed in [21, Section 5.2].

In light of Theorem 2.8, one can conclude from Corollary 2.11 that for all but countably many  $\varepsilon > 0$  the problem

$$\min_{A \in \mathcal{B}(\mathcal{X})} \int_{\mathcal{X}} \sup_{\tilde{x} \in B_\varepsilon(x)} \mathbb{1}_{A^c}(\tilde{x}) d\mu_1(x) + \int_{\mathcal{X}} \sup_{\tilde{x} \in B_\varepsilon(x)} \mathbb{1}_A(\tilde{x}) d\mu_2(x)$$

admits solutions; notice that the above is the open-ball version of the optimisation problem in Corollary 2.11. However, notice that the results in [15] guarantee existence of solutions for *all* values of  $\varepsilon$ . It is interesting to note that the technique used in [15] cannot be easily adapted to the multiclass case  $K > 2$ . Specifically, it does not seem to be straightforward to generalise [15, Lemma C.1] to the multiclass case since it implicitly relies on the fact that one can always find an optimal hard classifier, i.e. an optimal decision boundary in the binary setting. Since for  $f = (f_1, f_2)$ ,  $f_2 = 1 - f_1$ ,  $f$  is basically a real-valued function  $f_1$ . Tailoring  $f_1$  suitably, one can replace sup by ess sup (inf by ess inf similarly), and use  $L^p$  spaces, which provide an appropriate topology for the purpose. For example, however, if one used the aforementioned lemma to modify the coordinate functions  $f_i$  of a multiclass classifier  $f$ , one could end up producing functions for which their sum may be greater than one for some points in  $\mathcal{X}$ , thus violating one of the conditions for belonging to  $\mathcal{F}$ .

We observe, on the other hand, that the total variation regularisation interpretation for the open ball model in the binary setting discussed in Ref. [15] continues to hold in the multiclass case. To make this connection precise, let us introduce the non-local TV functionals:

$$\widetilde{\text{TV}}_\varepsilon(f_i, \mu_i) := \frac{1}{\varepsilon} \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \sup_{\tilde{x} \in B_\varepsilon(x)} \{f_i(x) - f_i(\tilde{x})\} d\mu_i(x).$$

It is then straightforward to show that Problem (3) is equivalent to

$$\inf_{f \in \mathcal{F}} \sum_{i=1}^K \int_{\mathcal{X}} (1 - f_i(x)) d\mu_i(x) + \varepsilon \sum_{i=1}^K \widetilde{\text{TV}}_\varepsilon(f_i, \mu_i), \quad (10)$$

which can be interpreted as a total variation minimisation problem with fidelity term. Indeed, the fidelity term in the above problems is the standard (unrobust) risk  $R(f, \mu)$ . On the other hand, the functional  $\widetilde{\text{TV}}_\varepsilon(\cdot, \mu_i)$  is a non-local total variation functional in the sense that it is convex, positive 1-homogeneous, invariant under addition of constants to the input function and is equal to zero when its input is a constant function. Moreover, in the case  $(\mathcal{X}, d) = (\mathbb{R}^d, \|\cdot\|)$  and when  $d\mu_i(x) = \rho_i(x)dx$  for a smooth function  $\rho_i$ , one can see that, for small  $\varepsilon > 0$ ,

$$\widetilde{\text{TV}}_\varepsilon(f_i, \mu_i) \approx \int_{\mathcal{X}} |\nabla f_i(x)| \rho_i(x) dx,$$

when  $f_i$  is a smooth enough function. The functional  $\widetilde{\text{TV}}_\varepsilon(f_i, \mu_i)$  is thus connected to more standard notions of (weighted) total variation in Euclidean space. This heuristic can be formalised further via variational tools, as has been done recently in Ref. [16].

Total variation regularisation with general TV functionals is an important methodology in imaging, and also in unsupervised and supervised learning on graphs, where it has been used for community detection, clustering, and graph trend-filtering; e.g., see [5, 14, 17, 18, 23, 25, 27, 29, 31, 37] and references therein.

### 3. Distributional-perturbing model and its generalised barycenter formulation

In this section we introduce some tools and develop a collection of technical results that we use in Section 4 when proving Theorem 2.5.

#### 3.1 Generalized barycenter and MOT problems

In our work [21], we introduced the following *generalised barycenter problem*. Given  $\mu \in \mathcal{P}(\mathcal{Z})$ , we consider the optimisation problem

$$\inf_{\lambda, \tilde{\mu}_1, \dots, \tilde{\mu}_K} \left\{ \lambda(\mathcal{Z}) + \sum_{i \in \mathcal{Y}} C(\mu_i, \tilde{\mu}_i) : \lambda \geq \tilde{\mu}_i \text{ for all } i \in \mathcal{Y} \right\}. \quad (11)$$

In the above, the infimum is taken over positive (Borel) measures  $\tilde{\mu}_1, \dots, \tilde{\mu}_K$  and  $\lambda$  satisfying the constraints  $\lambda \geq \tilde{\mu}_i$  for all  $i \in \mathcal{Y}$ . This constraint must be interpreted as:  $\lambda(A) \geq \tilde{\mu}_i(A)$  for all  $A \in \mathfrak{B}(\mathcal{X})$ . Problem (11) can be understood as a generalisation of the standard (Wasserstein) barycenter problem studied in [1]. Indeed, if all measures  $\mu_1, \dots, \mu_K$  had the same total mass and the term  $\lambda(\mathcal{X})$  in (11) was rescaled by a constant  $\alpha \in (0, \infty)$ , then, as  $\alpha \rightarrow \infty$ , the resulting problem would recover the classical barycenter problem with pairwise cost function  $c$ . As stated, one can regard (11) as a partial optimal transport barycenter problem: we transport each  $\mu_i$  to a part of  $\lambda$  while requiring the transported masses to overlap as much as possible (this is enforced by asking for the term  $\lambda(\mathcal{X})$  to be small). For the canonical example  $c$  defined in (7), allowable  $\tilde{\mu}_i$  is a positive measure with the same mass of  $\mu_i$  such that its support is away from the support of  $\mu_i$  as most  $\varepsilon$  almost-surely (otherwise, (11) becomes  $\infty$ ). Then, as long as  $C(\mu_i, \tilde{\mu}_i) = 0$ , the optimal strategy to minimise (11) is to reduce  $\lambda(\mathcal{X})$ , i.e. to enlarge the overlap between  $\tilde{\mu}_i$  and  $\tilde{\mu}_j$  as much as possible, which engenders smaller  $\lambda$ .

We recall a result from Ref. [21], which essentially states that the generalised barycenter Problem (11) is dual to (5).

**Theorem 3.1.** [21, Proposition 7 and Corollary 32]. *Suppose that  $c$  satisfies Assumption 2.3. Then*

$$(5) = 1 - (11).$$

*Furthermore, the infimum of (11) is attained. In other words, there exists  $(\lambda^*, \tilde{\mu}^*)$  which minimises (11).*

Like classical barycenter problems, (11) has an equivalent MOT formulation. To be precise, we use a *stratified* multimarginal optimal transport problem to obtain an equivalent reformulation of (11).

**Theorem 3.2.** [21, Proposition 14 and Proposition 15]. *Suppose that  $c$  satisfies Assumption 2.3. Let  $S_K := \{A \subseteq \mathcal{Y} : A \neq \emptyset\}$ . Given  $A \in S_K$ , define  $c_A : \mathcal{X}^K \rightarrow [0, \infty]$  as  $c_A(x_1, \dots, x_K) := \inf_{x' \in \mathcal{X}} \sum_{i \in A} c(x', x_i)$ .*

*Consider the problem:*

$$\begin{aligned} & \inf_{\{\pi_A : A \in S_K\}} \sum_{A \in S_K} \int_{\mathcal{X}^K} (c_A(x_1, \dots, x_K) + 1) d\pi_A(x_1, \dots, x_K) \\ & \text{s.t. } \sum_{A \in S_K(i)} \mathcal{P}_i \# \pi_A = \mu_i \text{ for all } i \in \mathcal{Y}, \end{aligned} \tag{12}$$

*where  $\mathcal{P}_i$  is the projection map  $\mathcal{P}_i : (x_1, \dots, x_K) \mapsto x_i$ , and  $S_K(i) := \{A \in S_K : i \in A\}$ . Then (11) = (12). Also, the infimum in (12) is attained.*

**Remark 3.3.** Even though  $c_A$  and  $\pi_A$  above are defined over  $\mathcal{X}^K$ , only the coordinates  $i$  where  $i \in A$  actually plays a role in the optimisation problem. Also, notice that (12) is not a standard MOT problem since in (12) we optimise over several couplings  $\pi_A$  (each with its own cost function  $c_A$ ) that are connected to each other via the marginal constraints. We refer to this type of problem as a stratified MOT problem.

In the following theorem, we discuss the duals of the generalised barycenter problem and its MOT formulation. The notions of  $c$ -transform and  $\bar{c}$ -transform, whose definition we revisit in Appendix B, play an important role in these results.

**Theorem 3.4.** [21, Proposition 22 and Proposition 24]. *Suppose that  $c$  satisfies Assumption 2.3. Let  $\mathcal{C}_b(\mathcal{X})$  be the set of bounded real-valued continuous functions over  $\mathcal{X}$ . The dual of (11) is*

$$\begin{aligned} & \sup_{f_1, \dots, f_K \in \mathcal{C}_b(\mathcal{X})} \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} f_i^c(x_i) d\mu_i(x_i) \\ & \text{s.t. } f_i(x) \geq 0, \quad \sum_{i \in \mathcal{Y}} f_i(x) \leq 1, \text{ for all } x \in \mathcal{X}, \quad i \in \{1, \dots, K\}, \end{aligned} \tag{13}$$

*and there is no duality gap between primal and dual problems. In other words, (11) = (13). In the above,  $f_i^c$  denotes the  $c$ -transform of  $f_i$  as introduced in Definition 3.*

The dual of (12) is

$$\begin{aligned} & \sup_{g_1, \dots, g_K \in \mathcal{C}_b(\mathcal{X})} \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i(x_i) d\mu_i(x_i) \\ \text{s.t. } & \sum_{i \in A} g_i(x_i) \leq 1 + c_A(x_1, \dots, x_K) \text{ for all } (x_1, \dots, x_K) \in \mathcal{X}^K, \quad A \in S_K, \end{aligned} \quad (14)$$

and there is no duality gap between primal and dual problems. In other words, (12)=(14).

If, in addition, the cost function  $c$  is bounded and Lipschitz, then (14) is achieved by  $g \in \mathcal{C}_b(\mathcal{X})^K$ . Also, for  $f$  feasible for (13),  $g' := f^c$  is feasible for (14). Similarly, for  $g$  feasible for (14),  $f' = \max\{g, 0\}^{\bar{c}}$  is feasible for (13). Therefore, the optimisation of (14) can be restricted to non-negative  $g$  satisfying  $g_i = \bar{g}_i^c$ , or  $0 \leq g_i \leq 1$  for all  $i \in \mathcal{Y}$ . The notion of  $\bar{c}$ -transform is also introduced in Definition 3.

**Remark 3.5.** By combining Theorem 3.1, Theorem 3.2 and Theorem 3.4, we conclude that 1-(14)=(5).

**Remark 3.6.** A standard argument in optimal transport theory shows that problem (14) is equivalent to

$$\sup_{g_1, \dots, g_K} \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i(x_i) d\mu_i(x_i), \quad (15)$$

where the sup is taken over all  $(g_1, \dots, g_K) \in \prod_{i \in \mathcal{Y}} L^\infty(\mathcal{X}; \mu_i)$  satisfying: for any  $A \in S_K$ ,

$$\sum_{i \in A} g_i(x_i) \leq 1 + c_A(x_1, \dots, x_K)$$

for  $\otimes_i \mu_i$ -almost every tuple  $(x_1, \dots, x_K)$ . Here,  $L^\infty(\mathcal{X}; \mu_i)$  is defined as:

$$L^\infty(\mathcal{X}; \mu_i) := \{f: \mathcal{X} \rightarrow \mathbb{R} : \text{measurable, ess sup}_{\mu_i} |f| < \infty\}.$$

Indeed, notice that since (14) has already been shown to be equal to (12), the claim follows from the observation that any feasible  $g_1, \dots, g_K$  for (15) satisfies the condition

$$\sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i(x_i) d\mu_i(x_i) \leq \sum_{A \in S_K} \int_{\mathcal{X}^K} (1 + c_A(x_1, \dots, x_K)) d\pi_A(x_1, \dots, x_K)$$

for every  $\{\pi_A\}_{A \in S_K}$  satisfying the constraints in (12).

### 3.2 Existence of optimal dual potentials $g$ for general lower semi-continuous costs

We already know from the last part in Theorem 3.4 that if  $c$  is bounded and Lipschitz, then there is a feasible  $g \in \mathcal{C}_b(\mathcal{X})^K$  that is optimal for (14). In this subsection, we prove an analogous existence result in the case of a general lower semi-continuous cost function  $c$  satisfying Assumption 2.3. More precisely, we prove existence of maximisers for (15). We start with an approximation result.

**Lemma 3.7.** *Let  $c$  be a cost function satisfying Assumption 2.3. For each  $n \in \mathbb{N}$ , let*

$$c_n(x, x') := \min\{\tilde{c}_n(x, x'), n\},$$

where

$$\tilde{c}_n(x, x') := \inf_{(\tilde{x}, \tilde{x}') \in \mathcal{X} \times \mathcal{X}} \{c(\tilde{x}, \tilde{x}') + nd(x, \tilde{x}) + nd(x', \tilde{x}')\}.$$

Then the following properties hold:

1.  $c_n$  is bounded and Lipschitz.
2.  $c_n \leq c_{n+1} \leq c$  and  $\tilde{c}_n \leq \tilde{c}_{n+1}$  for all  $n \in \mathbb{N}$ .
3.  $\lim_{n \rightarrow \infty} c_n(x, x') = c(x, x')$  for all  $(x, x') \in \mathcal{X} \times \mathcal{X}$ .

**Proof.** First, it is easy to see that  $\tilde{c}_n$  is Lipschitz. Truncating it by  $n$ , finally, we obtain a bounded Lipschitz cost function  $c_n$ . Hence, items (1) and (2) are straightforward to prove. To prove item (3), notice that due to the monotonicity of the cost functions, we know that  $\lim_{n \rightarrow \infty} c_n(x, x')$  exists in  $[0, \infty]$  and  $\lim_{n \rightarrow \infty} c_n(x, x') \leq c(x, x')$ . If  $\lim_{n \rightarrow \infty} c_n(x, x') = \infty$ , then we would be done. Hence, we may assume that  $\lim_{n \rightarrow \infty} c_n(x, x') < \infty$ . From the definition of  $c_n$ , it then holds that  $\lim_{n \rightarrow \infty} \tilde{c}_n(x, x') = \lim_{n \rightarrow \infty} c_n(x, x') < \infty$ . Let  $(x_n, x'_n) \in \mathcal{X} \times \mathcal{X}$  be such that

$$c(x_n, x'_n) + nd(x, x_n) + nd(x', x'_n) \leq \tilde{c}_n(x, x') + \frac{1}{n}.$$

Since  $c(x_n, x'_n) \geq 0$ , the above implies that  $\lim_{n \rightarrow \infty} d(x, x_n) = 0$  and  $\lim_{n \rightarrow \infty} d(x', x'_n) = 0$ . Indeed, if this was not the case, then we would contradict  $\lim_{n \rightarrow \infty} \tilde{c}_n(x, x') < \infty$ . By the lower semi-continuity of the cost function  $c$ , we then conclude that

$$\begin{aligned} c(x, x') &\leq \liminf_{n \rightarrow \infty} c(x_n, x'_n) \leq \liminf_{n \rightarrow \infty} c(x_n, x'_n) + nd(x, x_n) + nd(x', x'_n) \leq \liminf_{n \rightarrow \infty} \tilde{c}_n(x, x') \\ &= \lim_{n \rightarrow \infty} c_n(x, x') \leq c(x, x'), \end{aligned}$$

from where the desired claim follows.  $\square$

**Lemma 3.8.** *Let  $c$  be a cost function satisfying Assumption 2.3, and let  $c_n$  be the cost function defined in Lemma 3.7. For each  $A \in S_K$ , let*

$$c_{A,n}(x_A) := \inf_{x' \in \mathcal{X}} \sum_{i \in A} c_n(x', x_i), \text{ and } c_A(x_A) := \inf_{x' \in \mathcal{X}} \sum_{i \in A} c(x', x_i),$$

where we use the shorthand notation  $x_A = (x_i)_{i \in A}$ . Then  $c_{A,n}$  monotonically converges towards  $c_A$  pointwise for all  $A \in S_K$ , as  $n \rightarrow \infty$ .

**Proof.** Fix  $A \in S_K$  and  $x_A := (x_i)_{i \in A} \in \mathcal{X}^{|A|}$ . From Lemma 3.7, it follows  $c_{A,n} \leq c_{A,n+1} \leq c_A$ . Therefore, for a given  $x_A$ ,  $\lim_{n \rightarrow \infty} c_{A,n}(x_A)$  exists in  $[0, \infty]$  and is less than or equal to  $c_A(x_A)$ . If the limit is  $\infty$ , we are done. We can then assume without the loss of generality that  $\lim_{n \rightarrow \infty} c_{A,n}(x_A) < \infty$ . We can then find sequences  $\{x_{n,i}\}_{n \in \mathbb{N}}$ ,  $\{x_{n,i}'\}_{n \in \mathbb{N}}$  and  $\{x_n'\}_{n \in \mathbb{N}}$  such that for all large enough  $n \in \mathbb{N}$

$$\sum_{i \in A} c(x_{n,i}', x_{n,i}) + n \left( \sum_{i \in A} (d(x_{n,i}, x_i) + d(x_{n,i}', x_n')) \right) \leq c_{A,n}(x_A) + \frac{1}{n}.$$

From the above, we derive that  $\lim_{n \rightarrow \infty} d(x_{n,i}', x_n') = 0$  and  $\lim_{n \rightarrow \infty} d(x_{n,i}, x_i) = 0$ . Hence, it follows that  $\limsup_{n \rightarrow \infty} c(x_{n,i}', x_{n,i}) < \infty$ . Combining the previous facts with Assumption 2.3, we conclude that  $\{x_n'\}_{n \in \mathbb{N}}$  is precompact, and thus, up to subsequence of  $n$ 's (that we do not relabel here), we have  $\lim_{n \rightarrow \infty} d(x_n', \hat{x}) = 0$  for some  $\hat{x} \in \mathcal{X}$ . Combining with  $\lim_{n \rightarrow \infty} d(x_{n,i}', x_n') = 0$ , we conclude that  $\lim_{n \rightarrow \infty} d(x_{n,i}', \hat{x}) = 0$  for all  $i \in A$ . Using the lower semi-continuity of  $c$ , we conclude that

$$c_A(x_A) \leq \sum_{i \in A} c(\hat{x}, x_i) \leq \liminf_{n \rightarrow \infty} \sum_{i \in A} c(x_{n,i}', x_{n,i}) \leq \lim_{n \rightarrow \infty} c_{A,n}(x_A) \leq c_A(x_A).$$

$\square$

**Proposition 3.9.** *Let  $c$  be a cost function satisfying Assumption 2.3. Then there exists a Borel solution for (15).*

**Proof.** Let  $\{c_n\}_{n \in \mathbb{N}}$  be the sequence of cost functions introduced in Lemma 3.7. Notice that for each  $n \in \mathbb{N}$ , there is a solution  $g^n = (g_1^n, \dots, g_K^n) \in \mathcal{C}_b(\mathcal{X})^K$  for the Problem (14) (with cost  $c_n$ ) that can be assumed to satisfy  $0 \leq g_i^n \leq 1$  for each  $i \in \mathcal{Y}$ . Therefore, for each  $i \in \mathcal{Y}$ , the sequence  $\{g_i^n\}_{n \in \mathbb{N}}$  is weakly\* precompact in  $L^\infty(\mathcal{X}; \mu_i)$  by Lemma A.2. This implies that there exists a subsequence of  $\{g^n\}_{n \in \mathbb{N}}$  (not relabelled) for which  $g^n$  weakly\* converges towards some  $g^* \in \prod_{i \in \mathcal{Y}} L^\infty(\mathcal{X}; \mathbb{R}, \mu_i)$ , which would necessarily satisfy  $0 \leq g_i^* \leq 1$  for all  $i \in \mathcal{Y}$ ; see Section A for the definition of weak\* topologies. We claim that this  $g^*$  is feasible for (15). Indeed, by Lemma 3.8, we know that  $c_{A,n} \leq c_A$  for all  $A \in S_K$ . In particular, since  $c_{A,n} \leq c_A$ , and  $\sum_{i \in A} g_i^n(x_i) \leq 1 + c_{A,n} \leq 1 + c_A$  for all  $A \in S_K$  and all  $n \in \mathbb{N}$ , it follows that

$\sum_{i \in \mathcal{A}} g_i^*(x_i) \leq 1 + c_A$ ,  $\otimes_i \mu_i$ -almost everywhere, due to the weak\* convergence of  $g_i^n$  towards  $g_i^*$ . This verifies that  $g^*$  is indeed feasible for (15).

Let  $\alpha_n$  and  $\beta_n$  be the optimal values of (11) and (14), respectively, for the cost  $c_n$ . Likewise, let  $\alpha$  and  $\beta$  be the optimal values of (11) and (14), respectively, for the cost  $c$ . Recall that, thanks to Theorem 3.2 and Theorem 3.4, we have  $\alpha_n = \beta_n$  for all  $n \in \mathbb{N}$  and  $\alpha = \beta$ . Suppose for a moment that we have already proved that  $\lim_{n \rightarrow \infty} \alpha_n = \alpha$ . Then we would have

$$\sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i^*(x) d\mu_i(x) = \lim_{n \rightarrow \infty} \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i^n(x) d\mu_i(x) = \lim_{n \rightarrow \infty} \beta_n = \lim_{n \rightarrow \infty} \alpha_n = \alpha,$$

which would imply that  $g^*$  is optimal for (15).

It thus remains to show that  $\lim_{n \rightarrow \infty} \alpha_n = \alpha$ . Given that  $c_n \leq c_{n+1} \leq c$ , it follows that  $\alpha_n \leq \alpha_{n+1} \leq \alpha$ . In particular, the limit  $\lim_{n \rightarrow \infty} \alpha_n$  exists in  $[0, \infty]$  and must satisfy  $\lim_{n \rightarrow \infty} \alpha_n \leq \alpha$ . If the limit is  $\infty$ , then there is nothing to prove. Thus, we can assume without the loss of generality that  $\alpha_\infty := \lim_{n \rightarrow \infty} \alpha_n < \infty$ .

Let  $\lambda^n$  and  $\tilde{\mu}_1^n, \dots, \tilde{\mu}_K^n$  be an optimal solution of (11) with the cost  $c_n$  and let  $\pi_i^n$  be a coupling realising  $C(\mu_i, \tilde{\mu}_i^n)$ . We first claim that  $\{\tilde{\mu}_i^n\}_{n \in \mathbb{N}}$  is weakly precompact for each  $i \in \mathcal{Y}$ . To see this, notice that for every  $n$ , we have  $\tilde{\mu}_i^n(\mathcal{X}) = \mu_i(\mathcal{X}) \leq 1$ , for otherwise,  $C(\mu_i, \tilde{\mu}_i^n) = \infty$ . Thus, by Prokhorov's theorem, it is enough to show that for every  $\eta > 0$ , there exists a compact set  $\mathcal{K} \subseteq \mathcal{X}$  such that  $\tilde{\mu}_i^n(\mathcal{X} \setminus \mathcal{K}) \leq C\eta$  for all  $n \in \mathbb{N}$  and some  $C$  independent of  $n, \eta$  or  $\mathcal{K}$ . To see that this is true, let us start by considering a compact set  $G$  such that  $\mu_i(G^c) \leq \eta$ . Let  $n_0 \in \mathbb{N}$  be such that  $n_0 - 1 > \frac{1}{\eta}$ . For  $n \geq n_0$ , we have

$$\alpha_\infty \geq \alpha_n = \lambda_n(\mathcal{X}) + \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} c_n(x_i, \tilde{x}_i) d\pi_i^n(x_i, \tilde{x}_i) \geq \int_G \int_{\mathcal{X}} c_{n_0}(x_i, \tilde{x}_i) d\pi_i^n(x_i, \tilde{x}_i).$$

Consider the set

$$\tilde{\mathcal{K}} := \{x \in \mathcal{X} \text{ s.t. } \inf_{\tilde{x} \in G} c_{n_0}(x, \tilde{x}) \leq n_0 - 1\};$$

using the definition of  $c_{n_0}$  and Assumption 2.3, it is straightforward to show that  $\tilde{\mathcal{K}}$  is a compact subset of  $\mathcal{X}$ . Since  $n_0 - 1 > \frac{1}{\eta}$ , we see that  $\alpha_\infty \geq \frac{1}{\eta}(\tilde{\mu}_i^n(\tilde{\mathcal{K}}^c) - \mu_i(G^c))$ , from where we can conclude that  $\tilde{\mu}_i^n(\tilde{\mathcal{K}}^c) \leq (\alpha_\infty + 1)\eta$  for all  $n \geq n_0$ . We now consider a compact set  $\hat{\mathcal{K}}$  for which  $\tilde{\mu}_i^n(\hat{\mathcal{K}}^c) \leq \eta$  for all  $n = 1, \dots, n_0$ , and set  $\mathcal{K} := \tilde{\mathcal{K}} \cup \hat{\mathcal{K}}$ , which is compact. Then, for all  $n \in \mathbb{N}$ , we have  $\tilde{\mu}_i^n(\mathcal{K}^c) \leq (\alpha_\infty + 1)\eta$ . This proves the desired claim.

Now, without the loss of generality, we can assume that  $\lambda^n$  has the form

$$d\lambda^n(x) = \max_{i=1, \dots, K} \left\{ \frac{d\tilde{\mu}_i^n}{d\bar{\mu}^n}(x) \right\} d\bar{\mu}^n(x),$$

where  $\bar{\mu}^n(x) = \sum_{i=1}^K \tilde{\mu}_i^n$ . Indeed, notice that the above is the smallest positive measure greater than  $\tilde{\mu}_1^n, \dots, \tilde{\mu}_K^n$ . Given the form of  $\lambda^n$  and the weak precompactness of each of the sequences  $\{\tilde{\mu}_i^n\}_{n \in \mathbb{N}}$ , we can conclude that  $\{\lambda^n\}_{n \in \mathbb{N}}$  is weakly precompact and so are the sequences  $\{\pi_i^n\}_{n \in \mathbb{N}}$ . We can thus assume that, up to subsequence,  $\tilde{\mu}_i^n$  converges weakly towards some  $\tilde{\mu}_i$ ;  $\pi_i^n$  converges weakly towards some  $\pi_i \in \Gamma(\mu_i, \tilde{\mu}_i)$  and  $\lambda^n$  converges weakly towards some  $\lambda$  satisfying  $\lambda \geq \tilde{\mu}_i$  for each  $i \in \mathcal{Y}$ . In particular,  $\lambda, \tilde{\mu}_1, \dots, \tilde{\mu}_K$  is feasible for (11).

Therefore, for all  $n_0 \in \mathbb{N}$ , we have

$$\begin{aligned} \alpha \geq \alpha_\infty &= \lim_{n \rightarrow \infty} \left( \lambda^n(\mathcal{X}) + \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} c_n(x_i, \tilde{x}_i) d\pi_i^n(x_i, \tilde{x}_i) \right) \\ &\geq \lim_{n \rightarrow \infty} \left( \lambda^n(\mathcal{X}) + \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} c_{n_0}(x_i, \tilde{x}_i) d\pi_i^n(x_i, \tilde{x}_i) \right) \\ &\geq \lambda(\mathcal{X}) + \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} c_{n_0}(x_i, \tilde{x}_i) d\pi_i(x_i, \tilde{x}_i). \end{aligned}$$

Sending  $n_0 \rightarrow \infty$ , we can then use the monotone convergence theorem to conclude that

$$\alpha \geq \alpha_\infty \geq \lambda(\mathcal{X}) + \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} c(x_i, \tilde{x}_i) d\pi_i(x_i, \tilde{x}_i) \geq \lambda(\mathcal{X}) + \sum_{i \in \mathcal{Y}} C(\mu_i, \tilde{\mu}_i) \geq \alpha.$$

This proves that  $\alpha_\infty = \alpha$ . □

### 3.3 From dual potentials to robust classifiers for continuous cost functions

Having discussed the existence of solutions  $g^*$  for (15), we move on to discussing the connection between  $g^*$  and solutions  $f^*$  of Problem (5).

**Proposition 3.10** [21, Originally in Corollary 33]. [22, see correction in Corollary 4.7 and Remark 4.9] *Let  $c: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$  be a lower semi-continuous function and suppose that  $(\tilde{\mu}^*, g^*)$  is a solution pair for the generalised barycenter problem (11) and the dual of its MOT formulation (15). Let  $f^*$  be defined as:*

$$f_i^*(\tilde{x}) := \max \left\{ \sup_{x \in \text{spt}(\mu_i)} \{g_i^*(x) - c(x, \tilde{x})\}, 0 \right\}, \quad (16)$$

for each  $i \in \mathcal{Y}$ .

*If  $f^*$  is Borel measurable, then  $(f^*, \tilde{\mu}^*)$  is a saddle solution for the problem (5). In particular,  $f^*$  is a minimiser of (5).*

**Remark 3.11.** By the definition (16),  $0 \leq f_i^*(x) \leq 1$  since  $0 \leq g_i^*(x) \leq 1$ . However, it is not trivial that  $\sum f_i^*(x) \leq 1$ . In [22, Corollary 4.7] (originally [21, Corollary 33]), it is proved that  $\sum f_i^*(x) \leq 1$ , i.e.,  $f^* = (f_1^*, \dots, f_K^*)$  is feasible.

The reason why we cannot directly use Proposition 3.10 to prove existence of solutions to (5) for arbitrary  $c$  and  $\mu$  is because it is a priori not guaranteed that  $f_i^*$ , as defined in (16), is Borel measurable; notice that the statement in Proposition 3.10 is conditional. If  $\text{spt}(\mu_i)$  was finite for all  $i$ , then the Borel measurability of  $f_i^*$  would follow immediately from the fact that the maximum of finitely many lower semi-continuous functions is Borel; this is of course the case when working with empirical measures. Likewise, the Borel measurability of  $f_i^*$  is guaranteed when  $\mu$  is arbitrary and  $c$  is a bounded Lipschitz function (in fact, it is sufficient for the cost to be continuous), as is discussed in [39, Definitions 5.2 and 5.7 and Theorem 5.10]. However, nothing can be said about the Borel measurability of  $f_i^*$  without further information on  $g_i^*$  (which in general is unavailable) when  $c$  is only assumed to be lower semi-continuous (as is the case for the cost  $C$  from (7)) and  $\text{spt}(\mu_i)$  is an uncountable set.

Our strategy to prove Theorem 2.5 in Section 4.1 will be to approximate an arbitrary cost function  $c$  from below with a suitable sequence of bounded and Lipschitz cost functions  $c_n$  (the costs defined in Lemma 3.7), and, in turn, consider a limit of the robust classifiers  $f_n^*$  associated to each of the  $c_n$ . This limit (lim sup, to be precise) will be our candidate solution for (5).

## 4. Proofs of our main results

In this section, we prove the existence of a Borel measurable robust classifier for Problem (5) when  $c$  is an arbitrary lower semi-continuous cost function satisfying Assumption 2.3. We also establish the existence of minimisers of (9) and establish Theorem 2.8 and Corollary 2.11.

### 4.1 Well-posedness of the DRO model

**Proof of Theorem 2.5.** Let  $\{c_n\}_{n \in \mathbb{N}}$  be the sequence of cost functions converging to  $c$  from below defined in Lemma 3.7. For each  $n \in \mathbb{N}$ , we use Theorem 3.4 and let  $g^n = (g_1^n, \dots, g_K^n) \in \mathcal{C}_b(\mathcal{X})^K$  be a solution

of (14) with cost  $c_n$ ; recall that we can assume that  $0 \leq g_i^n \leq 1$ . In turn, we use  $g^n$  and the cost  $c_n$  to define  $f^n := (f_1^n, \dots, f_K^n)$  following (16). Since the  $g_i^n$  and  $c_n$  are continuous, and given that the pointwise supremum of a family of continuous functions is lower semi-continuous, we can conclude that  $f_i^n$  is lower semi-continuous and thus also Borel measurable for each  $n \in \mathbb{N}$ . Thanks to Proposition 3.10,  $f^n$  is optimal for (13) with cost function  $c_n$ .

From the proof of Proposition 3.9, we know that there exists a subsequence (that we do not relabel) such that the  $g_i^n$  converge in the weak\* topology, as  $n \rightarrow \infty$ , towards limits  $g_i^*$  that form a solution for (15) with cost  $c$ . Using this subsequence and recalling (16), we define  $f^* \in \mathcal{F}$  according to

$$f_i^*(\tilde{x}) := \limsup_{n \rightarrow \infty} f_i^n(\tilde{x}), \quad \tilde{x} \in \mathcal{X}. \quad (17)$$

Notice that each  $f_i^*$  is indeed Borel measurable since it is the lim sup of Borel measurable functions. In addition, notice that  $0 \leq f_i^* \leq 1$ , due to the fact that  $0 \leq f_i^n \leq 1$  for all  $i \in \mathcal{Y}$  and all  $n \in \mathbb{N}$ . We'll conclude by proving that  $f^*$  is a solution for (5).

Let  $\tilde{\mu} \in \mathcal{P}(\mathcal{Z})$  be an arbitrary Borel probability measure with  $C(\mu, \tilde{\mu}) < \infty$ . For each  $i \in \mathcal{Y}$ , let  $\pi_i$  be an optimal coupling realising the cost  $C(\mu_i, \tilde{\mu}_i)$ . Then

$$\begin{aligned} R(f^*, \tilde{\mu}) - C(\mu, \tilde{\mu}) &= 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} f_i^*(\tilde{x}) d\tilde{\mu}_i(\tilde{x}) - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) d\pi_i(x, \tilde{x}) \\ &= 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{X}} (f_i^*(\tilde{x}) + c(x, \tilde{x})) d\pi_i(x, \tilde{x}) \\ &= 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{X}} \left( \limsup_{n \rightarrow \infty} \max \left\{ \sup_{x' \in \text{spt}(\mu_i)} \{g_i^n(x') - c_n(x', \tilde{x})\}, 0 \right\} + c(x, \tilde{x}) \right) d\pi_i(x, \tilde{x}) \\ &\leq 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{X}} \left( \limsup_{n \rightarrow \infty} \sup_{x' \in \text{spt}(\mu_i)} \{g_i^n(x') - c_n(x', \tilde{x})\} + c(x, \tilde{x}) \right) d\pi_i(x, \tilde{x}) \end{aligned}$$

where the last inequality follows from the simple fact that  $-\max\{a, 0\} \leq -a$  for any  $a \in \mathbb{R}$ . Choosing  $x' = x$  in the sup term (notice that indeed  $x$  can be assumed to belong to  $\text{spt}(\mu_i)$  since  $\pi_i$  has first marginal equal to  $\mu_i$ ), and applying reverse Fatou's lemma, we find that

$$\begin{aligned} R(f^*, \tilde{\mu}) - C(\mu, \tilde{\mu}) &\leq 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{X}} \limsup_{n \rightarrow \infty} \{g_i^n(x) - c_n(x, \tilde{x}) + c_n(x, \tilde{x})\} d\pi_i(x, \tilde{x}) \\ &= 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{X}} \limsup_{n \rightarrow \infty} \{g_i^n(x)\} d\pi_i(x, \tilde{x}) \\ &= 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \limsup_{n \rightarrow \infty} \{g_i^n(x)\} d\mu_i(x) \\ &\leq 1 - \limsup_{n \rightarrow \infty} \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i^n(x) d\mu_i(x) \\ &= 1 - \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i^*(x) d\mu_i(x) \\ &= R_{DRO}^*(\varepsilon), \end{aligned}$$

where the third equality follows from the weak\* convergence of  $g_i^n$  towards  $g_i^*$  and the last equality follows from Remark 3.5 and the fact that  $g^*$  is a solution for (15) (combined with Remark 3.6). Taking the sup

over  $\tilde{\mu} \in \mathcal{P}(\mathcal{X})$ , we conclude that

$$\sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{X})} \{R(f^*, \tilde{\mu}) - C(\mu, \tilde{\mu})\} \leq R_{DRO}^*(\varepsilon),$$

and thus  $f^*$  is indeed a minimiser of (5).

Let now  $\tilde{\mu}^*$  be a solution of (11) (which exists due to Theorem 3.1). The fact that  $(\tilde{\mu}^*, f^*)$  is a saddle for (5) follows from the above computations and the fact that by Theorem 3.1 and in [21, Corollary 32], we have

$$R_{DRO}^*(\varepsilon) = \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{X})} \inf_{f \in \mathcal{F}} \{R(f, \tilde{\mu}) - C(\mu, \tilde{\mu})\} = \inf_{f \in \mathcal{F}} \{R(f, \tilde{\mu}^*) - C(\mu, \tilde{\mu}^*)\}.$$

□

The next proposition states that the function  $g_i^*$  constructed in the proof of Proposition 3.9 is a Borel measurable version of the  $c$ -transform of  $f_i^*$ , where  $f_i^*$  was defined in (17).

**Lemma 4.1.** *Let  $\{g^n\}_{n \in \mathbb{N}}$  and  $\{f^n\}_{n \in \mathbb{N}}$  be as in the proof of Theorem 2.5,  $g^*$  be the weak\* limit of the  $g^n$  and  $f^*$  be as defined in (17). Then, for every  $i \in \mathcal{Y}$ ,*

$$g_i^*(x) = \inf_{\tilde{x} \in \mathcal{X}} \{f_i^*(\tilde{x}) + c(x, \tilde{x})\} \quad (18)$$

for  $\mu_i$ -a.e.  $x \in \mathcal{X}$ . This statement must be interpreted as: the set in which (18) is violated is contained in a Borel measurable set with zero  $\mu_i$  measure.

**Proof.** From the proof of Theorem 2.5, it holds that for each  $i \in \mathcal{Y}$

$$\int_{\mathcal{X}} f_i^*(\tilde{x}) d\tilde{\mu}_i^*(\tilde{x}) + \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) d\pi_i^*(x, \tilde{x}) = \int_{\mathcal{X}} g_i^*(x) d\mu_i(x). \quad (19)$$

On the other hand, from the definition of  $f_i^n$ , it follows that

$$g_i^n(x) \leq f_i^n(\tilde{x}) + c_n(x, \tilde{x}), \quad \forall \tilde{x} \in \mathcal{X}, \text{ and } \mu_i\text{-a.e. } x \in \mathcal{X}.$$

We can then combine the above with Lemma A.3 to conclude that for  $\mu_i$ -a.e.  $x \in \mathcal{X}$  and every  $\tilde{x} \in \mathcal{X}$ , we have

$$g_i^*(x) \leq \limsup_{n \rightarrow \infty} g_i^n(x) \leq \limsup_{n \rightarrow \infty} f_i^n(\tilde{x}) + c_n(x, \tilde{x}) = f_i^*(\tilde{x}) + c(x, \tilde{x}).$$

Taking the inf over  $\tilde{x} \in \mathcal{X}$ , we conclude that for  $\mu_i$ -a.e.  $x \in \mathcal{X}$ , we have

$$g_i^*(x) \leq \inf_{\tilde{x} \in \mathcal{X}} \{f_i^*(\tilde{x}) + c(x, \tilde{x})\}. \quad (20)$$

From this and (19), we see that  $g_i^* \in L^1(\mu_i)$  and  $-f_i^* \in L^1(\tilde{\mu}_i)$  are optimal dual potentials for the optimal transport problem  $C(\mu_i, \tilde{\mu}_i^*)$ . If (20) did not hold with equality for  $\mu_i$ -a.e.  $x \in \mathcal{X}$ , then we would be able to construct a Borel measurable version  $h_i$  of the right-hand side of (20) (see Lemma C.1 in the Appendix), which would be strictly greater than  $g_i^*$  in a set of positive  $\mu_i$ -measure. In addition, we would have that  $(h_i, -f_i^*)$  is a feasible dual pair for the OT problem  $C(\mu_i, \tilde{\mu}_i)$ . However, the above would contradict the optimality of the dual potentials  $(g_i^*, -f_i^*)$ . We thus conclude that (20) holds with equality except on a set contained in a set of  $\mu_i$  measure zero. □

#### 4.2 Well-posedness of the closed-ball model (9)

**Proof of Theorem 2.7.** We actually prove that for arbitrary cost  $c$  satisfying Assumption 2.3, the solution  $f^*$  to (5) constructed in the proof of Theorem 2.5 is also a solution for the problem:

$$\inf_{f \in \mathcal{F}} \left\{ \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \sup_{\tilde{x} \in \mathcal{X}} \{1 - f_i(\tilde{x}) - c(x, \tilde{x})\} d\tilde{\mu}_i(x) \right\}. \quad (21)$$

Theorem 2.7 will then be an immediate consequence of this more general result when applied to  $c = c_\varepsilon$  as in (7).

Let  $f^*$  be the Borel solution of (5) constructed in the proof of Theorem 2.5. It suffices to show that for any  $f \in \mathcal{F}$

$$\sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \inf_{\tilde{x} \in \mathcal{X}} \{f_i^*(\tilde{x}) + c(x, \tilde{x})\} d\overline{\mu}_i(x) \geq \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \inf_{\tilde{x} \in \mathcal{X}} \{f_i(\tilde{x}) + c(x, \tilde{x})\} d\overline{\mu}_i(x).$$

Suppose for the sake of contradiction that the above inequality does not hold. Then there exists some  $\hat{f} \in \mathcal{F}$  which provides a strict inequality in the opposite direction. Now, on one hand, (18) of Lemma 4.1 implies

$$\sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \inf_{\tilde{x} \in \mathcal{X}} \{f_i^*(\tilde{x}) + c(x, \tilde{x})\} d\overline{\mu}_i(x) = \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i^*(x) d\mu_i(x).$$

On the other hand, by Lemma C.1, for each  $i \in \mathcal{Y}$ , there exists a Borel measurable function  $\hat{g}_i$  equal to  $\inf_{\tilde{x} \in \mathcal{X}} \{\hat{f}_i(\tilde{x}) + c(x, \tilde{x})\}$   $\overline{\mu}_i$ -almost everywhere. Let  $\hat{g} := (\hat{g}_1, \dots, \hat{g}_K)$ . Combining the existence of such  $\hat{g}$  with the above equation, and using (1), it follows that  $\hat{g}$  satisfies

$$\sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} g_i^*(x) d\mu_i(x) < \sum_{i \in \mathcal{Y}} \int_{\mathcal{X}} \hat{g}_i(x) d\mu_i(x). \quad (22)$$

Notice that for each  $A \in S_K$  and  $\otimes \overline{\mu}_i$ -almost everywhere  $x_1, \dots, x_K$ , we have

$$\sum_{i \in A} \inf_{\tilde{x} \in \mathcal{X}} \{\hat{f}_i(\tilde{x}) + c(x_i, \tilde{x})\} \leq \inf_{\tilde{x} \in \mathcal{X}} \left\{ \sum_{i \in A} \hat{f}_i(\tilde{x}) + c(x_i, \tilde{x}) \right\} \leq 1 + c_A(x_1, \dots, x_K).$$

From the above, we conclude that  $\hat{g}$  is feasible for (15). However, this and (22) combined contradict the fact that  $g^*$  is optimal for (15), as had been shown in Proposition 3.9.  $\square$

**Proof of Corollary 2.11.** It is straightforward to verify (e.g., see [15]) that for  $(f_1, 1 - f_1) \in \mathcal{F}$ , we can write

$$R_{closed}^*((f_1, 1 - f_1); \varepsilon) = \int_0^1 R_{closed}^*((\mathbb{1}_{\{f_1 \geq t\}}, \mathbb{1}_{\{f_1 \geq t\}^c}); \varepsilon) dt. \quad (23)$$

It is also straightforward to see that

$$R_{closed}^*((\mathbb{1}_{\{f_1 \geq t\}}, \mathbb{1}_{\{f_1 \geq t\}^c}); \varepsilon) = \int_{\mathcal{X}} \sup_{\tilde{x} \in \overline{B}_\varepsilon(x)} \mathbb{1}_{A^c}(\tilde{x}) d\overline{\mu}_1(x) + \int_{\mathcal{X}} \sup_{\tilde{x} \in \overline{B}_\varepsilon(x)} \mathbb{1}_A(\tilde{x}) d\overline{\mu}_2(x).$$

Let  $(f_1, 1 - f_1)$  be a solution to (9) (which by Remark 2.2 can indeed be taken of this form). It follows from (23) that for almost every  $t \in [0, 1]$ , the pair  $(\mathbb{1}_{\{f_1 \geq t\}}, \mathbb{1}_{\{f_1 \geq t\}^c})$  is also a solution for that same problem and thus also for the problem restricted to hard classifiers. This proves the desired result.  $\square$

#### 4.3 Connection between closed-ball model and open-ball model

**Proof of Theorem 2.8.** One can easily observe that for any fixed  $\varepsilon > 0$  and  $\delta > 0$ , we have

$$\sup_{\tilde{x} \in \overline{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\} \leq \sup_{\tilde{x} \in \overline{B}_\varepsilon(x)} \{1 - f_i(\tilde{x})\} \leq \sup_{\tilde{x} \in B_{\varepsilon+\delta}(x)} \{1 - f_i(\tilde{x})\}$$

for all  $x \in \mathcal{X}$  and all  $f \in \mathcal{F}$ . This simple observation leads to  $R_{open}^*(f; \varepsilon) \leq R_{closed}^*(f; \varepsilon) \leq R_{open}^*(f; \varepsilon + \delta)$  for all  $f \in \mathcal{F}$ . Thus, we also have  $R_{open}^*(\varepsilon) \leq R_{closed}^*(\varepsilon) \leq R_{open}^*(\varepsilon + \delta)$ , and, in particular,  $R_{open}^*(\varepsilon) \leq R_{closed}^*(\varepsilon) \leq \liminf_{\delta \rightarrow 0} R_{open}^*(\varepsilon + \delta)$ . From the above, we can also see that the function  $\varepsilon \mapsto R_{open}^*(\varepsilon)$  is non-decreasing and, as such, is continuous for all but at most countably many values of  $\varepsilon > 0$ . Therefore, for all but at most countably many  $\varepsilon$ , we have  $R_{open}^*(\varepsilon) = R_{closed}^*(\varepsilon)$ .

Now, let  $f^*$  be solution of (9) and assume we have  $R_{open}^*(\varepsilon) = R_{closed}^*(\varepsilon)$ . Then

$$R_{open}^*(f^*; \varepsilon) \leq R_{closed}^*(f^*; \varepsilon) = R_{closed}^*(\varepsilon) = R_{open}^*(\varepsilon),$$

which means  $f^*$  is a solution of (3).  $\square$

## 5. Conclusion and future works

In this paper, we establish the equivalence of three popular models of multiclass adversarial training: the open ball model, the closed ball model and the DRO model, and, for the first time, (with the exception of partial results in [15]) we prove the existence of Borel measurable optimal robust classifiers in the agnostic-classifier setting for any number of classes. We are able to unify these models via a framework we have developed that connects these problems to optimal transport. Notably, our results show that it is unnecessary to grapple with the cumbersome machinery of universal sigma algebras, which was needed to prove existence of classifiers in past results.

Although our analysis sheds light on this area, many open questions still remain on both the theoretical and practical side. One of the most important practical questions is how to extend these results when the set of classifiers  $\mathcal{F}$  is some parametric family, for example, neural networks. In particular, one would like to specify the properties a parametric family must satisfy in order to approximate robust classifiers to some desired degree of accuracy. In the case of neural networks, one might ask for the number of neurons or number of layers that are required for robust classification.

Related to the above practical question is the following geometric/theoretical question: given an optimal robust classifier  $f^*$ , can we give a characterisation of the regularity of  $f^*$  as in [15]? In particular, one would like to quantify the smoothness of the interfaces between the different classes. In general, we cannot guarantee that  $f^*$  is a hard classifier, thus, this problem is best posed as a question about the smoothness of the level sets of  $f^*$ . Since optimal classifiers need not be unique, one can also pose the more general question of when it is possible to find at least one optimal Borel robust classifier with some specified regularity property. Due to the connection between approximation and regularity, answering this question will provide insights to the previous question of how well one can approximate optimal robust classifiers using certain parametric families.

A final question is how to extend our framework to other more general settings. In this paper, we have assumed throughout that the loss function is the 0–1 loss. However, most practitioners prefer strongly convex loss functions, for example, the cross entropy function, which allows for faster optimisation and has other desirable properties. As a result, one would like to establish the analogue of these results in this more general setting. This would be crucial for bringing these theoretical insights closer to the models favoured by working practitioners.

**Financial support.** NGT is supported by the NSF grants DMS-2005797 and DMS-2236447. MJ is supported by the NSF grant DMS-2400641. JK thanks to PIMS Kantorovich Initiative supported through a PIMS PRN and NSF-DMS 2133244.

**Competing interests.** The authors declare none.

## References

- [1] Aguech, M. & Carlier, G. (2011) Barycenters in the wasserstein space. *SIAM J. Math. Anal.* **43**(2), 904–924.
- [2] Awasthi, P., Frank, N. & Mohri, M. (2021) On the existence of the adversarial bayes classifier. *Adv. Neural Inform. Process. Syst.* **34**, 2978–2990.
- [3] Awasthi, P., Frank, N. S. & Mohri, M. (2021) On the existence of the adversarial bayes classifier. In: Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P. S. & Wortman Vaughan, J. (editors), *Advances in Neural Information Processing Systems*, Vol. 34, Curran Associates, Inc, 2978–2990. [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/172ef5a94b4dd0aa120c6878fc29f70c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/172ef5a94b4dd0aa120c6878fc29f70c-Paper.pdf)
- [4] Balcan, M.-F., Pukdee, R., Ravikumar, P. & Zhang, H. Nash equilibria and pitfalls of adversarial training in adversarial robustness games. In F. Ruiz, J. Dy, and J.-W. van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, 206 *Proceedings of Machine Learning Research*, 9607–9636. PMLR, (25–27 Apr, 2023).
- [5] Bertozzi, A. L. & Flenner, A. (2016) Diffuse interface models on graphs for classification of high dimensional data. *SIAM Rev.* **58** (2), 293–328.
- [6] Bertsekas, D. P. & Shreve, S. E. (1978). *Stochastic Optimal Control: The Discrete-Time Case*, USA, Academic Press, Inc.
- [7] Bhagoji, A. N., Cullina, D. & Mittal, P. (2019). Lower bounds on adversarial robustness from optimal transport. In: Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. & Garnett, R. (editors), *Advances in Neural Information Processing Systems*, Vol. 32, Curran Associates, Inc. [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/02bf86214e264535e3412283e817deaa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/02bf86214e264535e3412283e817deaa-Paper.pdf)

- [8] Biggio, B. & Roli, F. (2018) Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recogn.* **84**, 317–331.
- [9] Blanchet, J., Kang, Y. & Murthy, K. (2019) Robust Wasserstein profile inference and applications to machine learning. *J. Appl. Probab.* **56** (3), 830–857.
- [10] Blanchet, J., Kuhn, D., Li, J. & Taskesen, B. Unifying distributionally robust optimization via optimal transport theory. arXiv preprint arXiv: 2308.05414
- [11] Blanchet, J. & Murthy, K. (2019) Quantifying distributional model risk via optimal transport. *Math. Oper. Res.* **44**(2), 565–600.
- [12] Bose, J., Gidel, G., Berard, H., et al (2020) Adversarial example games. *Adv. Neur. Inform. Process. Syst.* **33**, 8921–8934.
- [13] Bousquet, O., Boucheron, S. & Lugosi, G. (2004). *Introduction to Statistical Learning Theory*, Springer, Berlin Heidelberg, pp. 169–207.
- [14] Boyd, Z. M., Porter, M. A. & Bertozzi, A. L. (2020) Stochastic block models are a discrete surface tension. *J. Nonlinear Sci.* **30** (5), 2429–2462.
- [15] Bungert, L., Trillos, N., García & Murray, R. (2023) *The geometry of adversarial training in binary classification* Information and Inference. *JIMA* **12** (2), 921–968.
- [16] Bungert, L. & Stinson, K. (2022) Gamma-convergence of a nonlocal perimeter arising in adversarial machine learning. arXiv preprint arXiv:2211.15223.
- [17] Caroccia, M., Chambolle, A. & Slepčev, D. (2020) Mumford–Shah functionals on graphs and their asymptotics. *Nonlinearity* **33** (8), 3846.
- [18] Cristofari, A., Rinaldi, F. & Tudisco, F. (2020) Total variation based community detection using a nonlinear optimization approach. *SIAM J. Appl. Math.* **80**(3), 1392–1419.
- [19] Frank, N. & Niles-Weed, J. (2023). The adversarial consistency of surrogate risks for binary classification. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M. & Levine, S. (editors), *Advances in Neural Information Processing Systems*, Vol. 36, Curran Associates, Inc, 41343–41354. [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/8185858b55a8c63763cfe088090242a-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/8185858b55a8c63763cfe088090242a-Paper-Conference.pdf)
- [20] Frank, N. S. & Niles-Weed, J. (2024). Existence and minimax theorems for adversarial surrogate risks in binary classification. *J. Mach. Learn. Res.* **25**(58), 1–41. <http://jmlr.org/papers/v25/23-0456.html>
- [21] García Trillo, N., Jacobs, M. & Kim, J. (2023) The multimarginal optimal transport formulation of adversarial multiclass classification. *J. Mach. Learn. Res.* **24** (45), 1–56.
- [22] García Trillo, N., Jacobs, M. & Kim, J. (2023) The multimarginal optimal transport formulation of adversarial multiclass classification. arXiv preprint, <https://arxiv.org/abs/2204.12676>
- [23] García Trillo, N. & Murray, R. (2017) A new analytical approach to consistency and overfitting in regularized empirical risk minimization. *Eur. J. Appl. Math.* **28** (6), 886–921.
- [24] García Trillo, N. & Murray, R. (2022) Adversarial classification: Necessary conditions and geometric flows. *J. Mach. Learn. Res.* **23**(187), 1–38.
- [25] García Trillo, N., Murray, R. & Thorpe, M. (2022) From graph cuts to isoperimetric inequalities: Convergence rates of Cheeger cuts on data clouds. *Arch. Ration. Mech. Anal.* **3**, 541–598.
- [26] Goodfellow, I., Shlens, J. & Szegedy, C. (2015). Explaining and harnessing adversarial examples. In International Conference on Learning Representations.
- [27] Hu, H., Laurent, T., Porter, M. A. & Bertozzi, A. L. (2013) A method based on total variation for network modularity optimization using the MBO scheme. *SIAM J. Appl. Math.* **73**(6), 2224–2246.
- [28] Javanmard, A., Soltanolkotabi, M. & Hassani, H. (2020). Precise tradeoffs in adversarial training for linear regression. In: Conference on Learning Theory, PMLR, 2034–2078.
- [29] Luo, X. & Bertozzi, A. L. (2017) Convergence of the graph Allen–Cahn scheme. *J. Stat. Phys.* **167** (3), 934–958.
- [30] Luzin, N. N. & Sierpiński, W. (1919). Sur quelques propriétés des ensembles (a).
- [31] Merkurjev, E., Kostić, T. & Bertozzi, A. L. (2013, 1903–1930) An MBO scheme on graphs for classification and image processing. *Siam J. Imaging Sci.* **6** (4), 1903–1930.
- [32] Meunier, L., Scetbon, M., Pinot, R. B., Atif, J. & Chevaleyre, Y. (2021) Mixed nash equilibria in the adversarial examples game. In: Meila, M., & Zhang, T. (editors), *Proceedings of the 38th International Conference on Machine Learning*, Vol. 139, *Proceedings of Machine Learning Research*, 7677–7687. PMLR, 18–24 Jul, 2021.
- [33] Nishiura, T. (2008). *Absolute Measurable Spaces*. 120 *Encyclopedia of Mathematics and Its Applications*, xiv+274, Cambridge University Press, Cambridge.
- [34] Pydi, M. S. & Jog, V. (2021) Adversarial risk via optimal transport and optimal couplings. *IEEE Trans. Inform Theory* **67**(9), 6031–6052.
- [35] Pydi, M. S. & Jog, V. (2021). The many faces of adversarial risk. In: Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P. & Vaughan, J. W. (editors), *Advances in Neural Information Processing Systems*, Vol. **34**, Curran Associates, Inc, 10000–10012.
- [36] Trillo, N. G., Jacobs, M., Kim, J. & Werenski, M. (2024). An optimal transport approach for computing adversarial training lower bounds in multiclass classification.
- [37] van Gennip, Y., Guillen, N., Osting, B. & Bertozzi, A. L. (2014) Mean curvature, threshold dynamics, and phase field theory on finite graphs. *Milan J. Math.* **82**, 3–65.
- [38] Villani, C. (2003). Topics in optimal transportation, *Graduate Studies in Mathematics*, xvi+370, American Mathematical Society, Providence, RI

- [39] Villani, C. (2009). Optimal transport, *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*, *Xxii+973*, Springer-Verlag, Berlin, Old and new
- [40] von Luxburg, U. & Schölkopf, B. (2011). Statistical learning theory: Models, concepts, and results. In: Gabbay, D. M., Hartmann, S. & Woods, J. (editors), *Inductive Logic, 10 Handbook of the History of Logic*, North-Holland, 651–706.

## Appendix A. Weak<sup>\*</sup> topology

**Definition 2** (Weak<sup>\*</sup> topology). Let  $\mu = (\mu_1, \dots, \mu_K) \in \prod_{i=1}^K \mathcal{M}_+(\mathcal{X})$ . For a sequence  $\{h^n\}_{n \in \mathbb{N}} \subseteq \prod_{i \in \mathcal{Y}} L^\infty(\mathcal{X}; \mu_i)$ , we say that  $\{h^n\}$  weak<sup>\*</sup>-converges to  $h \in \prod_{i \in \mathcal{Y}} L^\infty(\mathcal{X}; \mu_i)$  if for any  $q \in \prod_{i \in \mathcal{Y}} L^1(\mathcal{X}; \mu_i)$ , it holds that

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} h_i^n(x) q_i(x) d\mu_i(x) = \int_{\mathcal{X}} h_i(x) q_i(x) d\mu_i(x) \quad (\text{A1})$$

for all  $i \in \mathcal{Y}$ .

**Remark A.1.** Note that for a Borel positive measure  $\rho$ , which is either finite or  $\sigma$ -finite over a Polish space, the dual of  $L^1(\rho)$  is  $L^\infty(\rho)$ , which justifies the definition (1).

The following lemma is the weak<sup>\*</sup> precompactness of the closed unit ball of the dual space  $\mathcal{X}'$ . In our case,  $\mathcal{X} = L^1(\mathcal{X}; \mu_i)$  and  $\mathcal{X}' = L^\infty(\mathcal{X}; \mu_i)$  with weak<sup>\*</sup> topology.

**Lemma A.2.** (Banach–Alaoglu theorem). *If  $\mathcal{X}$  is a normed space, then the closed unit ball in the continuous dual space  $\mathcal{X}'$  (endowed with its usual operator norm) is compact with respect to the weak<sup>\*</sup> topology.*

**Lemma A.3.** Suppose  $\{g_i^n\}_{n \in \mathbb{N}}$  is a sequence of measurable real-valued functions over  $\mathcal{X}$  satisfying  $0 \leq g_i^n \leq 1$  for every  $n \in \mathbb{N}$ . Suppose that  $g_i^n$  converges in the weak<sup>\*</sup> topology of  $L^\infty(\mathcal{X}; \mu_i)$  towards  $g_i$ , where  $\mu_i$  is a finite positive measure. Then, for  $\mu_i$ -a.e.  $x \in X$ , we have

$$\limsup_{n \rightarrow \infty} g_i^n(x) \geq g_i(x).$$

**Proof.** Let  $E$  be a measurable subset of  $\mathcal{X}$ . Then

$$\int_{\mathcal{X}} (\limsup_{n \rightarrow \infty} g_i^n(x) - g_i(x)) \mathbb{1}_E(x) d\mu_i(x) \geq \limsup_{n \rightarrow \infty} \int_{\mathcal{X}} (g_i^n(x) - g_i(x)) \mathbb{1}_E(x) d\mu_i(x) = 0,$$

by the reverse Fatou inequality and the assumption that the sequence  $\{g_i^n\}_{n \in \mathbb{N}}$  converges in the weak<sup>\*</sup> sense towards  $g_i$ . Since  $E$  was arbitrary, the result follows.  $\square$

## Appendix B. *c*-transform

*c*-transform has an important role in optimal transport theory. One can characterise an optimiser of a dual problem by iterating *c*-transform: see [38, 39] for more details.

**Definition 3** (39, Definition 5.2 and Definition 5.8). Let  $\mathcal{X}, \mathcal{X}'$  be measurable spaces, and let  $c: \mathcal{X} \times \mathcal{X}' \rightarrow (-\infty, \infty]$ . Given a measurable function  $h: \mathcal{X} \rightarrow \mathbb{R} \cup \{\infty, -\infty\}$ , its *c*-transform is defined as:

$$h^c(x') := \inf_{x \in \mathcal{X}} \{h(x) + c(x, x')\}.$$

Similarly, for  $g: \mathcal{X}' \rightarrow \mathbb{R} \cup \{\infty, -\infty\}$ , its  $\bar{c}$ -transform is defined as:

$$g^{\bar{c}}(x) := \sup_{x' \in \mathcal{X}'} \{g(x') - c(x, x')\}.$$

**Proposition B.1.** For any measurable functions  $h$  over  $\mathcal{X}$  and  $g$  over  $\mathcal{X}'$ , and cost function  $c: \mathcal{X} \times \mathcal{X}' \rightarrow (-\infty, \infty]$ , it holds that for every  $(x, x') \in \mathcal{X} \times \mathcal{X}'$ ,

$$h^c(x') - h(x) \leq c(x, x'), \quad g(x') - g^{\bar{c}}(x) \leq c(x, x').$$

**Theorem B.2.** [39, Theorem 5.10] Let  $\mathcal{X}$  be a Polish space and  $c(\cdot, \cdot)$  be a cost function bounded from below and lower semi-continuous. Then, for  $v, \tilde{v} \in \mathcal{P}(\mathcal{X})$ ,

$$\begin{aligned} \inf_{\pi_i \in \Gamma(v, \tilde{v})} \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) d\pi_i(x, \tilde{x}) &= \sup_{g_i, f_i \in \mathcal{C}_b, g_i - f_i \leq c} \left\{ \int_{\mathcal{X}} g_i(x) d\nu(x) - \int_{\mathcal{X}} f_i(\tilde{x}) d\tilde{v}(\tilde{x}) \right\} \\ &= \sup_{f_i \in L^1(\tilde{v})} \left\{ \int_{\mathcal{X}} (f_i)^c(x) d\nu(x) - \int_{\mathcal{X}} f_i(\tilde{x}) d\tilde{v}(\tilde{x}) \right\} \\ &= \sup_{g_i \in L^1(v)} \left\{ \int_{\mathcal{X}} g_i(x) d\nu(x) - \int_{\mathcal{X}} (g_i)^{\bar{c}}(\tilde{x}) d\tilde{v}(\tilde{x}) \right\}. \end{aligned}$$

Furthermore, the infimum is indeed a minimum. However, the supremum may not be achieved.

### Appendix C. Decomposition of universally measurable functions

The following lemma is a well known fact about measure theory. For the sake of completeness, we write the full proof.

**Lemma C.1.** Let  $\mathcal{X}$  be a polish space, and let  $\mu$  and  $\bar{\mu}$  be a Borel probability measure and its extension to the universal  $\sigma$ -algebra, respectively. Let  $f$  be a universally measurable function for which  $\int_{\mathcal{X}} |f(x)| d\bar{\mu}(x) < \infty$ . Then there exists a Borel measurable function  $g$  such that  $f = g$   $\bar{\mu}$ -almost everywhere. Also,

$$\int_{\mathcal{X}} f(x) d\bar{\mu}(x) = \int_{\mathcal{X}} g(x) d\mu(x). \quad (\text{C1})$$

**Proof.** Without the loss of generality, we can assume that  $f \geq 0$ . Since  $f$  is universally measurable, we can write

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) := \lim_{n \rightarrow \infty} \sum_{k=1}^n c_k^n \mathbb{1}_{A_k^n}(x),$$

for positive coefficients  $c_1^n, \dots, c_n^n$  and  $A_1^n, \dots, A_n^n$  universally measurable and pairwise disjoint sets. By the definition of universally measurable sets, for each  $A_k^n$ , there exists a Borel set  $B_k^n$  such that  $\bar{\mu}(A_k^n \setminus B_k^n) = 0$ . Hence, for each  $n \in \mathbb{N}$ , we can write

$$f_n(x) = \sum_{k=1}^n c_k^n \mathbb{1}_{B_k^n}(x) + \sum_{k=1}^n c_k^n \mathbb{1}_{C_k^n}(x),$$

where  $C_k^n = A_k^n \setminus B_k^n$ . We conclude that

$$f(x) = g(x) + h(x) := \limsup_{n \rightarrow \infty} \sum_{k=1}^n c_k^n \mathbb{1}_{B_k^n}(x) + \liminf_{n \rightarrow \infty} \sum_{k=1}^n c_k^n \mathbb{1}_{C_k^n}(x)$$

where  $g$  is Borel measurable,  $h$  is universally measurable and  $h = 0$   $\bar{\mu}$ -almost everywhere.

Since  $f = g$   $\bar{\mu}$ -almost everywhere and  $g$  is Borel measurable, then

$$\int_{\mathcal{X}} f(x) d\bar{\mu}(x) = \int_{\mathcal{X}} g(x) d\bar{\mu}(x) = \int_{\mathcal{X}} g(x) d\mu(x),$$

from which (1) follows.  $\square$