

JACOBIAN DESCENT FOR MULTI-OBJECTIVE OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Many optimization problems require balancing multiple conflicting objectives. As gradient descent is limited to single-objective optimization, we introduce its direct generalization: Jacobian descent (JD). This algorithm iteratively updates parameters using the Jacobian matrix of a vector-valued objective function, in which each row is the gradient of an individual objective. While several methods to combine gradients already exist in the literature, they are generally hindered when the objectives conflict. In contrast, we propose projecting gradients to fully avoid conflict while ensuring that they preserve an influence proportional to their norm. We prove significantly stronger convergence guarantees with this approach, supported by our empirical results. Our method also enables instance-wise risk minimization (IWRM), a novel learning paradigm in which the loss of each training example is considered a separate objective. Applied to simple image classification tasks, IWRM exhibits promising results compared to the direct minimization of the average loss. Additionally, we outline an efficient implementation of JD using the Gramian of the Jacobian matrix to reduce time and memory requirements.

1 INTRODUCTION

The field of multi-objective optimization studies minimization of vector-valued objective functions (Sawaragi et al., 1985; Ehrgott, 2005; Branke, 2008; Deb et al., 2016). In deep learning, a widespread approach to train a model with multiple objectives is to combine those into a scalar loss function minimized by stochastic gradient descent. While this method is simple, it comes at the expense of potentially degrading some individual objectives. Without prior knowledge of their relative importance, this is undesirable. In opposition, multi-objective optimization methods typically attempt to optimize all objectives simultaneously, without making arbitrary compromises: the goal is to find points for which no improvement can be made on some objectives without degrading others.

Early works have attempted to extend gradient descent (GD) to consider several objectives simultaneously, and thus several gradients (Fliege & Svaiter, 2000; Désidéri, 2012). Essentially, they propose a heuristic to prevent the degradation of any individual objective. Several other works have built upon this method, analyzing its convergence properties or extending it to a stochastic setting (Fliege et al., 2019; Poirion et al., 2017; Mercier et al., 2018). Later, this has been applied to multi-task learning to tackle conflict between tasks, illustrated by contradicting gradient directions (Sener & Koltun, 2018). Many studies have followed, proposing various other algorithms for the training of multi-task models (Yu et al., 2020; Liu et al., 2021a;b; Lin et al., 2021; Navon et al., 2022; Senushkin et al., 2023; Chen et al., 2020). They commonly rely on an aggregator that maps a collection of task-specific gradients (a Jacobian matrix) to a shared parameter update.

We propose to unify all such methods under the *Jacobian descent* (JD) algorithm, specified by an aggregator.¹ This algorithm aims to minimize a differentiable vector-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ iteratively without relying on a scalarization of the objective. Under this formulation, the existing methods are simply distinguished by their aggregator. Consequently, studying its properties is essential for understanding the behavior and convergence of JD. Under significant conflict, existing aggregators often fail to provide strong convergence guarantees. To address this, we propose $\mathcal{A}_{\text{UPGrad}}$, specifically designed to resolve conflicts while naturally preserving the relative influence of individual gradients.

¹Our library enabling JD with PyTorch is available at https://github.com/***/***

Furthermore, we introduce a novel stochastic variant of JD that enables the training of neural networks with a large number of objectives. This unlocks a particularly interesting perspective: considering the minimization of instance-wise loss vectors rather than the usual minimization of the average training loss. As this paradigm is a direct generalization of the well-known empirical risk minimization (ERM) (Vapnik, 1995), we name it *instance-wise risk minimization* (IWRM).

Our contributions are organized as follows: In Section 2, we formalize the JD algorithm and its stochastic variants. We then introduce three important aggregator properties and define $\mathcal{A}_{\text{UPGrad}}$ to satisfy them. In the smooth convex case, we show convergence of JD with $\mathcal{A}_{\text{UPGrad}}$ to the Pareto front. We present applications for JD and aggregators in Section 3, emphasizing the IWRM paradigm. We then discuss existing aggregators and analyze their properties in Section 4. In Section 5, we report experiments with IWRM optimized with stochastic JD with various aggregators. Lastly, we address computational efficiency in Section 6, giving a path towards an efficient implementation.

2 THEORETICAL FOUNDATION

A suitable partial order between vectors must be considered to enable multi-objective optimization. Throughout this paper, \leq denotes the relation defined for any pair of vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ as $\mathbf{u} \leq \mathbf{v}$ whenever $u_i \leq v_i$ for all coordinates i . Similarly, $<$ is the relation defined by $\mathbf{u} < \mathbf{v}$ whenever $u_i < v_i$ for all coordinates i . Furthermore, $\mathbf{u} \preceq \mathbf{v}$ indicates that both $\mathbf{u} \leq \mathbf{v}$ and $\mathbf{u} \neq \mathbf{v}$ hold. The Euclidean vector norm and the Frobenius matrix norm are denoted by $\|\cdot\|$ and $\|\cdot\|_F$, respectively. Finally, for any $m \in \mathbb{N}$, the symbol $[m]$ represents the range $\{i \in \mathbb{N} : 1 \leq i \leq m\}$.

2.1 JACOBIAN DESCENT

In the following, we introduce Jacobian descent, a natural extension of gradient descent supporting the optimization of vector-valued functions.

Suppose that $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuously differentiable. Let $\mathcal{J}\mathbf{f}(\mathbf{x}) \in \mathbb{R}^{m \times n}$ be the Jacobian matrix of \mathbf{f} at \mathbf{x} , i.e.

$$\mathcal{J}\mathbf{f}(\mathbf{x}) = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \nabla f_2(\mathbf{x})^\top \\ \vdots \\ \nabla f_m(\mathbf{x})^\top \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(\mathbf{x}) & \frac{\partial}{\partial x_2} f_1(\mathbf{x}) & \cdots & \frac{\partial}{\partial x_n} f_1(\mathbf{x}) \\ \frac{\partial}{\partial x_1} f_2(\mathbf{x}) & \frac{\partial}{\partial x_2} f_2(\mathbf{x}) & \cdots & \frac{\partial}{\partial x_n} f_2(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} f_m(\mathbf{x}) & \frac{\partial}{\partial x_2} f_m(\mathbf{x}) & \cdots & \frac{\partial}{\partial x_n} f_m(\mathbf{x}) \end{bmatrix} \quad (1)$$

Given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, Taylor’s theorem yields

$$\mathbf{f}(\mathbf{x} + \mathbf{y}) = \mathbf{f}(\mathbf{x}) + \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y} + o(\|\mathbf{y}\|), \quad (2)$$

where $o(\|\mathbf{y}\|)$ indicates that $\lim_{\|\mathbf{y}\| \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + \mathbf{y}) - \mathbf{f}(\mathbf{x}) - \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}}{\|\mathbf{y}\|} = \mathbf{0}$. The term $\mathbf{f}(\mathbf{x}) + \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}$ is the first-order Taylor approximation of $\mathbf{f}(\mathbf{x} + \mathbf{y})$. Via this approximation, we aim to select a small update \mathbf{y} that reduces $\mathbf{f}(\mathbf{x} + \mathbf{y})$, ideally achieving $\mathbf{f}(\mathbf{x} + \mathbf{y}) \leq \mathbf{f}(\mathbf{x})$. As the approximation depends on \mathbf{y} only through $\mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}$, selecting the update based on the Jacobian is natural. A mapping $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^n$ reducing such a matrix into a vector is called an *aggregator*. For any $J \in \mathbb{R}^{m \times n}$, $\mathcal{A}(J)$ is called the *aggregation* of J by \mathcal{A} .

To minimize \mathbf{f} , consider the update $\mathbf{y} = -\eta \mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x}))$, where η is an appropriate step size, and \mathcal{A} is an appropriate aggregator. Jacobian descent simply consists in applying this update iteratively, as shown in Algorithm 1. To put it into perspective, we also provide a minimal version of GD in Algorithm 2. Remarkably, when $m = 1$, the Jacobian has a single row, so GD is a special case of JD where the aggregator is the identity.

Algorithm 1: Jacobian descent with aggregator \mathcal{A}

Input: $\mathbf{x} \in \mathbb{R}^n, 0 < \eta, T \in \mathbb{N}, \mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^n$
for $t \leftarrow 1$ **to** T **do**
 $\mathbf{x} \leftarrow \mathbf{x} - \eta \mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x}))$
Output: \mathbf{x}

Algorithm 2: Gradient descent

Input: $\mathbf{x} \in \mathbb{R}^n, 0 < \eta, T \in \mathbb{N}$
for $t \leftarrow 1$ **to** T **do**
 $\mathbf{x} \leftarrow \mathbf{x} - \eta \nabla f(\mathbf{x})$
Output: \mathbf{x}

Note that other gradient-based optimization algorithms, e.g. Adam (Kingma & Ba, 2014), can similarly be extended to the multi-objective case.

In some settings, the exact computation of the update can be prohibitively slow or even intractable. When dealing with a single objective, stochastic gradient descent (SGD) replaces the gradient $\nabla f(\mathbf{x})$ with some estimation. More generally, *stochastic Jacobian descent* (SJD) relies on estimates of the aggregation of the Jacobian. One approach, that we call *stochastically estimated Jacobian descent* (SEJD), is to compute and aggregate an estimation of the Jacobian. Alternatively, when the number of objectives is very large, we propose to aggregate a matrix whose rows are a random subset of the rows of the true Jacobian. We call this approach *stochastic sub-Jacobian descent* (SSJD).

2.2 DESIRABLE PROPERTIES FOR AGGREGATORS

An inherent challenge of multi-objective optimization is to manage conflicting objectives (Sener & Koltun, 2018; Yu et al., 2020; Liu et al., 2021a). Substituting the update $\mathbf{y} = -\eta \mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x}))$ into the first-order Taylor approximation $\mathbf{f}(\mathbf{x}) + \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}$ yields $\mathbf{f}(\mathbf{x}) - \eta \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x}))$. In particular, if $\mathbf{0} \leq \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x}))$, then no coordinate of the approximation of \mathbf{f} will increase. A pair of vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ is said to *conflict* if $\mathbf{x}^\top \mathbf{y} < 0$. Hence, for a sufficiently small η , if any row of $\mathcal{J}\mathbf{f}(\mathbf{x})$ conflicts with $\mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x}))$, the corresponding coordinate of \mathbf{f} will increase. When minimizing \mathbf{f} , avoiding conflict between the aggregation and any gradient is thus desirable, motivating the first property.

Definition 1 (Non-conflicting). Let $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^n$ be an aggregator. If for all $J \in \mathbb{R}^{m \times n}$, $\mathbf{0} \leq J \cdot \mathcal{A}(J)$, then \mathcal{A} is said to be *non-conflicting*.

For any collection of vectors $C \subseteq \mathbb{R}^n$, the *dual cone* of C is $\{\mathbf{x} \in \mathbb{R}^n : \forall \mathbf{y} \in C, 0 \leq \mathbf{x}^\top \mathbf{y}\}$ (Boyd & Vandenberghe, 2004). Notice that an aggregator \mathcal{A} is non-conflicting if and only if for any J , $\mathcal{A}(J)$ is in the dual cone of the rows of J .

In a step of GD, the update scales proportionally to the gradient norm. Small gradients thus lead to small updates, and conversely, large gradients lead to large updates. To maintain coherence with GD, it would be natural that the rows of the Jacobian also contribute to the aggregation proportionally to their norm. Scaling each row of $\mathcal{J}\mathbf{f}(\mathbf{x})$ by the corresponding element of some vector $\mathbf{c} \in \mathbb{R}^m$ yields $\text{diag}(\mathbf{c}) \cdot \mathcal{J}\mathbf{f}(\mathbf{x})$. This insight can then be formalized as the following property.

Definition 2 (Linear under scaling). Let $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^n$ be an aggregator. If for all $J \in \mathbb{R}^{m \times n}$, the mapping from any $\mathbf{0} < \mathbf{c} \in \mathbb{R}^m$ to $\mathcal{A}(\text{diag}(\mathbf{c}) \cdot J)$ is linear in \mathbf{c} , then \mathcal{A} is said to be *linear under scaling*.

Finally, as $\|\mathbf{y}\|$ decreases asymptotically to 0, the precision of the first-order Taylor approximation $\mathbf{f}(\mathbf{x}) + \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}$ improves, as highlighted in (2). The projection \mathbf{y}' of any candidate update \mathbf{y} onto the span of the rows of $\mathcal{J}\mathbf{f}(\mathbf{x})$ satisfies $\mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}' = \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathbf{y}$ and $\|\mathbf{y}'\| \leq \|\mathbf{y}\|$, so this projection decreases the norm of the update while preserving the value of the approximation. Without additional information about \mathbf{f} , it is thus reasonable to select \mathbf{y} directly in the row span of $\mathcal{J}\mathbf{f}(\mathbf{x})$, i.e. to have a vector of weights $\mathbf{w} \in \mathbb{R}^m$ satisfying $\mathbf{y} = \mathcal{J}\mathbf{f}(\mathbf{x})^\top \cdot \mathbf{w}$. This yields the last desirable property.

Definition 3 (Weighted). Let $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^n$ be an aggregator. If for all $J \in \mathbb{R}^{m \times n}$, there exists $\mathbf{w} \in \mathbb{R}^m$ satisfying $\mathcal{A}(J) = J^\top \cdot \mathbf{w}$, then \mathcal{A} is said to be *weighted*.

2.3 UNCONFLICTING PROJECTION OF GRADIENTS

We now define the *unconflicting projection of gradients* aggregator $\mathcal{A}_{\text{UPGrad}}$, specifically designed to be non-conflicting, linear under scaling, and weighted. In essence, it projects each gradient onto the dual cone of the rows of the Jacobian and averages the results, as illustrated in Figure 1a.

For any $J \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$, the *projection* of \mathbf{x} onto the dual cone of the rows of J is

$$\pi_J(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathbb{R}^n : \mathbf{0} \leq J\mathbf{y}} \|\mathbf{y} - \mathbf{x}\|^2. \quad (3)$$

Denoting by $\mathbf{e}_i \in \mathbb{R}^m$ the i th standard basis vector, $J^\top \mathbf{e}_i$ is the i th row of J . $\mathcal{A}_{\text{UPGrad}}$ is defined as

$$\mathcal{A}_{\text{UPGrad}}(J) = \frac{1}{m} \sum_{i \in [m]} \pi_J(J^\top \mathbf{e}_i). \quad (4)$$

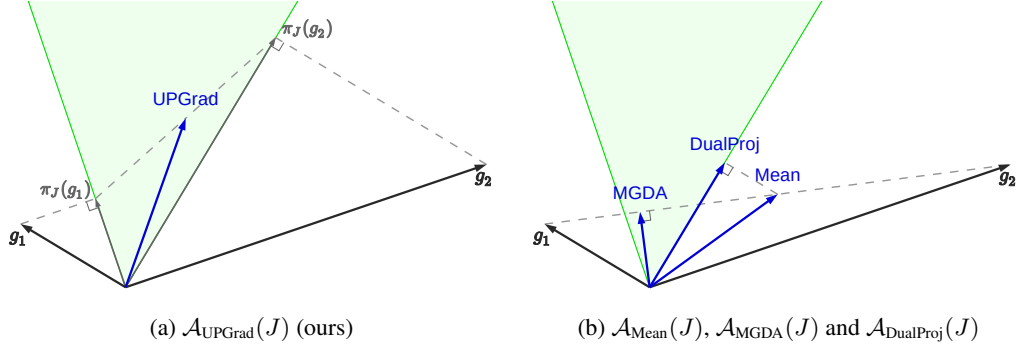


Figure 1: Aggregation of $J = [\mathbf{g}_1 \ \mathbf{g}_2]^\top \in \mathbb{R}^{2 \times 2}$ by four different aggregators. The dual cone of $\{\mathbf{g}_1, \mathbf{g}_2\}$ is represented in green.

(a) $\mathcal{A}_{\text{UPGrad}}$ projects \mathbf{g}_1 and \mathbf{g}_2 onto the dual cone and averages the results.

(b) The mean $\mathcal{A}_{\text{Mean}}(J) = \frac{1}{2}(\mathbf{g}_1 + \mathbf{g}_2)$ conflicts with \mathbf{g}_1 . $\mathcal{A}_{\text{DualProj}}$ projects this mean onto the dual cone, so it lies on its boundary. $\mathcal{A}_{\text{MGDA}}(J)$ is almost orthogonal to \mathbf{g}_2 because of its larger norm.

Since the dual cone is convex, it is closed under positive combinations of its elements. For any J , $\mathcal{A}_{\text{UPGrad}}(J)$ is thus always in the dual cone of the rows of J , so $\mathcal{A}_{\text{UPGrad}}$ is non-conflicting. Note that if no pair of gradients conflicts, $\mathcal{A}_{\text{UPGrad}}$ simply averages the rows of the Jacobian.

Since π_J is a projection onto a closed convex cone, if $\mathbf{x} \in \mathbb{R}^n$ and $0 < a \in \mathbb{R}$, then $\pi_J(a \cdot \mathbf{x}) = a \cdot \pi_J(\mathbf{x})$. By (4), $\mathcal{A}_{\text{UPGrad}}$ is thus linear under scaling.

When n is large, the projection in (3) is prohibitively expensive to compute. An alternative but equivalent approach is to use its dual formulation, which is independent of n .

Proposition 1. Let $J \in \mathbb{R}^{m \times n}$. For any $\mathbf{u} \in \mathbb{R}^m$, $\pi_J(J^\top \mathbf{u}) = J^\top \mathbf{w}$ with

$$\mathbf{w} \in \arg \min_{\mathbf{v} \in \mathbb{R}^n: \mathbf{u} \leq \mathbf{v}} \mathbf{v}^\top J J^\top \mathbf{v}. \quad (5)$$

Proof. See Appendix A.2.

The problem defined in (5) can be solved efficiently using a quadratic programming solver, such as those bundled in `qp solvers` (Caron et al., 2024). For any $i \in [m]$, let \mathbf{w}_i be given by (5) when substituting \mathbf{u} with \mathbf{e}_i . Then, by Proposition 1,

$$\mathcal{A}_{\text{UPGrad}}(J) = J^\top \left(\frac{1}{m} \sum_{i \in [m]} \mathbf{w}_i \right). \quad (6)$$

This provides an efficient implementation of $\mathcal{A}_{\text{UPGrad}}$ and proves that it is weighted. $\mathcal{A}_{\text{UPGrad}}$ can also be easily extended to incorporate a vector of preferences by replacing the average in (4) and (6) by a weighted sum with positive weights. This extension remains non-conflicting, linear under scaling, and weighted.

2.4 CONVERGENCE TO THE PARETO FRONT

We now provide theoretical convergence guarantees of JD with $\mathcal{A}_{\text{UPGrad}}$ when minimizing some $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ satisfying standard assumptions. If for a given $\mathbf{x} \in \mathbb{R}^n$, there exists no $\mathbf{y} \in \mathbb{R}^n$ satisfying $\mathbf{f}(\mathbf{y}) \preceq \mathbf{f}(\mathbf{x})$, then \mathbf{x} is said to be *Pareto optimal*. The set $X^* \subseteq \mathbb{R}^n$ of Pareto optimal points is called the *Pareto set*, and its image $\mathbf{f}(X^*)$ is called the *Pareto front*.

Whenever $\mathbf{f}(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda \mathbf{f}(\mathbf{x}) + (1 - \lambda)\mathbf{f}(\mathbf{y})$ holds for any pair of vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and any $\lambda \in [0, 1]$, \mathbf{f} is said to be \leq -convex. Moreover, \mathbf{f} is said to be β -smooth whenever $\|\mathcal{J}\mathbf{f}(\mathbf{x}) - \mathcal{J}\mathbf{f}(\mathbf{y})\|_{\text{F}} \leq \beta \|\mathbf{x} - \mathbf{y}\|$ holds for any pair of vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

Theorem 1. Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a β -smooth and \leq -convex function. Suppose that the Pareto front $\mathbf{f}(X^*)$ is bounded and that for any $\mathbf{x} \in \mathbb{R}^n$, there is $\mathbf{x}^* \in X^*$ satisfying $\mathbf{f}(\mathbf{x}^*) \leq \mathbf{f}(\mathbf{x})$.² Let $\mathbf{x}_1 \in \mathbb{R}^n$, and for all $t \geq 1$, $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathcal{A}_{\text{UPGrad}}(\mathcal{J}\mathbf{f}(\mathbf{x}_t))$, with $\eta = \frac{1}{\beta\sqrt{m}}$. Let \mathbf{w}_t be the weights defining $\mathcal{A}_{\text{UPGrad}}(\mathcal{J}\mathbf{f}(\mathbf{x}_t))$ as per (6), i.e. $\mathcal{A}_{\text{UPGrad}}(\mathcal{J}\mathbf{f}(\mathbf{x}_t)) = \mathcal{J}\mathbf{f}(\mathbf{x}_t)^\top \cdot \mathbf{w}_t$. If \mathbf{w}_t is bounded, then $\mathbf{f}(\mathbf{x}_t)$ converges to $\mathbf{f}(\mathbf{x}^*)$ for some $\mathbf{x}^* \in X^*$. In other words, $\mathbf{f}(\mathbf{x}_t)$ converges to the Pareto front.

Proof. See Appendix A.3.

Empirically, \mathbf{w}_t appears to converge to some $\mathbf{w}^* \in \mathbb{R}^m$ satisfying both $\mathbf{0} < \mathbf{w}^*$ and $\mathcal{J}\mathbf{f}(\mathbf{x}^*)^\top \mathbf{w}^* = \mathbf{0}$. This suggests that the boundedness of \mathbf{w}_t could be relaxed or even removed from the set of assumptions of Theorem 1.

Another commonly studied type of convergence for multi-objective optimization is convergence to a stationary point. If for a given $\mathbf{x} \in \mathbb{R}^n$, there exists $\mathbf{0} \preceq \mathbf{w}$ satisfying $\mathcal{J}\mathbf{f}(\mathbf{x})^\top \mathbf{w} = \mathbf{0}$ then \mathbf{x} is said to be *Pareto stationary*. Even though every Pareto optimal point is Pareto stationary, the converse does not hold, even in the convex case. The function $[x \ y]^\top \mapsto [x^2 \ y^2]^\top$ illustrates this discrepancy. Its Pareto set only contains the origin, but its set of Pareto stationary points is the union of the two axes. Despite being necessary, convergence to a Pareto stationary point is thus not a sufficient condition for optimality and, hence, constitutes a rather weak guarantee. To the best of our knowledge, $\mathcal{A}_{\text{UPGrad}}$ is the first non-conflicting aggregator that provably converges to the Pareto front in the smooth convex case.

In addition to the asymptotic convergence guarantees of Theorem 1, Appendix A.3 provides the following rate of convergence for any number of iterations $T \in \mathbb{N}$:

$$\frac{1}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \leq \frac{\sqrt{m}}{T} \left(\|\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)\| + \frac{\beta}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \right).$$

Although this result does not directly provide a convergence rate for \mathbf{f} , the bound $\mathbf{1} \leq \mathbf{w}_t$, suggests a convergence rate of $\mathcal{O}\left(\frac{1}{T}\right)$, offering valuable insight into the algorithm’s asymptotic behavior.

3 APPLICATIONS

Instance-wise risk minimization. In machine learning, we generally have access to a training set consisting of m examples. The goal of empirical risk minimization (ERM) (Vapnik, 1995) is simply to minimize the average loss over the whole training set. More generally, instance-wise risk minimization (IWRM) considers the loss associated with each training example as a distinct objective. Formally, if $\mathbf{x} \in \mathbb{R}^n$ are the parameters of the model and $f_i(\mathbf{x})$ is the loss associated to the i th example, the respective objective functions of ERM and IWRM are:

$$\text{(Empirical risk)} \quad \bar{f}(\mathbf{x}) = \frac{1}{m} \sum_{i \in [m]} f_i(\mathbf{x}) \quad (7)$$

$$\text{(Instance-wise risk)} \quad \mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}) \ f_2(\mathbf{x}) \ \cdots \ f_m(\mathbf{x})]^\top \quad (8)$$

Naively using GD for ERM is inefficient in most practical cases, so a prevalent alternative is to use SGD or one of its variants. Similarly, using JD for IWRM is typically intractable. Indeed, it would require computing a Jacobian matrix with one row per training example at each iteration. In contrast, we can use the Jacobian of a random batch of training example losses. Since it consists of a subset of the rows of the full Jacobian, this approach is a form of stochastic sub-Jacobian descent, as introduced in Section 2.1. IWRM can also be extended to cases where each f_i is a vector-valued function. The objective would then be the concatenation of the losses of all examples.

Multi-task learning. In multi-task learning, a single model is trained to perform several related tasks simultaneously, leveraging shared representations to improve overall performance (Ruder, 2017). At its core, multi-task learning is a multi-objective optimization problem (Sener & Koltun,

²This condition is a generalization to the case $m \geq 1$ of the existence of a minimizer $\mathbf{x}^* \in \mathbb{R}^n$ when $m = 1$.

270 2018), making it a straightforward application for Jacobian descent. Yet, the conflict between
 271 tasks is often too limited to justify the overhead of computing all task-specific gradients, i.e. the
 272 whole Jacobian (Kurin et al., 2022; Xin et al., 2022). In such cases, a practical approach is to
 273 minimize some linear scalarization of the objectives using an SGD-based method. Nevertheless,
 274 we believe that a setting with inherent conflict between tasks naturally prescribes Jacobian descent
 275 with a non-conflicting aggregator. We analyze several related works applied to multi-task learning in
 276 Section 4.

277
 278 **Adversarial training.** In adversarial domain adaptation, the feature extractor of a model is trained
 279 with two conflicting objectives: The features should be helpful for the main task and should be
 280 unable to discriminate the domain of the input (Ganin et al., 2016). Likewise, in adversarial fairness,
 281 the feature extractor is trained to both minimize the predictability of sensitive attributes, such as
 282 race or gender, and maximize the performance on the main task (Adel et al., 2019). Combining
 283 the corresponding gradients with a non-conflicting aggregator could enhance the optimization of
 284 such methods. We believe that the training of generative adversarial networks (Goodfellow et al.,
 285 2014) could be similarly formulated as a multi-objective optimization problem. The generator and
 286 discriminator could then be jointly optimized with JD.

287
 288 **Momentum-based optimization.** In gradient-based single-objective optimization, several methods
 289 use some form of gradient momentum to improve their convergence speed (Polyak, 1964). Essentially,
 290 their updates consider an exponential moving average of past gradients rather than just the last one. An
 291 appealing idea is to modify those algorithms to make them combine the gradient and the momentum
 292 with some aggregator, such as \mathcal{A}_{UPGrad} , instead of summing them. This would apply to many popular
 293 optimizers, like SGD with Nesterov momentum (Nesterov, 1983), Adam (Kingma & Ba, 2014),
 294 AdamW (Loshchilov & Hutter, 2019) and NAdam (Dozat, 2016).

295
 296 **Distributed optimization.** In a distributed data-parallel setting with multiple machines or multiple
 297 GPUs, model updates are computed in parallel. This can be viewed as multi-objective optimization
 298 with one objective per data share. Rather than the typical averaging, a specialized aggregator, such as
 299 \mathcal{A}_{UPGrad} , could thus combine the model updates. This consideration can even be extended to federated
 300 learning, in which multiple entities participate in the training of a common model from their own
 301 private data by sharing model updates (Kairouz et al., 2021). In this setting, as security is one of the
 302 main challenges, the non-conflicting property of the aggregator could be key.

303 4 EXISTING AGGREGATORS

304
 305 In the context of multi-task learning, several works have proposed iterative optimization algorithms
 306 based on the combination of task-specific gradients (Sener & Koltun, 2018; Yu et al., 2020; Liu
 307 et al., 2021b;a; Lin et al., 2021; Navon et al., 2022; Senushkin et al., 2023). These methods can be
 308 formulated as variants of JD parameterized by different aggregators. More specifically, since the
 309 gradients are stochastically estimated from batches of data, these are cases of what we call SEJD.
 310 In the following, we briefly present the most prominent aggregators and summarize their properties
 311 in Table 1. As a baseline, we also consider \mathcal{A}_{Mean} , which simply averages the rows of the Jacobian.
 312 Their formal definitions are provided in Appendix B. Some of them are also illustrated in Figure 1b.

313 \mathcal{A}_{RGW} aggregates the matrix using a random vector of weights (Lin et al., 2021). \mathcal{A}_{MGDA} gives the
 314 aggregation that maximizes the smallest improvement (Désidéri, 2012; Sener & Koltun, 2018; Fliege
 315 & Svaiter, 2000). \mathcal{A}_{CAGrad} maximizes the smallest improvement in a ball around the average gradient
 316 whose radius is parameterized by $c \in [0, 1[$ (Liu et al., 2021a). \mathcal{A}_{PCGrad} projects each gradient onto
 317 the orthogonal hyperplane of other gradients in case of conflict, iteratively and in a random order (Yu
 318 et al., 2020). It is, however, only non-conflicting when $m \leq 2$, in which case $\mathcal{A}_{PCGrad} = m \cdot \mathcal{A}_{UPGrad}$.
 319 IMTL-G is a method to balance some gradients with impartiality (Liu et al., 2021b). It is only defined
 320 for linearly independent gradients, but we generalize it as a formal aggregator, denoted \mathcal{A}_{IMTL-G} ,
 321 in Appendix B.6. Aligned-MTL orthonormalizes the Jacobian and weights its rows according to
 322 some preferences (Senushkin et al., 2023). We denote by $\mathcal{A}_{Aligned-MTL}$ this method with uniform
 323 preferences. $\mathcal{A}_{Nash-MTL}$ aggregates Jacobians by finding the Nash equilibrium between task-specific
 gradients (Navon et al., 2022). Lastly, the GradDrop layer (Chen et al., 2020) defines a custom
 backward pass that combines gradients with respect to some internal activation. The corresponding

aggregator, denoted $\mathcal{A}_{\text{GradDrop}}$, randomly drops out some gradient coordinates based on their sign and sums the remaining ones.

In the context of continual learning, to limit forgetting, an idea is to project the gradient onto the dual cone of gradients computed with past examples (Lopez-Paz & Ranzato, 2017). This idea can be translated into an aggregator that projects the mean gradient onto the dual cone of the rows of the Jacobian. We name this $\mathcal{A}_{\text{DualProj}}$.

Several other works consider the gradients to be noisy when making their theoretical analysis (Liu & Vicente, 2021; Zhou et al., 2022; Fernando et al., 2022; Chen et al., 2024; Xiao et al., 2024). Their solutions for combining gradients are typically stateful. Although this could enhance practical convergence rates, we have restricted our focus to the analysis of stateless aggregators. Exploring and analyzing a generalized Jacobian descent algorithm, that would preserve some state over the iterations, is a promising future direction.

In the federated learning setting, several aggregators have been proposed to combine the model updates while being robust to adversaries (Blanchard et al., 2017; Guerraoui et al., 2018; Chen et al., 2017; Yin et al., 2018). We do not study them here as they mainly focus on security aspects.

Table 1: Properties satisfied for any number of objectives. Proofs are provided in Appendix B.

Ref.	Aggregator	Non-conflicting	Linear under scaling	Weighted
—	$\mathcal{A}_{\text{Mean}}$	✗	✓	✓
Désidéri (2012)	$\mathcal{A}_{\text{MGDA}}$	✓	✗	✓
Lopez-Paz & Ranzato (2017)	$\mathcal{A}_{\text{DualProj}}$	✓	✗	✓
Yu et al. (2020)	$\mathcal{A}_{\text{PCGrad}}$	✗	✓	✓
Chen et al. (2020)	$\mathcal{A}_{\text{GradDrop}}$	✗	✗	✗
Liu et al. (2021b)	$\mathcal{A}_{\text{IMTL-G}}$	✗	✗	✓
Liu et al. (2021a)	$\mathcal{A}_{\text{CAGrad}}$	✗	✗	✓
Lin et al. (2021)	\mathcal{A}_{RGW}	✗	✓	✓
Navon et al. (2022)	$\mathcal{A}_{\text{Nash-MTL}}$	✓	✗	✓
Senushkin et al. (2023)	$\mathcal{A}_{\text{Aligned-MTL}}$	✗	✗	✓
(ours)	$\mathcal{A}_{\text{UPGrad}}$	✓	✓	✓

5 EXPERIMENTS

In the following, we present empirical results for instance-wise risk minimization on some simple image classification datasets. IWRM is performed by stochastic sub-Jacobian descent, as described in Section 3. A key consideration is that when the aggregator is $\mathcal{A}_{\text{Mean}}$, this approach becomes equivalent to empirical risk minimization with SGD. It is thus used as a baseline for comparison.

We train convolutional neural networks on subsets of SVHN (Netzer et al., 2011), CIFAR-10 (Krizhevsky et al., 2009), EuroSAT (Helber et al., 2019), MNIST (LeCun et al., 1998), Fashion-MNIST (Xiao et al., 2017) and Kuzushiji-MNIST (Clanuwat et al., 2018). To make the comparisons as fair as possible, we have tuned the learning rate very precisely for each aggregator, as explained in detail in Appendix C.1. We have also run the same experiments several times independently to gain confidence in our results. Since this leads to a total of 43776 training runs across all of our experiments, we have limited the size of each training dataset to 1024 images, greatly reducing computational costs. Note that this is strictly an optimization problem: we are not studying the generalization of the model, which would be captured by some performance metric on a test set. Other experimental settings, such as the network architectures and the total computational budget used to run our experiments, are given in Appendix C. Figure 2 reports the main results on SVHN and CIFAR-10, two of the datasets exhibiting the most substantial performance gap. Results on the other datasets and aggregators are reported in Appendix D.1. They also demonstrate a significant performance gap.

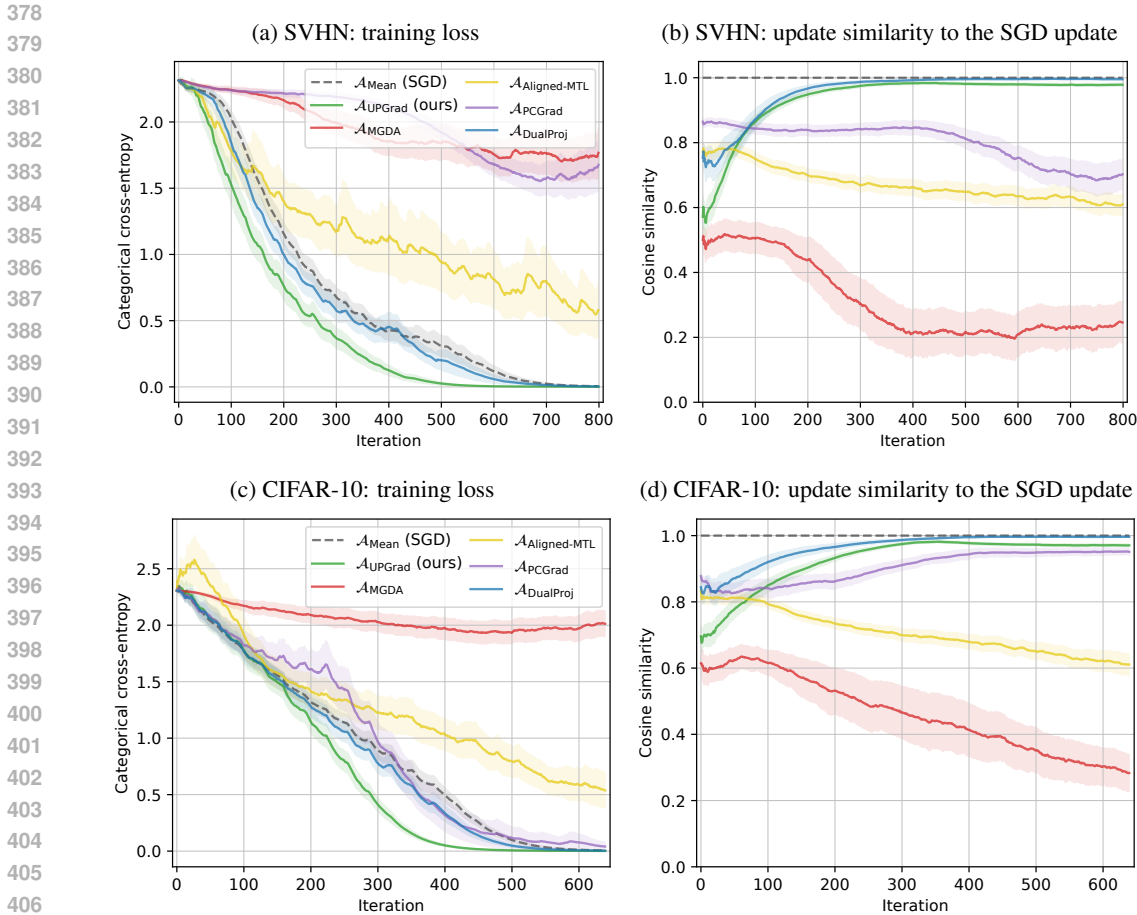


Figure 2: Optimization metrics obtained with IWRM with 1024 training examples and a batch size of 32, averaged over 8 independent runs. The shaded area around each curve shows the estimated standard error of the mean over the 8 runs. Curves are smoothed for readability. Best viewed in color.

Here, we compare the aggregators in terms of their average loss over the training set: the goal of ERM. For this reason, it is rather surprising that $\mathcal{A}_{\text{Mean}}$, which directly optimizes this objective, exhibits a slower convergence rate than some other aggregators. In particular, $\mathcal{A}_{\text{UPGrad}}$, and to a lesser extent $\mathcal{A}_{\text{DualProj}}$, provide improvements on all datasets.

Figures 2b and 2d show the similarity between the update of each aggregator and the update given by $\mathcal{A}_{\text{Mean}}$. For $\mathcal{A}_{\text{UPGrad}}$, a low similarity indicates that there are some conflicting gradients with imbalanced norms (a setting illustrated in Figure 1). Our interpretation is thus that $\mathcal{A}_{\text{UPGrad}}$ prevents gradients of hard examples from being dominated by those of easier examples early into the training. Since fitting those is more complex and time-consuming, it is beneficial to consider them earlier. We believe the similarity increases later on because the gradients become more balanced. This further suggests a greater stability of $\mathcal{A}_{\text{UPGrad}}$ compared to $\mathcal{A}_{\text{Mean}}$, which may allow it to perform effectively at a higher learning rate and, consequently, accelerate its convergence.

The sub-optimal performance of $\mathcal{A}_{\text{MGDA}}$ in this setting can be attributed to its sensitivity to small gradients. If any row of the Jacobian approaches zero, the aggregation by $\mathcal{A}_{\text{MGDA}}$ will also approach zero. This observation illustrates the discrepancy between stationarity and optimality, as discussed in Section 2.4. A notable advantage of linearity under scaling is to explicitly prevent this from happening.

Overall, these experiments demonstrate a high potential for the IWRM paradigm and confirm the relevance of JD, and more specifically of SSJD, as multi-objective optimization algorithms. Besides, the superiority of $\mathcal{A}_{\text{UPGrad}}$ in such a simple setting supports our theoretical results.

While increasing the batch size in SGD reduces variance, the effect of doing so in SSJD combined with $\mathcal{A}_{\text{UPGrad}}$ is non-trivial, as it also tightens the dual cone. Additional results obtained when varying the batch size or updating the parameters with the Adam optimizer are available in Appendices D.2 and D.3, respectively.

While an iteration of SSJD is more expensive than an iteration of SGD, its runtime is influenced by several factors, including the choice of aggregator, the parallelization capabilities of the hardware used for Jacobian computation, and the implementation. Appendix E provides memory usage and computation time considerations for our methods. Additionally, we propose a path towards a more efficient implementation in the next section.

6 GRAMIAN-BASED JACOBIAN DESCENT

When the number of objectives is dominated by the number of parameters of the model, the main overhead of JD comes from the usage of a Jacobian matrix rather than a single gradient. In the following, we motivate an alternative implementation of JD that only uses the inner products between each pair of gradients.

For any $J \in \mathbb{R}^{m \times n}$, the matrix $G = JJ^\top$ is called the *Gramian* of J and is positive semi-definite. Let $\mathcal{M}_m \subseteq \mathbb{R}^{m \times m}$ be the set of positive semi-definite matrices. The Gramian of the Jacobian, denoted $\mathcal{G}\mathbf{f}(\mathbf{x}) = \mathcal{J}\mathbf{f}(\mathbf{x}) \cdot \mathcal{J}\mathbf{f}(\mathbf{x})^\top \in \mathcal{M}_m$, captures the relations – including conflicts – between all pairs of gradients. Whenever \mathcal{A} is a weighted aggregator, the update of JD is $\mathbf{y} = -\eta \mathcal{J}\mathbf{f}(\mathbf{x})^\top \mathbf{w}$ for some vector of weights $\mathbf{w} \in \mathbb{R}^m$. Substituting this into the Taylor approximation of (2) gives

$$\mathbf{f}(\mathbf{x} + \mathbf{y}) = \mathbf{f}(\mathbf{x}) - \eta \mathcal{G}\mathbf{f}(\mathbf{x}) \cdot \mathbf{w} + o\left(\eta \sqrt{\mathbf{w}^\top \cdot \mathcal{G}\mathbf{f}(\mathbf{x}) \cdot \mathbf{w}}\right). \quad (9)$$

This expression only depends on the Jacobian through its Gramian. It is thus sensible to focus on aggregators whose weights are only a function of the Gramian. Denoting this function as $\mathcal{W} : \mathcal{M}_m \rightarrow \mathbb{R}^m$, those aggregators satisfy $\mathcal{A}(J) = J^\top \cdot \mathcal{W}(G)$. Remarkably, all weighted aggregators of Table 1 can be expressed in this form. In the case of $\mathcal{A}_{\text{UPGrad}}$, this is clearly demonstrated in Proposition 1, which shows that the weights depend on G . For such aggregators, substitution and linearity of differentiation³ then yield

$$\mathcal{A}(\mathcal{J}\mathbf{f}(\mathbf{x})) = \nabla\left(\mathcal{W}(\mathcal{G}\mathbf{f}(\mathbf{x}))^\top \cdot \mathbf{f}\right)(\mathbf{x}). \quad (10)$$

After computing $\mathcal{W}(\mathcal{G}\mathbf{f}(\mathbf{x}))$, a step of JD would thus only require the backpropagation of a scalar function. The computational cost of applying \mathcal{W} depends on the aggregator and is often dominated by the cost of computing the Gramian.

We now outline a method to compute the Gramian of the Jacobian without ever having to store the full Jacobian in memory. Similarly to the backpropagation algorithm, we can leverage the chain rule. Let $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and $\mathbf{f} : \mathbb{R}^k \rightarrow \mathbb{R}^m$, then for any $\mathbf{x} \in \mathbb{R}^n$, the chain rule for Gramians is

$$\mathcal{G}(\mathbf{f} \circ \mathbf{g})(\mathbf{x}) = \mathcal{J}\mathbf{f}(\mathbf{g}(\mathbf{x})) \cdot \mathcal{G}\mathbf{g}(\mathbf{x}) \cdot \mathcal{J}\mathbf{f}(\mathbf{g}(\mathbf{x}))^\top. \quad (11)$$

Moreover, when the function has multiple inputs, the Gramian can be computed as a sum of individual Gramians. Let $\mathbf{f} : \mathbb{R}^{n_1 + \dots + n_k} \rightarrow \mathbb{R}^m$ and $\mathbf{x} = [\mathbf{x}_1^\top \dots \mathbf{x}_k^\top]^\top$. We can write $\mathcal{J}\mathbf{f}(\mathbf{x})$ as the concatenation of Jacobians $[\mathcal{J}_{\mathbf{x}_1}\mathbf{f}(\mathbf{x}) \dots \mathcal{J}_{\mathbf{x}_k}\mathbf{f}(\mathbf{x})]$, where $\mathcal{J}_{\mathbf{x}_i}\mathbf{f}(\mathbf{x})$ is the Jacobian of \mathbf{f} with respect to \mathbf{x}_i evaluated at \mathbf{x} . For any $i \in [k]$, let $\mathcal{G}_{\mathbf{x}_i}\mathbf{f}(\mathbf{x}) = \mathcal{J}_{\mathbf{x}_i}\mathbf{f}(\mathbf{x}) \cdot \mathcal{J}_{\mathbf{x}_i}\mathbf{f}(\mathbf{x})^\top$. Then

$$\mathcal{G}\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_k) = \sum_{i \in [k]} \mathcal{G}_{\mathbf{x}_i}\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_k). \quad (12)$$

When a function is made of compositions and concatenations of elementary functions, the Gramian of the Jacobian can thus be expressed with sums and products of partial Jacobians.

We now provide an example algorithm to compute the Gramian of a sequence of layers. For $0 \leq i < k$, let $\mathbf{f}_i : \mathbb{R}^{n_i} \times \mathbb{R}^{\ell_i} \rightarrow \mathbb{R}^{n_{i+1}}$ be a layer parameterized by $\mathbf{p}_i \in \mathbb{R}^{\ell_i}$. Given $\mathbf{x}_0 \in \mathbb{R}^{n_0}$, for $0 \leq i < k$,

³For any $\mathbf{x} \in \mathbb{R}^n$ and any $\mathbf{w} \in \mathbb{R}^m$, $\mathcal{J}\mathbf{f}(\mathbf{x})^\top \mathbf{w} = \nabla(\mathbf{w}^\top \mathbf{f})(\mathbf{x})$

the activations are recursively defined as $\mathbf{x}_{i+1} = \mathbf{f}_i(\mathbf{x}_i, \mathbf{p}_i)$. Algorithm 3 illustrates how (11) and (12) can be combined to compute the Gramian of the network with respect to its parameters.

Algorithm 3: *Gramian reverse accumulation* for a sequence of layers

```

490  $J_x \leftarrow I$  # Identity matrix of size  $n_k \times n_k$ 
491  $G \leftarrow 0$  # Zero matrix of size  $n_k \times n_k$ 
492 for  $i \leftarrow k - 1$  to 0 do
493    $J_p \leftarrow \mathcal{J}_{\mathbf{p}_i} \mathbf{f}_i(\mathbf{x}_i, \mathbf{p}_i) \cdot J_x$  # Jacobian of  $\mathbf{x}_k$  w.r.t.  $\mathbf{p}_i$ 
494    $J_x \leftarrow \mathcal{J}_{\mathbf{x}_i} \mathbf{f}_i(\mathbf{x}_i, \mathbf{p}_i) \cdot J_x$  # Jacobian of  $\mathbf{x}_k$  w.r.t.  $\mathbf{x}_i$ 
495    $G \leftarrow G + J_p J_p^\top$ 

```

Output: G

Generalizing Algorithm 3 to any computational graph and implementing it efficiently remains an open challenge extending beyond the scope of this work.

7 CONCLUSION

In this paper, we introduced Jacobian descent (JD), a multi-objective optimization algorithm defined by some aggregator that maps the Jacobian to an update direction. We identified desirable properties for aggregators and proposed $\mathcal{A}_{\text{UPGrad}}$, addressing the limitations of existing methods while providing stronger convergence guarantees. We also highlighted potential applications of JD and proposed IWRM, a novel learning paradigm considering the loss of each training example as a distinct objective. Given its promising empirical results, we believe this paradigm deserves further attention. Additionally, we see potential for $\mathcal{A}_{\text{UPGrad}}$ beyond JD, as a linear algebra tool for combining conflicting vectors in broader contexts. As speed is the primary limitation of JD, we have outlined an algorithm for efficiently computing the Gramian of the Jacobian, which could unlock JD’s full potential. We hope this work serves as a foundation for future research in multi-objective optimization and encourages a broader adoption of these methods.

Limitations and future directions. Our experimentation has some limitations. First, we only evaluate JD on IWRM, a setting with moderately conflicting objectives. It would be essential to develop proper benchmarks to compare aggregators on a wide variety of problems. Ideally, such problems should involve substantially conflicting objectives, e.g. multi-task learning with inherently competing or even adversarial tasks. Then, we have limited our scope to the comparison of optimization speeds, disregarding generalization. While this simplifies the experiments and makes the comparison rigorous, optimization and generalization are sometimes intertwined. We thus believe that future works should focus on both aspects.

REFERENCES

- Tameem Adel, Isabel Valera, Zoubin Ghahramani, and Adrian Weller. One-network adversarial fairness. In *AAAI Conference on Artificial Intelligence*, 2019.
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Advances in Neural Information Processing Systems*, 2017.
- Stephen P Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge university press, 2004.
- Jürgen Branke. *Multiobjective Optimization: Interactive and Evolutionary Approaches*. Springer Science & Business Media, 2008.
- Stéphane Caron, Daniel Arnström, Suraj Bonagiri, Antoine Dechaume, Nikolai Flowers, Adam Heins, Takuma Ishikawa, Dustin Kenefake, Giacomo Mazzamuto, Donato Meoli, Brendan O’Donoghue, Adam A. Oppenheimer, Abhishek Pandala, Juan José Quiroz Omaña, Nikitas Rontsis, Paarth Shah, Samuel St-Jean, Nicola Vitucci, Soeren Wolfers, Fengyu Yang, @bdelhaisse, @MeindertHH, @rimaddo, @urob, and @shaoanlu. qpsolvers: Quadratic Programming Solvers in Python, 2024. URL <https://github.com/qpsolvers/qpsolvers>.

- 540 Lisha Chen, Heshan Fernando, Yiming Ying, and Tianyi Chen. Three-way trade-off in multi-objective
541 learning: Optimization, generalization and conflict-avoidance. In *Advances in Neural Information*
542 *Processing Systems*, 2024.
- 543 Yudong Chen, Lili Su, and Jiaming Xu. Distributed statistical machine learning in adversarial
544 settings: Byzantine gradient descent. In *Proceedings of the ACM on Measurement and Analysis of*
545 *Computing Systems*, 2017.
- 547 Zhao Chen, Jiquan Ngiam, Yanping Huang, Thang Luong, Henrik Kretschmar, Yuning Chai, and
548 Dragomir Anguelov. Just pick a sign: Optimizing deep multitask models with gradient sign
549 dropout. In *Advances in Neural Information Processing Systems*, 2020.
- 550 Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David
551 Ha. Deep learning for classical japanese literature. In *NeurIPS Workshop on Machine Learning*
552 *for Creativity and Design*, 2018.
- 554 Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network
555 learning by exponential linear units (ELUs). *arXiv preprint arXiv:1511.07289*, 2015.
- 556 Kalyanmoy Deb, Karthik Sindhya, and Jussi Hakanen. Multi-objective optimization. In *Decision*
557 *sciences*. CRC Press, 2016.
- 559 Timothy Dozat. Incorporating Nesterov momentum into Adam. In *International Conference on*
560 *Learning Representations Workshop*, 2016.
- 561 Jean-Antoine Désidéri. Multiple-gradient descent algorithm (MGDA) for multiobjective optimization.
562 *Comptes Rendus Mathématique*, 2012.
- 563 Matthias Ehrgott. *Multicriteria Optimization*. Springer Science & Business Media, 2005.
- 564 Heshan Devaka Fernando, Han Shen, Miao Liu, Subhajit Chaudhury, Keerthiram Murugesan, and
565 Tianyi Chen. Mitigating gradient bias in multi-objective learning: A provably convergent approach.
566 In *International Conference on Learning Representations*, 2022.
- 569 Jörg Fliege, A Ismael F Vaz, and Luís Nunes Vicente. Complexity of gradient descent for multiobjec-
570 tive optimization. *Optimization Methods and Software*, 2019.
- 571 Jörg Fliege and Benar Fux Svaiter. Steepest descent methods for multicriteria optimization. *Mathe-*
572 *matical Methods of Operations Research*, 2000.
- 574 Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
575 Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks.
576 *Journal of machine learning research*, 2016.
- 577 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
578 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural*
579 *Information Processing Systems*, 2014.
- 580 Rachid Guerraoui, Sébastien Rouault, et al. The hidden vulnerability of distributed learning in
581 byzantium. In *International Conference on Machine Learning*, 2018.
- 583 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. EuroSAT: A novel dataset
584 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected*
585 *Topics in Applied Earth Observations and Remote Sensing*, 2019.
- 586 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
587 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-
588 vances and open problems in federated learning. *Foundations and Trends® in Machine Learning*,
589 2021.
- 591 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
592 *arXiv:1412.6980*, 2014.
- 593 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images, 2009.

- 594 Vitaly Kurin, Alessandro De Palma, Ilya Kostrikov, Shimon Whiteson, and Pawan K Mudigonda. In
595 defense of the unitary scalarization for deep multi-task learning. In *Advances in Neural Information*
596 *Processing Systems*, 2022.
- 597 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
598 document recognition. In *Proceedings of the IEEE*, 1998.
- 600 Baijiong Lin, Feiyang Ye, Yu Zhang, and Ivor W Tsang. Reasonable effectiveness of random
601 weighting: A litmus test for multi-task learning. *arXiv preprint arXiv:2111.10603*, 2021.
- 602 Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for
603 multi-task learning. In *Advances in Neural Information Processing Systems*, 2021a.
- 604 Liyang Liu, Yi Li, Zhanghui Kuang, Jing-Hao Xue, Yimin Chen, Wenming Yang, Qingmin Liao, and
605 Wayne Zhang. Towards impartial multi-task learning. In *International Conference on Learning*
606 *Representations*, 2021b.
- 607 Suyun Liu and Luis Nunes Vicente. The stochastic multi-gradient algorithm for multi-objective
608 optimization and its application to supervised machine learning. *Annals of Operations Research*,
609 2021.
- 610 David Lopez-Paz and Marc’ Aurelio Ranzato. Gradient episodic memory for continual learning. In
611 *Advances in Neural Information Processing Systems*, 2017.
- 612 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
613 *ence on Learning Representations*, 2019.
- 614 Quentin Mercier, Fabrice Poirion, and Jean-Antoine Désidéri. A stochastic multiple gradient descent
615 algorithm. *European Journal of Operational Research*, 2018.
- 616 Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik, and
617 Ethan Fetaya. Multi-task learning as a bargaining game. In *International Conference on Machine*
618 *Learning*, 2022.
- 619 Yuri Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^{**2})$.
620 *Proceedings of the USSR Academy of Sciences*, 1983.
- 621 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.
622 Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep*
623 *learning and unsupervised feature learning*, 2011.
- 624 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
625 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style,
626 high-performance deep learning library. In *Advances in Neural Information Processing Systems*,
627 2019.
- 628 Fabrice Poirion, Quentin Mercier, and Jean-Antoine Désidéri. Descent algorithm for nonsmooth
629 stochastic multiobjective optimization. *Computational Optimization and Applications*, 2017.
- 630 Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR*
631 *computational mathematics and mathematical physics*, 1964.
- 632 Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint*
633 *arXiv:1706.05098*, 2017.
- 634 Yoshikazu Sawaragi, Hirotaka Nakayama, and Tetsuzo Tanino. *Theory of Multiobjective Optimization*.
635 Elsevier, 1985.
- 636 Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in*
637 *Neural Information Processing Systems*, 2018.
- 638 Dmitry Senushkin, Nikolay Patakin, Arseny Kuznetsov, and Anton Konushin. Independent component
639 alignment for multi-task learning. In *IEEE/CVF Conference on Computer Vision and Pattern*
640 *Recognition*, 2023.

648 Vladimir Naumovich Vapnik. *The Nature of Statistical learning theory*. Wiley New York, 1995.
649

650 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
651 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

652 Peiyao Xiao, Hao Ban, and Kaiyi Ji. Direction-oriented multi-objective learning: Simple and provable
653 stochastic algorithms. In *Advances in Neural Information Processing Systems*, 2024.
654

655 Derrick Xin, Behrooz Ghorbani, Justin Gilmer, Ankush Garg, and Orhan Firat. Do current multi-task
656 optimization methods in deep learning even help? In *Advances in Neural Information Processing*
657 *Systems*, 2022.

658 Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed
659 learning: Towards optimal statistical rates. In *International Conference on Machine Learning*,
660 2018.

661 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
662 Gradient surgery for multi-task learning. In *Advances in Neural Information Processing Systems*,
663 2020.
664

665 Shiji Zhou, Wenpeng Zhang, Jiyang Jiang, Wenliang Zhong, Jinjie Gu, and Wenwu Zhu. On the
666 convergence of stochastic multi-objective gradient manipulation and beyond. In *Advances in*
667 *Neural Information Processing Systems*, 2022.
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A PROOFS

703 A.1 SUPPLEMENTARY THEORETICAL RESULTS

704 Recall that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is \leq -convex if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and any $\lambda \in [0, 1]$,

$$705 f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

706 **Lemma 1.** If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a continuously differentiable \leq -convex function, then for any pair of
707 vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\mathcal{J}f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x})$.

708 *Proof.*

$$709 \begin{aligned} \mathcal{J}f(\mathbf{x})(\mathbf{y} - \mathbf{x}) &= \lim_{\lambda \rightarrow 0^+} \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda} && \text{(differentiation)} \\ &\leq \lim_{\lambda \rightarrow 0^+} \frac{f(\mathbf{x}) + \lambda(f(\mathbf{y}) - f(\mathbf{x})) - f(\mathbf{x})}{\lambda} && (\leq\text{-convexity}) \\ &= f(\mathbf{y}) - f(\mathbf{x}), \end{aligned}$$

710 which concludes the proof. \square

711 **Lemma 2.** Let $J \in \mathbb{R}^{m \times n}$, let $\mathbf{u} \in \mathbb{R}^m$ and let $\mathbf{x} \in \mathbb{R}^n$, then

$$712 \mathbf{u}^\top J \mathbf{x} \leq \|\mathbf{u}\| \cdot \|J\|_{\text{F}} \cdot \|\mathbf{x}\|$$

713 *Proof.* Let J_i be the i th row of J , then

$$714 \begin{aligned} (\mathbf{u}^\top J \mathbf{x})^2 &\leq \|\mathbf{u}\|^2 \cdot \|J \mathbf{x}\|^2 && \text{(Cauchy-Schwartz)} \\ &= \|\mathbf{u}\|^2 \cdot \sum_{i \in [m]} (J_i^\top \mathbf{x})^2 && \text{inequality} \\ &\leq \|\mathbf{u}\|^2 \cdot \sum_{i \in [m]} \|J_i\|^2 \cdot \|\mathbf{x}\|^2 && \text{(Cauchy-Schwartz)} \\ &= \|\mathbf{u}\|^2 \cdot \|J\|_{\text{F}}^2 \cdot \|\mathbf{x}\|^2, && \text{inequality} \end{aligned}$$

715 which concludes the proof. \square

716 Recall that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is β -smooth if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$717 \|\mathcal{J}f(\mathbf{x}) - \mathcal{J}f(\mathbf{y})\|_{\text{F}} \leq \beta \|\mathbf{x} - \mathbf{y}\| \quad (13)$$

718 **Lemma 3.** Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be β -smooth, then for any $\mathbf{w} \in \mathbb{R}^m$ and any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$719 \mathbf{w}^\top (f(\mathbf{x}) - f(\mathbf{y}) - \mathcal{J}f(\mathbf{y})(\mathbf{x} - \mathbf{y})) \leq \frac{\beta}{2} \|\mathbf{w}\| \cdot \|\mathbf{x} - \mathbf{y}\|^2 \quad (14)$$

720 *Proof.*

$$721 \begin{aligned} &\mathbf{w}^\top (f(\mathbf{x}) - f(\mathbf{y}) - \mathcal{J}f(\mathbf{y})(\mathbf{x} - \mathbf{y})) \\ &= \mathbf{w}^\top \left(\int_0^1 \mathcal{J}f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))(\mathbf{x} - \mathbf{y}) dt - \mathcal{J}f(\mathbf{y})(\mathbf{x} - \mathbf{y}) \right) && \text{(fundamental theorem of calculus)} \\ &= \int_0^1 \mathbf{w}^\top \left(\mathcal{J}f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \mathcal{J}f(\mathbf{y}) \right) (\mathbf{x} - \mathbf{y}) dt \\ &\leq \int_0^1 \|\mathbf{w}\| \cdot \left\| \mathcal{J}f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \mathcal{J}f(\mathbf{y}) \right\|_{\text{F}} \cdot \|\mathbf{x} - \mathbf{y}\| dt && \text{(Lemma 2)} \\ &\leq \int_0^1 \|\mathbf{w}\| \cdot \beta t \cdot \|\mathbf{x} - \mathbf{y}\|^2 dt && (\beta\text{-smoothness 13}) \\ &= \frac{\beta}{2} \|\mathbf{w}\| \cdot \|\mathbf{x} - \mathbf{y}\|^2, \end{aligned}$$

722 which concludes the proof. \square

A.2 PROPOSITION 1

Proposition 1. Let $J \in \mathbb{R}^{m \times n}$. For any $\mathbf{u} \in \mathbb{R}^m$, $\pi_J(J^\top \mathbf{u}) = J^\top \mathbf{w}$ with

$$\mathbf{w} \in \arg \min_{\mathbf{v} \in \mathbb{R}^m: \mathbf{u} \leq \mathbf{v}} \mathbf{v}^\top J J^\top \mathbf{v}. \quad (5)$$

Proof. This is a direct consequence of Lemma 4. \square

Lemma 4. Let $J \in \mathbb{R}^{m \times n}$, $G = J J^\top$, $\mathbf{u} \in \mathbb{R}^m$. For any $\mathbf{w} \in \mathbb{R}^m$ satisfying

$$\begin{cases} \mathbf{u} \leq \mathbf{w} & (15a) \\ \mathbf{0} \leq G\mathbf{w} & (15b) \\ \mathbf{u}^\top G\mathbf{w} = \mathbf{w}^\top G\mathbf{w} & (15c) \end{cases}$$

we have $\pi_J(J^\top \mathbf{u}) = J^\top \mathbf{w}$. Such a \mathbf{w} is the solution to

$$\mathbf{w} \in \arg \min_{\mathbf{u} \leq \mathbf{v}} \mathbf{v}^\top G\mathbf{v}.$$

Proof. The projection

$$\pi_J(J^\top \mathbf{u}) = \arg \min_{\substack{\mathbf{x} \in \mathbb{R}^n: \\ \mathbf{0} \leq J\mathbf{x}}} \frac{1}{2} \|\mathbf{x} - J^\top \mathbf{u}\|^2$$

is a convex program. Consequently, the KKT conditions are both necessary and sufficient. The Lagrangian is given by $\mathcal{L}(\mathbf{x}, \mathbf{v}) = \frac{1}{2} \|\mathbf{x} - J^\top \mathbf{u}\|^2 - \mathbf{v}^\top J\mathbf{x}$. The KKT conditions are then given by

$$\begin{aligned} & \begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{v}) = \mathbf{0} \\ \mathbf{0} \leq \mathbf{v} \\ \mathbf{0} \leq J\mathbf{x} \\ 0 = \mathbf{v}^\top J\mathbf{x} \end{cases} \\ \Leftrightarrow & \begin{cases} \mathbf{x} = J^\top(\mathbf{u} + \mathbf{v}) \\ \mathbf{0} \leq \mathbf{v} \\ \mathbf{0} \leq G(\mathbf{u} + \mathbf{v}) \\ 0 = \mathbf{v}^\top G(\mathbf{u} + \mathbf{v}) \end{cases} \\ \Leftrightarrow & \begin{cases} \mathbf{x} = J^\top(\mathbf{u} + \mathbf{v}) \\ \mathbf{u} \leq \mathbf{u} + \mathbf{v} \\ \mathbf{0} \leq G(\mathbf{u} + \mathbf{v}) \\ \mathbf{u}^\top G(\mathbf{u} + \mathbf{v}) = (\mathbf{u} + \mathbf{v})^\top G(\mathbf{u} + \mathbf{v}) \end{cases} \end{aligned}$$

The simple change of variable $\mathbf{w} = \mathbf{u} + \mathbf{v}$ finishes the proof of the first part.

Since $\mathbf{x} = J^\top(\mathbf{u} + \mathbf{v})$, the Wolfe dual program of $\pi_J(J^\top \mathbf{u})$ gives

$$\begin{aligned} \mathbf{w} & \in \mathbf{u} + \arg \max_{\mathbf{v} \in \mathbb{R}^m: \mathbf{0} \leq \mathbf{v}} \mathcal{L}(J^\top(\mathbf{u} + \mathbf{v}), \mathbf{v}) \\ & = \mathbf{u} + \arg \max_{\mathbf{v} \in \mathbb{R}^m: \mathbf{0} \leq \mathbf{v}} \frac{1}{2} \|J^\top \mathbf{v}\|^2 - \mathbf{v}^\top J J^\top (\mathbf{u} + \mathbf{v}) \\ & = \mathbf{u} + \arg \max_{\mathbf{v} \in \mathbb{R}^m: \mathbf{0} \leq \mathbf{v}} -\frac{1}{2} \mathbf{v}^\top G\mathbf{v} - \mathbf{v}^\top G\mathbf{u} \\ & = \mathbf{u} + \arg \min_{\mathbf{v} \in \mathbb{R}^m: \mathbf{u} \leq \mathbf{u} + \mathbf{v}} \frac{1}{2} (\mathbf{u} + \mathbf{v})^\top G(\mathbf{u} + \mathbf{v}) \\ & = \arg \min_{\mathbf{v}' \in \mathbb{R}^m: \mathbf{u} \leq \mathbf{v}'} \frac{1}{2} \mathbf{v}'^\top G\mathbf{v}', \end{aligned}$$

which concludes the proof. \square

A.3 THEOREM 1

Theorem 1. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a β -smooth and \leq -convex function. Suppose that the Pareto front $f(X^*)$ is bounded and that for any $\mathbf{x} \in \mathbb{R}^n$, there is $\mathbf{x}^* \in X^*$ satisfying $f(\mathbf{x}^*) \leq f(\mathbf{x})$. Let $\mathbf{x}_1 \in \mathbb{R}^n$, and for all $t \in \mathbb{N}$, $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathcal{A}_{\text{UPGrad}}(\mathcal{J}f(\mathbf{x}_t))$, with $\eta = \frac{1}{\beta\sqrt{m}}$. Let \mathbf{w}_t be the weights defining $\mathcal{A}_{\text{UPGrad}}(\mathcal{J}f(\mathbf{x}_t))$ as per (6), i.e. $\mathcal{A}_{\text{UPGrad}}(\mathcal{J}f(\mathbf{x}_t)) = \mathcal{J}f(\mathbf{x}_t)^\top \cdot \mathbf{w}_t$. If \mathbf{w}_t is bounded, then $f(\mathbf{x}_t)$ converges to $f(\mathbf{x}^*)$ for some $\mathbf{x}^* \in X^*$. In other words, $f(\mathbf{x}_t)$ converges to the Pareto front.

To prove the theorem we will need Lemmas 5, 6 and 7 below.

Lemma 5. Let $J \in \mathbb{R}^{m \times n}$ and $\mathbf{w} = \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i$ be the weights defining $\mathcal{A}_{\text{UPGrad}}(J)$ as per (6). Let, as usual, $G = JJ^\top$, then,

$$\mathbf{w}^\top G \mathbf{w} \leq \mathbf{1}^\top G \mathbf{w}.$$

Proof. Observe that if, for any $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$, $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top G \mathbf{v}$, then $\langle \cdot, \cdot \rangle$ is an inner product. In this Hilbert space, the Cauchy-Schwartz inequality reads as

$$\begin{aligned} (\mathbf{u}^\top G \mathbf{v})^2 &= \langle \mathbf{u}, \mathbf{v} \rangle^2 \\ &\leq \langle \mathbf{u}, \mathbf{u} \rangle \cdot \langle \mathbf{v}, \mathbf{v} \rangle \\ &= \mathbf{u}^\top G \mathbf{u} \cdot \mathbf{v}^\top G \mathbf{v}. \end{aligned}$$

Therefore

$$\begin{aligned} &\mathbf{w}^\top G \mathbf{w} \\ &= \frac{1}{m^2} \sum_{i,j} \mathbf{w}_i^\top G \mathbf{w}_j \\ &\leq \frac{1}{m^2} \sum_{i,j} \sqrt{\mathbf{w}_i^\top G \mathbf{w}_i} \cdot \sqrt{\mathbf{w}_j^\top G \mathbf{w}_j} && \left(\begin{array}{l} \text{Cauchy-Schwartz} \\ \text{inequality} \end{array} \right) \\ &= \left(\sum_i \frac{1}{m} \sqrt{\mathbf{w}_i^\top G \mathbf{w}_i} \right)^2 \\ &\leq \sum_i \frac{1}{m} \left(\sqrt{\mathbf{w}_i^\top G \mathbf{w}_i} \right)^2 && \left(\text{Jensen's inequality} \right) \\ &= \frac{1}{m} \sum_i \mathbf{w}_i^\top G \mathbf{w}_i && \left(G \text{ positive semi-definite} \right) \\ &= \frac{1}{m} \sum_i \mathbf{e}_i^\top G \mathbf{w}_i && \left(\text{Lemma 4, (15c)} \right) \\ &\leq \frac{1}{m} \sum_i \mathbf{1}^\top G \mathbf{w}_i && \left(\begin{array}{l} \text{Lemma 4, (15b)} \\ \mathbf{e}_i \leq \mathbf{1} \end{array} \right) \\ &= \mathbf{1}^\top G \mathbf{w}, \end{aligned}$$

which concludes the proof. \square

Lemma 6. Under the assumptions of Theorem 1, for any $\mathbf{w} \in \mathbb{R}^m$ and any $t \in \mathbb{N}$,

$$\mathbf{w}^\top (f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) \leq \frac{\|\mathbf{w}\|}{\beta\sqrt{m}} \left(\frac{\mathbf{1}}{2\sqrt{m}} - \frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^\top G_t \mathbf{w}_t.$$

864 *Proof.* For all $t \in \mathbb{N}$, let $J_t = \mathcal{J}f(\mathbf{x}_t)$, $G_t = J_t J_t^\top$. Then $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathcal{A}_{\text{UPGrad}}(J_t) = \mathbf{x}_t - \eta J_t^\top \mathbf{w}_t$.
 865 Therefore

$$\begin{aligned}
 & \mathbf{w}^\top (f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) \\
 & \leq -\eta \mathbf{w}^\top J_t J_t^\top \mathbf{w}_t + \frac{\beta \eta^2}{2} \|\mathbf{w}\| \cdot \|J_t^\top \mathbf{w}_t\|^2 && \text{(Lemma 3)} \\
 & = -\frac{1}{\beta \sqrt{m}} \mathbf{w}^\top G_t \mathbf{w}_t + \frac{1}{2\beta m} \|\mathbf{w}\| \cdot \mathbf{w}_t^\top G_t \mathbf{w}_t && (\eta = \frac{1}{\beta \sqrt{m}}) \\
 & \leq -\frac{1}{\beta \sqrt{m}} \mathbf{w}^\top G_t \mathbf{w}_t + \frac{1}{2\beta m} \|\mathbf{w}\| \cdot \mathbf{1}^\top G_t \mathbf{w}_t && \text{(Lemma 5)} \\
 & = \frac{\|\mathbf{w}\|}{\beta \sqrt{m}} \left(\frac{\mathbf{1}}{2\sqrt{m}} - \frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^\top G_t \mathbf{w}_t,
 \end{aligned}$$

878 which concludes the proof. \square

880 **Lemma 7.** Under the assumptions of Theorem 1, if $\mathbf{x}^* \in X^*$ satisfies $\mathbf{1}^\top f(\mathbf{x}^*) \leq \mathbf{1}^\top f(\mathbf{x}_t)$ for all
 881 $t \in \mathbb{N}$, then

$$\frac{1}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{1}{T} \left(\mathbf{1}^\top (f(\mathbf{x}_1) - f(\mathbf{x}^*)) + \frac{\beta \sqrt{m}}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \right). \quad (16)$$

888 *Proof.* We first bound, for any $t \in \mathbb{N}$, $\mathbf{1}^\top (f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t))$ as follows

$$\begin{aligned}
 & \mathbf{1}^\top (f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) \\
 & \leq -\frac{1}{2\beta \sqrt{m}} \cdot \mathbf{1}^\top G_t \mathbf{w}_t && \text{(Lemma 6 with } \mathbf{w}=\mathbf{1}\text{)} \\
 & \leq -\frac{1}{2\beta \sqrt{m}} \cdot \mathbf{w}_t^\top G_t \mathbf{w}_t. && \text{(Lemma 5)}
 \end{aligned}$$

896 Summing this over $t \in [T]$ yields

$$\begin{aligned}
 & \frac{1}{2\beta \sqrt{m}} \sum_{t \in [T]} \mathbf{w}_t^\top G_t \mathbf{w}_t \\
 & \leq \sum_{t \in [T]} \mathbf{1}^\top (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})) \\
 & = \mathbf{1}^\top (f(\mathbf{x}_1) - f(\mathbf{x}_{T+1})) && \text{(Telescoping sum)} \\
 & \leq \mathbf{1}^\top (f(\mathbf{x}_1) - f(\mathbf{x}^*)). && \left(\mathbf{1}^\top f(\mathbf{x}^*) \leq \mathbf{1}^\top f(\mathbf{x}_{T+1}) \right) \quad (17)
 \end{aligned}$$

908 Since $\mathbf{0} \leq \mathbf{w}_t$,

$$\begin{aligned}
 & \mathbf{w}_t^\top (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \\
 & \leq \mathbf{w}_t^\top J_t (\mathbf{x}_t - \mathbf{x}^*) && \text{(Lemma 1)} \\
 & = \frac{1}{\eta} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*) && (\mathbf{x}_{t+1} = \mathbf{x}_t - \eta J_t^\top \mathbf{w}_t) \\
 & = \frac{1}{2\eta} \left(\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right) && \text{(Parallelogram law)} \\
 & = \frac{1}{2\beta \sqrt{m}} \mathbf{w}_t^\top G_t \mathbf{w}_t + \frac{\beta \sqrt{m}}{2} \left(\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right). && (\eta = \frac{1}{\beta \sqrt{m}})
 \end{aligned}$$

Summing this over $t \in [T]$ yields

$$\begin{aligned}
& \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \\
& \leq \frac{1}{2\beta\sqrt{m}} \sum_{t \in [T]} \mathbf{w}_t^\top G_t \mathbf{w}_t + \frac{\beta\sqrt{m}}{2} (\|\mathbf{x}_1 - \mathbf{x}^*\|^2 - \|\mathbf{x}_{T+1} - \mathbf{x}^*\|^2) \quad (\text{Telescoping sum}) \\
& \leq \frac{1}{2\beta\sqrt{m}} \sum_{t \in [T]} \mathbf{w}_t^\top G_t \mathbf{w}_t + \frac{\beta\sqrt{m}}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \\
& \leq \mathbf{1}^\top (\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)) + \frac{\beta\sqrt{m}}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2. \quad (\text{By (17)})
\end{aligned}$$

Scaling down this inequality by T yields

$$\frac{1}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \leq \frac{1}{T} \left(\mathbf{1}^\top (\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)) + \frac{\beta\sqrt{m}}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \right),$$

which concludes the proof. \square

We are now ready to prove Theorem 1.

Proof. For all $t \in \mathbb{N}$, let $J_t = \mathcal{J}\mathbf{f}(\mathbf{x}_t)$, $G_t = J_t J_t^\top$. Then

$$\begin{aligned}
\mathbf{x}_{t+1} &= \mathbf{x}_t - \eta \mathcal{A}_{\text{UPGrad}}(J_t) \\
&= \mathbf{x}_t - \eta J_t^\top \mathbf{w}_t.
\end{aligned}$$

Substituting $\mathbf{w} = \mathbf{1}$ in the term $\frac{\mathbf{1}}{2\sqrt{m}} - \frac{\mathbf{w}}{\|\mathbf{w}\|}$ of Lemma 6 yields

$$\begin{aligned}
\frac{\mathbf{1}}{2\sqrt{m}} - \frac{\mathbf{w}}{\|\mathbf{w}\|} &= -\frac{\mathbf{1}}{2\sqrt{m}} \\
&< \mathbf{0}.
\end{aligned}$$

Therefore there exists some $\varepsilon > 0$ such that any $\mathbf{w} \in \mathbb{R}^m$ with $\|\mathbf{1} - \mathbf{w}\| < \varepsilon$ satisfies $\frac{\mathbf{1}}{2\sqrt{m}} < \frac{\mathbf{w}}{\|\mathbf{w}\|}$. Denote by $B_\varepsilon(\mathbf{1}) = \{\mathbf{w} \in \mathbb{R}^m : \|\mathbf{1} - \mathbf{w}\| < \varepsilon\}$, i.e. for all $\mathbf{w} \in B_\varepsilon(\mathbf{1})$, $\frac{\mathbf{1}}{2\sqrt{m}} < \frac{\mathbf{w}}{\|\mathbf{w}\|}$. By the non-conflicting property of $\mathcal{A}_{\text{UPGrad}}$, $\mathbf{0} \leq G_t \mathbf{w}_t$ and therefore for all $\mathbf{w} \in B_\varepsilon(\mathbf{1})$,

$$\begin{aligned}
& \mathbf{w}^\top (\mathbf{f}(\mathbf{x}_{t+1}) - \mathbf{f}(\mathbf{x}_t)) \\
& \leq \frac{\|\mathbf{w}\|}{\beta\sqrt{m}} \left(\frac{\mathbf{1}}{2\sqrt{m}} - \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) G_t \mathbf{w}_t \quad (\text{Lemma 6}) \\
& \leq 0.
\end{aligned}$$

Since $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t)$ is bounded and non-increasing, it converges. Since $B_\varepsilon(\mathbf{1})$ contains a basis of \mathbb{R}^m , $\mathbf{f}(\mathbf{x}_t)$ converges to some $\mathbf{f}^* \in \mathbb{R}^m$. By assumption on \mathbf{f} , there exists \mathbf{x}^* in the Pareto set satisfying $\mathbf{f}(\mathbf{x}^*) \leq \mathbf{f}^*$.

We now prove that $\mathbf{f}(\mathbf{x}^*) = \mathbf{f}^*$. Since $\mathbf{f}(\mathbf{x}^*) \leq \mathbf{f}^*$, it is sufficient to show that $\mathbf{1}^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}^*)) \leq 0$.

First, the additional assumption of Lemma 7 applies since $\mathbf{1}^\top \mathbf{f}(\mathbf{x}_t)$ decreases to $\mathbf{1}^\top \mathbf{f}^*$ which is larger than $\mathbf{1}^\top \mathbf{f}(\mathbf{x}^*)$. Therefore

$$\begin{aligned}
& \mathbf{1}^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}^*)) \\
& \leq \left(\frac{m}{T} \sum_{t \in [T]} \mathbf{w}_t \right)^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}^*)) \quad \left(\begin{array}{l} \mathbf{f}(\mathbf{x}^*) \leq \mathbf{f}^* \\ \mathbf{1} \leq m \mathbf{w}_t \\ \text{by (15a)} \end{array} \right) \\
& = \frac{m}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}_t) + \mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \\
& = \frac{m}{T} \left(\sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}_t)) + \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \right) \\
& \leq \frac{m}{T} \left(\sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}_t)) + \mathbf{1}^\top (\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)) + \frac{\beta \sqrt{m}}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \right) \quad (\text{Lemma 7})
\end{aligned} \tag{18}$$

Taking the limit as $T \rightarrow \infty$, we get

$$\begin{aligned}
& \mathbf{1}^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}^*)) \\
& \leq \lim_{T \rightarrow \infty} \frac{m}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}^* - \mathbf{f}(\mathbf{x}_t)) \\
& \leq \lim_{T \rightarrow \infty} \frac{m}{T} \sum_{t \in [T]} \|\mathbf{w}_t\| \cdot \|\mathbf{f}^* - \mathbf{f}(\mathbf{x}_t)\| \quad (\text{Cauchy-Schwartz inequality}) \\
& = 0, \quad (\mathbf{w}_t \text{ bounded}, \mathbf{f}(\mathbf{x}_t) \rightarrow \mathbf{f}^*)
\end{aligned}$$

which concludes the proof. \square

The proof of Theorem 1 provides some additional insights about the convergence rate as well as some notion of convergence in the non-convex case.

Convergence rate. Combining the equality $\mathbf{f}^* = \mathbf{f}(\mathbf{x}^*)$ and (18), we have

$$\frac{1}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \leq \frac{1}{T} \mathbf{1}^\top (\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)) + \frac{\beta \sqrt{m}}{2T} \|\mathbf{x}_1 - \mathbf{x}^*\|^2$$

Furthermore, the Cauchy-Schwartz inequality yields $\mathbf{1}^\top (\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)) \leq \sqrt{m} \cdot \|\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)\|$, so

$$\frac{1}{T} \sum_{t \in [T]} \mathbf{w}_t^\top (\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}^*)) \leq \frac{\sqrt{m}}{T} \left(\|\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)\| + \frac{\beta}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \right) \tag{19}$$

This hints a convergence rate of order $\mathcal{O}(\frac{1}{T})$, whose constant depends on the initial point \mathbf{x}_1 .

Non-convex setting. The proof of (17) does not require the convexity of the objective function. This bound can be equivalently formulated as

$$\frac{1}{T} \sum_{t=1}^T \|J_t^\top \mathbf{w}_t\|^2 \leq \frac{2\beta \sqrt{m}}{T} \mathbf{1}^\top (\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}^*)) \tag{20}$$

This shows the convergence of the updates in the non-convex setting, under the smoothness condition of Theorem 1. Note that this does not prove the convergence of \mathbf{x}_t , which is in line with the current limitations of gradient descent in the single-objective setting.

B PROPERTIES OF EXISTING AGGREGATORS

In the following, we prove the properties of the aggregators from Table 1. Some aggregators, e.g. \mathcal{A}_{RGW} , $\mathcal{A}_{\text{GradDrop}}$ and $\mathcal{A}_{\text{PCGrad}}$, are non-deterministic and are thus not technically functions but rather random variables whose distribution depends on the matrix $J \in \mathbb{R}^{m \times n}$ to aggregate. Still, the properties of Section 2.2 can be easily adapted to a random setting. If \mathcal{A} is a random aggregator, then for any J , $\mathcal{A}(J)$ is a random vector in \mathbb{R}^n . The aggregator is non-conflicting if $\mathcal{A}(J)$ is in the dual cone of the rows of J with probability 1. It is linear under scaling if for all $J \in \mathbb{R}^{m \times n}$, there is a – possibly random – matrix $\tilde{\mathcal{J}} \in \mathbb{R}^{m \times n}$ such that for all $\mathbf{0} < \mathbf{c} \in \mathbb{R}^m$, $\mathcal{A}(\text{diag}(\mathbf{c}) \cdot J) = \tilde{\mathcal{J}}^\top \cdot \mathbf{c}$. Finally, \mathcal{A} is weighted if for any $J \in \mathbb{R}^{m \times n}$ there is a – possibly random – weighting $\mathbf{w} \in \mathbb{R}^m$ satisfying $\mathcal{A}(J) = J^\top \cdot \mathbf{w}$.

B.1 MEAN

$\mathcal{A}_{\text{Mean}}$ simply averages the rows of the input matrix, i.e. for all $J \in \mathbb{R}^{m \times n}$,

$$\mathcal{A}_{\text{Mean}}(J) = \frac{1}{m} J^\top \cdot \mathbf{1} \quad (21)$$

✗ **Non-conflicting.** $\mathcal{A}_{\text{Mean}}\left(\begin{bmatrix} -2 \\ 4 \end{bmatrix}\right) = [1]$, which conflicts with $[-2]$, so $\mathcal{A}_{\text{Mean}}$ is not non-conflicting.

✓ **Linear under scaling.** For any $\mathbf{c} \in \mathbb{R}^m$, $\mathcal{A}_{\text{Mean}}(\text{diag}(\mathbf{c}) \cdot J) = \frac{1}{m} J^\top \cdot \mathbf{c}$, which is linear in \mathbf{c} . $\mathcal{A}_{\text{Mean}}$ is therefore linear under scaling.

✓ **Weighted.** By (21), $\mathcal{A}_{\text{Mean}}$ is weighted with constant weighting equal to $\frac{1}{m} \mathbf{1}$.

B.2 MGDA

The optimization algorithm presented in [Désidéri \(2012\)](#), called MGDA, is tied to a particular method for aggregating the gradients. We thus refer to this aggregator as $\mathcal{A}_{\text{MGDA}}$. The dual problem of this method was also introduced independently in [Fliege & Svaiter \(2000\)](#). We show the equivalence between the two solutions to make the analysis of $\mathcal{A}_{\text{MGDA}}$ easier.

For all $J \in \mathbb{R}^{m \times n}$, the aggregation described in [Désidéri \(2012\)](#) is defined as

$$\mathcal{A}_{\text{MGDA}}(J) = J^\top \cdot \mathbf{w} \quad (22)$$

$$\text{with } \mathbf{w} \in \arg \min_{\substack{\mathbf{0} \leq \mathbf{v}: \\ \mathbf{1}^\top \mathbf{v} = 1}} \|J^\top \mathbf{v}\|^2 \quad (23)$$

In Equation (3) of [Fliege & Svaiter \(2000\)](#), the following problem is studied:

$$\min_{\substack{\alpha \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^n: \\ J\mathbf{x} \leq \alpha \mathbf{1}}} \alpha + \frac{1}{2} \|\mathbf{x}\|^2 \quad (24)$$

We show that the problems in (23) and (24) are dual to each other. Furthermore, the duality gap is null since this is a convex problem. The Lagrangian of the problem in (24) is given by $\mathcal{L}(\alpha, \mathbf{x}, \boldsymbol{\mu}) = \alpha + \frac{1}{2} \|\mathbf{x}\|^2 - \boldsymbol{\mu}^\top (\alpha \mathbf{1} - J\mathbf{x})$. Differentiating w.r.t. α and \mathbf{x} gives respectively $1 - \mathbf{1}^\top \boldsymbol{\mu}$ and $\mathbf{x} + J^\top \boldsymbol{\mu}$. The dual problem is obtained by setting those two to $\mathbf{0}$ and then maximizing the

Lagrangian on $\mathbf{0} \leq \boldsymbol{\mu}$ and α , i.e.

$$\begin{aligned} & \arg \max_{\substack{\alpha, \mathbf{0} \leq \boldsymbol{\mu}: \\ \mathbf{1}^\top \boldsymbol{\mu} = 1}} \alpha + \frac{1}{2} \|J^\top \boldsymbol{\mu}\|^2 - \boldsymbol{\mu}^\top (\alpha \mathbf{1} + J J^\top \boldsymbol{\mu}) \\ &= \arg \max_{\substack{\alpha, \mathbf{0} \leq \boldsymbol{\mu}: \\ \mathbf{1}^\top \boldsymbol{\mu} = 1}} \alpha + \frac{1}{2} \|J^\top \boldsymbol{\mu}\|^2 - \alpha \boldsymbol{\mu}^\top \mathbf{1} - \boldsymbol{\mu}^\top J J^\top \boldsymbol{\mu} \\ &= \arg \min_{\substack{\mathbf{0} \leq \boldsymbol{\mu}: \\ \mathbf{1}^\top \boldsymbol{\mu} = 1}} \frac{1}{2} \|J^\top \boldsymbol{\mu}\|^2 \end{aligned}$$

Therefore, (23) and (24) are equivalent, with $\boldsymbol{x} = -J^\top \boldsymbol{w}$.

✓ **Non-conflicting.** Observe that since in (24), $\alpha = 0$ and $\boldsymbol{x} = \mathbf{0}$ is feasible, the objective is non-positive and therefore $\alpha \leq 0$. Substituting $\boldsymbol{x} = -J^\top \boldsymbol{w}$ in $J \cdot \boldsymbol{x} \leq \alpha \mathbf{1} \leq \mathbf{0}$ yields $\mathbf{0} \leq J J^\top \boldsymbol{w}$, i.e. $\mathbf{0} \leq J \cdot \mathcal{A}_{\text{MGDA}}(J)$, so $\mathcal{A}_{\text{MGDA}}$ is non-conflicting.

✗ **Linear under scaling.** With $J = \begin{bmatrix} 2 & 0 \\ 0 & 2 \\ a & a \end{bmatrix}$, if $0 \leq a \leq 1$, $\mathcal{A}_{\text{MGDA}}(J) = \begin{bmatrix} a \\ a \end{bmatrix}$. However, if $a \geq 1$, $\mathcal{A}_{\text{MGDA}}(J) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. This is not affine in a , so $\mathcal{A}_{\text{MGDA}}$ is not linear under scaling. In particular, if any row of J is $\mathbf{0}$, $\mathcal{A}_{\text{MGDA}}(J) = \mathbf{0}$. This implies that the optimization will stop whenever one objective has converged.

✓ **Weighted.** By (22), $\mathcal{A}_{\text{MGDA}}$ is weighted.

B.3 DUALPROJ

The projection of a gradient of interest onto a dual cone was first described in Lopez-Paz & Ranzato (2017). When this gradient is the average of the rows of the Jacobian, we call this aggregator $\mathcal{A}_{\text{DualProj}}$. Formally,

$$\mathcal{A}_{\text{DualProj}}(J) = \frac{1}{m} \cdot \pi_J (J^\top \cdot \mathbf{1}) \quad (25)$$

where π_J is the projection operator defined in (3).

✓ **Non-conflicting.** By the constraint in (3), $\mathcal{A}_{\text{DualProj}}$ is non-conflicting.

✗ **Linear under scaling.** With $J = \begin{bmatrix} 2 & 0 \\ -2a & 2a \end{bmatrix}$, if $a \geq 1$, $\mathcal{A}_{\text{DualProj}}(J) = \begin{bmatrix} 0 \\ a \end{bmatrix}$. However, if $0.5 \leq a \leq 1$, $\mathcal{A}_{\text{DualProj}}(J) = \begin{bmatrix} 1-a \\ a \end{bmatrix}$. This is not affine in a , so $\mathcal{A}_{\text{DualProj}}$ is not linear under scaling.

✓ **Weighted.** By Proposition 1, $\mathcal{A}_{\text{DualProj}}(J) = \frac{1}{m} J^\top \cdot \boldsymbol{w}$, with $\boldsymbol{w} \in \arg \min_{\mathbf{1} \leq \boldsymbol{v}} \boldsymbol{v}^\top J J^\top \boldsymbol{v}$. $\mathcal{A}_{\text{DualProj}}$ is thus weighted.

B.4 PCGRAD

$\mathcal{A}_{\text{PCGrad}}$ is described in Yu et al. (2020). It projects each gradient onto the orthogonal hyperplane of other gradients in case of conflict with them, iteratively and in random order. When $m \leq 2$, $\mathcal{A}_{\text{PCGrad}}$ is deterministic and satisfies $\mathcal{A}_{\text{PCGrad}} = m \cdot \mathcal{A}_{\text{UPGrad}}$. Therefore, in this case, it satisfies all three properties. When $m > 2$, $\mathcal{A}_{\text{PCGrad}}$ is non-deterministic, so $\mathcal{A}_{\text{PCGrad}}(J)$ is a random vector.

For any index $i \in [m]$, let $\boldsymbol{g}_i = J^\top \cdot \boldsymbol{e}_i$ and let $\mathbf{p}(i)$ be a random vector distributed uniformly on the set of permutations of the elements in $[m] \setminus \{i\}$. For instance, if $m = 3$, $\mathbf{p}(2) = [1 \ 3]^\top$ with probability 0.5 and $\mathbf{p}(2) = [3 \ 1]^\top$ with probability 0.5. For notation convenience, whenever i

is clear from context, we denote $j_k = \mathbf{p}(i)_k$. The iterative projection of $\mathcal{A}_{\text{PCGrad}}$ is then defined recursively as:

$$\mathbf{g}_{i,1}^{\text{PC}} = \mathbf{g}_i \quad (26)$$

$$\mathbf{g}_{i,k+1}^{\text{PC}} = \mathbf{g}_{i,k}^{\text{PC}} - \mathbb{1}\{\mathbf{g}_{i,k}^{\text{PC}} \cdot \mathbf{g}_{j_k} < 0\} \frac{\mathbf{g}_{i,k}^{\text{PC}} \cdot \mathbf{g}_{j_k}}{\|\mathbf{g}_{j_k}\|^2} \mathbf{g}_{j_k} \quad (27)$$

We noticed that an equivalent formulation to the conditional projection of (27) is the projection onto the dual cone of $\{\mathbf{g}_{j_k}\}$:

$$\mathbf{g}_{i,k+1}^{\text{PC}} = \pi_{\mathbf{g}_{j_k}^\top}(\mathbf{g}_{i,k}^{\text{PC}}) \quad (28)$$

Finally, the aggregation is given by

$$\mathcal{A}_{\text{PCGrad}}(J) = \sum_{i=1}^m \mathbf{g}_{i,m}^{\text{PC}}. \quad (29)$$

✗ Non-conflicting. If $J = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -0.5 & -1 \end{bmatrix}$, the only non-conflicting direction is $\mathbf{0}$. However, $\mathcal{A}_{\text{PCGrad}}(J)$ is uniform over the set $\left\{ \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.4 \\ -0.2 \end{bmatrix}, \begin{bmatrix} 0.8 \\ -0.2 \end{bmatrix} \right\}$, i.e. $\mathcal{A}_{\text{PCGrad}}(J)$ is in the dual cone of the rows of J with probability 0. $\mathcal{A}_{\text{PCGrad}}$ is thus not non-conflicting. Here, $\mathbb{E}[\mathcal{A}_{\text{PCGrad}}(J)] = [0.6 \ 0]^\top$, so $\mathcal{A}_{\text{PCGrad}}$ is neither non-conflicting in expectation.

✓ Linear under scaling. To show that $\mathcal{A}_{\text{PCGrad}}$ is linear under scaling, let $\mathbf{0} < \mathbf{c} \in \mathbb{R}^m$, $\mathbf{g}'_i = c_i \mathbf{g}_i$, $\mathbf{g}'_{i,1}{}^{\text{PC}} = \mathbf{g}'_i$ and $\mathbf{g}'_{i,k+1}{}^{\text{PC}} = \pi_{\mathbf{g}'_{j_k}{}^\top}(\mathbf{g}'_{i,k}{}^{\text{PC}})$. We show by induction that $\mathbf{g}'_{i,k}{}^{\text{PC}} = c_i \mathbf{g}_{i,k}^{\text{PC}}$.

The base case is given by $\mathbf{g}'_{i,1}{}^{\text{PC}} = \mathbf{g}'_i = c_i \mathbf{g}_i = c_i \mathbf{g}_{i,1}^{\text{PC}}$.

Then, assuming the induction hypothesis $\mathbf{g}'_{i,k}{}^{\text{PC}} = c_i \mathbf{g}_{i,k}^{\text{PC}}$, we show $\mathbf{g}'_{i,k+1}{}^{\text{PC}} = c_i \mathbf{g}_{i,k+1}^{\text{PC}}$:

$$\mathbf{g}'_{i,k+1}{}^{\text{PC}} = \pi_{c_{j_k} \mathbf{g}_{j_k}^\top} (c_i \mathbf{g}_{i,k}^{\text{PC}}) \quad (\text{Induction hypothesis})$$

$$\mathbf{g}'_{i,k+1}{}^{\text{PC}} = c_i \pi_{\mathbf{g}_{j_k}^\top} (\mathbf{g}_{i,k}^{\text{PC}}) \quad (0 < c_i \text{ and } 0 < c_{j_k})$$

$$\mathbf{g}'_{i,k+1}{}^{\text{PC}} = c_i \mathbf{g}_{i,k+1}^{\text{PC}} \quad (\text{By (28)})$$

Therefore $\mathcal{A}_{\text{PCGrad}}(\text{diag}(\mathbf{c}) \cdot J) = \sum_{i=1}^m c_i \mathbf{g}_{i,m}^{\text{PC}}$, so it can be written as $\mathcal{A}_{\text{PCGrad}}(\text{diag}(\mathbf{c}) \cdot J) = \mathfrak{J}^\top \cdot \mathbf{c}$ with $\mathfrak{J} = [\mathbf{g}_{1,m}^{\text{PC}} \ \cdots \ \mathbf{g}_{m,m}^{\text{PC}}]^\top$. Therefore, $\mathcal{A}_{\text{PCGrad}}$ is linear under scaling.

✓ Weighted. For all i , $\mathbf{g}_{i,m}^{\text{PC}}$ is always a random linear combination of rows of J . $\mathcal{A}_{\text{PCGrad}}$ is thus weighted.

B.5 GRADDROP

The aggregator used by the GradDrop layer, which we denote $\mathcal{A}_{\text{GradDrop}}$, is described in [Chen et al. \(2020\)](#). It is non-deterministic, so $\mathcal{A}_{\text{GradDrop}}(J)$ is a random vector. Given $J \in \mathbb{R}^{m \times n}$, let $|J| \in \mathbb{R}^{m \times n}$ be the element-wise absolute value of J . Let $P = \frac{1}{2} \left(\mathbf{1} + \frac{J^\top \cdot \mathbf{1}}{|J|^\top \cdot \mathbf{1}} \right) \in \mathbb{R}^n$, where the division is element-wise. Each coordinate $i \in [n]$ is independently assigned to the set \mathcal{I}_+ with probability P_i and to the set \mathcal{I}_- otherwise. The aggregation at coordinate $i \in \mathcal{I}_+$ is given by the sum of all positive J_{ji} , for $j \in [m]$. The aggregation at coordinate $i \in \mathcal{I}_-$ is given by the sum of all negative J_{ji} , for $j \in [m]$. Formally,

$$\mathcal{A}_{\text{GradDrop}}(J) = \left(\sum_{i \in \mathcal{I}_+} \mathbf{e}_i \sum_{\substack{j \in [m]: \\ J_{ji} > 0}} J_{ji} \right) + \left(\sum_{i \in \mathcal{I}_-} \mathbf{e}_i \sum_{\substack{j \in [m]: \\ J_{ji} < 0}} J_{ji} \right) \quad (30)$$

1188 **✗ Non-conflicting.** If $J = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$, then $P = [1/3]$. Therefore, $\mathbb{P}[\mathcal{A}_{\text{GradDrop}}(J) = [-2]] = 2/3$ and
 1189 $\mathbb{P}[\mathcal{A}_{\text{GradDrop}}(J) = [1]] = 1/3$, i.e. $\mathcal{A}_{\text{GradDrop}}(J)$ is in the dual cone of the rows of J with probability
 1190 0. Therefore, $\mathcal{A}_{\text{GradDrop}}$ is not non-conflicting. Here, $\mathbb{E}[\mathcal{A}_{\text{GradDrop}}(J)] = [-1]^\top$, so $\mathcal{A}_{\text{GradDrop}}$ is
 1191 neither non-conflicting in expectation.
 1192

1193
 1194 **✗ Linear under scaling.** If $J = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$, then $P = \frac{1}{2} \cdot \mathbf{1}$ and the aggregation is one of the four
 1195 vectors $[\pm 1 \quad \pm 1]^\top$ with equal probability. Scaling the first line of J by 2 yields $J = \begin{bmatrix} 2 & -2 \\ -1 & 1 \end{bmatrix}$
 1196 and $P = [2/3 \quad 1/3]^\top$, which cannot lead to a uniform distribution over four elements. Therefore,
 1197 $\mathcal{A}_{\text{GradDrop}}$ is not linear under scaling.
 1198

1199
 1200 **✗ Weighted.** With $J = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$, the span of J does not include $[1 \quad 1]^\top$ nor $[-1 \quad -1]^\top$.
 1201 Therefore, $\mathcal{A}_{\text{GradDrop}}$ is not weighted.
 1202

1203 B.6 IMTL-G

1204 In Liu et al. (2021b), the authors describe a method to impartially balance gradients by weight-
 1205 ing them. Let \mathbf{g}_i be the i 'th row of J and let $\mathbf{u}_i = \frac{\mathbf{g}_i}{\|\mathbf{g}_i\|}$. They want to find a combina-
 1206 tion $\mathbf{g} = \sum_{i=1}^m \alpha_i \mathbf{g}_i$ such that $\mathbf{g}^\top \mathbf{u}_i$ is equal for all i . Let $U = [\mathbf{u}_1 - \mathbf{u}_2 \quad \dots \quad \mathbf{u}_1 - \mathbf{u}_m]^\top$,
 1207 $D = [\mathbf{g}_1 - \mathbf{g}_2 \quad \dots \quad \mathbf{g}_1 - \mathbf{g}_m]^\top$. If $\alpha_{2:m} = [\alpha_2 \quad \dots \quad \alpha_m]^\top$, then $\alpha_{2:m} = (UD^\top)^{-1} U \cdot \mathbf{g}_1$ and
 1208 $\alpha_1 = 1 - \sum_{i=2}^m \alpha_i$. Notice that this is defined only when the gradients are linearly independent.
 1209 Thus, this is not strictly speaking an aggregator since it can only be computed on matrices of rank
 1210 m . We thus propose a generalization defined for matrices of any rank that is equivalent when the
 1211 matrix has rank m . In the original formulation, requiring $\mathbf{g}^\top \mathbf{u}_i$ to be equal to some $c \in \mathbb{R}$ for all
 1212 i , is equivalent to requiring that for all i , $\mathbf{g}^\top \mathbf{g}_i$ is equal to $c \|\mathbf{g}_i\|$. Writing $\mathbf{g} = J^\top \alpha$ and letting
 1213 $\mathbf{d} \in \mathbb{R}^m$ be the vector of norms of the rows of J , the objective is thus to find α satisfying $JJ^\top \alpha \propto \mathbf{d}$.
 1214 Besides, to match the original formulation, the elements of α should sum to 1.
 1215

1216 Letting $(JJ^\top)^\dagger$ be the Moore-Penrose pseudo inverse of JJ^\top , we define

$$1217 \mathcal{A}_{\text{IMTL-G}}(J) = J^\top \cdot \mathbf{w} \tag{31}$$

$$1218 \text{with } \mathbf{w} = \begin{cases} \frac{\mathbf{v}}{\mathbf{1}^\top \mathbf{v}}, & \text{if } \mathbf{1}^\top \mathbf{v} \neq \mathbf{0} \\ \mathbf{0}, & \text{otherwise} \end{cases} \tag{32}$$

$$1219 \text{and } \mathbf{v} = (JJ^\top)^\dagger \cdot \mathbf{d}. \tag{33}$$

1220 **✗ Non-conflicting.** If $J = [1 \quad -1 \quad -1]^\top$, then $\mathbf{d} = [1 \quad 1 \quad 1]^\top$, $JJ^\top = \begin{bmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \\ -1 & 1 & 1 \end{bmatrix}$,
 1221 and thus $\mathbf{v} = \frac{1}{9} [-1 \quad 1 \quad 1]^\top$. Therefore, $\mathbf{w} = [-1 \quad 1 \quad 1]^\top$, $\mathcal{A}_{\text{IMTL-G}}(J) = [-3]^\top$ and $J \cdot$
 1222 $\mathcal{A}_{\text{IMTL-G}}(J) = [-3 \quad 3 \quad 3]^\top$. $\mathcal{A}_{\text{IMTL-G}}$ is thus not non-conflicting.
 1223

1224 It should be noted that when J has rank m , $\mathcal{A}_{\text{IMTL-G}}$ seems to be non-conflicting. Thus, it would be
 1225 possible to make a different non-conflicting generalization, for instance, by deciding $\mathbf{0}$ when J is not
 1226 full rank.
 1227

1228 **✗ Linear under scaling.** With $J = \begin{bmatrix} a & 0 \\ 0 & 1 \end{bmatrix}$ and $a > 0$, we have $\mathbf{v} = \begin{bmatrix} 1/a^2 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} a \\ 1 \end{bmatrix} = \begin{bmatrix} 1/a \\ 1 \end{bmatrix}$
 1229 and $\mathcal{A}_{\text{IMTL-G}}(J) = \begin{bmatrix} a & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1/a \\ 1 \end{bmatrix} \cdot \frac{1}{\frac{1}{a}+1} = \frac{1}{\frac{1}{a}+1} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. This is not affine in a , so $\mathcal{A}_{\text{IMTL-G}}$ is not
 1230 linear under scaling.
 1231

1232 **✓ Weighted.** By (31), $\mathcal{A}_{\text{IMTL-G}}$ is weighted.
 1233

B.7 CAGRAD

$\mathcal{A}_{\text{CAGrad}}$ is described in Liu et al. (2021a). It is parameterized by $c \in [0, 1[$. If $c = 0$, this is equivalent to $\mathcal{A}_{\text{Mean}}$. Therefore, we restrict our analysis to the case $c > 0$. For any $J \in \mathbb{R}^{m \times n}$, let $\bar{\mathbf{g}}$ be the average gradient $\frac{1}{m} J^\top \cdot \mathbf{1}$, and let $\mathbf{e}_i^\top J$ denote the i 'th row of J . The aggregation is then defined as

$$\mathcal{A}_{\text{CAGrad}}(J) \in \arg \max_{\substack{\mathbf{d} \in \mathbb{R}^n: \\ \|\mathbf{d} - \bar{\mathbf{g}}\| \leq c \|\bar{\mathbf{g}}\|}} \min_{i \in [m]} \mathbf{e}_i^\top J \mathbf{d} \quad (34)$$

✗ Non-conflicting. Let $J = \begin{bmatrix} 2 & 0 \\ -2a - 2 & 2 \end{bmatrix}$, with a satisfying $-a + c\sqrt{a^2 + 1} < 0$. We have $\bar{\mathbf{g}} = [-a \quad 1]^\top$ and $\|\bar{\mathbf{g}}\| = \sqrt{a^2 + 1}$. Observe that any $\mathbf{d} \in \mathbb{R}^n$ satisfying the constraint $\|\mathbf{d} - \bar{\mathbf{g}}\| \leq c \|\bar{\mathbf{g}}\|$ has first coordinate at most $-a + c\sqrt{a^2 + 1}$. Because $-a + c\sqrt{a^2 + 1} < 0$, any feasible \mathbf{d} has a negative first coordinate, making \mathbf{d} conflict with the first row of J . For any $c \in [0, 1[$, $-a + c\sqrt{a^2 + 1} < 0$ is equivalent to $\sqrt{\frac{c^2}{1-c^2}} < a$. Thus, this provides a counter-example to the non-conflicting property for any $c \in [0, 1[$, i.e. $\mathcal{A}_{\text{CAGrad}}$ is not non-conflicting.

If we generalize to the case $c \geq 1$, as suggested in the original paper, then $\mathbf{d} = \mathbf{0}$ becomes feasible, which yields $\min_{i \in [m]} \mathbf{e}_i^\top J \mathbf{d} = 0$. Therefore the optimal \mathbf{d} satisfies $0 \leq \min_{i \in [m]} \mathbf{e}_i^\top J \mathbf{d}$, i.e. $\mathbf{0} \leq J \mathbf{d}$. With $c \geq 1$, $\mathcal{A}_{\text{CAGrad}}$ would thus be non-conflicting.

✗ Linear under scaling (sketch of proof). Let $J = \begin{bmatrix} 2 & 0 \\ 0 & 2a \end{bmatrix}$, then $\bar{\mathbf{g}} = [1 \quad a]^\top$ and $\|\bar{\mathbf{g}}\| = \sqrt{1 + a^2}$. One can show that the constraint $\|\mathbf{d} - \bar{\mathbf{g}}\| \leq c \|\bar{\mathbf{g}}\|$ needs to be satisfied with equality since, otherwise, we can scale \mathbf{d} to make the objective larger. Substituting J in $\min_{i \in [m]} \mathbf{e}_i^\top J \mathbf{d}$ yields $2 \min(d_1, ad_2)$. For any a satisfying $c\sqrt{1 + a^2} + 1 < a^2$, it can be shown that the optimal \mathbf{d} satisfies $d_1 < ad_2$. In that case the inner minimum over i is $2d_1$ and, to satisfy $\|\mathbf{d} - \bar{\mathbf{g}}\| = c \|\bar{\mathbf{g}}\|$, the KKT conditions over the Lagrangian yield $\mathbf{d} - \bar{\mathbf{g}} \propto \nabla_{\mathbf{d}} d_1 = [1 \quad 0]^\top$. This yields $\mathbf{d} = \begin{bmatrix} c \cdot \|\bar{\mathbf{g}}\| + 1 \\ a \end{bmatrix} = \begin{bmatrix} c\sqrt{1 + a^2} + 1 \\ a \end{bmatrix}$. This is not affine in a ; therefore, $\mathcal{A}_{\text{CAGrad}}$ is not linear under scaling.

✓ Weighted. In Liu et al. (2021a), $\mathcal{A}_{\text{CAGrad}}$ is formulated via its dual: $\mathcal{A}_{\text{CAGrad}}(J) = \frac{1}{m} J^\top \left(\mathbf{1} + \frac{c \|J^\top \mathbf{1}\|}{\|J^\top \mathbf{w}\|} \mathbf{w} \right)$, with $\mathbf{w} \in \arg \min_{\mathbf{w} \in \Delta(m)} \mathbf{1}^\top J J^\top \mathbf{w} + c \cdot \|J^\top \mathbf{1}\| \cdot \|J^\top \mathbf{w}\|$, where $\Delta(m)$ is the probability simplex of dimension m . Therefore, $\mathcal{A}_{\text{CAGrad}}$ is weighted.

B.8 RGW

\mathcal{A}_{RGW} is defined in Lin et al. (2021) as the weighted sum of the rows of the input matrix, with a random weighting. The weighting is obtained by sampling m i.i.d. normally distributed random variables and applying a softmax. Formally,

$$\mathcal{A}_{\text{RGW}}(J) = J^\top \cdot \text{softmax}(\mathbf{w}) \quad (35)$$

$$\text{with } \mathbf{w} \sim \mathcal{N}(\mathbf{0}, I) \quad (36)$$

✗ Non-conflicting. When $J = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$, the only non-conflicting solution is $\mathbf{0}$. However, $\mathbb{P}[\mathcal{A}_{\text{RGW}}(J) = \mathbf{0}] = 0$, i.e. $\mathcal{A}_{\text{RGW}}(J)$ is in the dual cone of the rows of J with probability 0. \mathcal{A}_{RGW} is thus not non-conflicting. Here,

$$\begin{aligned} \mathbb{E}[\mathcal{A}_{\text{RGW}}(J)] &= \mathbb{E} \left[\frac{e^{w_1} - 2e^{w_2}}{e^{w_1} + e^{w_2}} \right] && \text{(By (35))} \\ &= -\mathbb{E} \left[\frac{e^{w_1}}{e^{w_1} + e^{w_2}} \right] && (\mathbf{w} \sim \mathcal{N}(\mathbf{0}, I)) \\ &< 0 && (0 < e^{\mathbf{w}}) \end{aligned}$$

so \mathcal{A}_{RGW} is neither non-conflicting in expectation.

✓ **Linear under scaling.** $\mathcal{A}_{\text{RGW}}(\text{diag}(\mathbf{c}) \cdot J) = (\text{diag}(\mathbf{c}) \cdot J)^\top \cdot \text{softmax}(\mathbf{w}) = J^\top \cdot \text{diag}(\mathbf{c}) \cdot \text{softmax}(\mathbf{w}) = J^\top \cdot \text{diag}(\text{softmax}(\mathbf{w})) \cdot \mathbf{c}$. We thus have $\mathcal{A}_{\text{RGW}}(\text{diag}(\mathbf{c}) \cdot J) = \mathfrak{J}^\top \cdot \mathbf{c}$ with $\mathfrak{J}^\top = J^\top \cdot \text{diag}(\text{softmax}(\mathbf{w}))$. Therefore, \mathcal{A}_{RGW} is linear under scaling.

✓ **Weighted.** By (35), \mathcal{A}_{RGW} is weighted.

B.9 NASH-MTL

Nash-MTL is described in Navon et al. (2022). Unfortunately, we were not able to verify the proof of Claim 3.1, and we believe that the official implementation of Nash-MTL may mismatch the desired objective by which it is defined. Therefore, we only analyze the initial objective even though our experiments for this aggregator are conducted with the official implementation.

Let $J \in \mathbb{R}^{m \times n}$ and $\varepsilon > 0$. Let also $B_\varepsilon = \{\mathbf{d} \in \mathbb{R}^n : \|\mathbf{d}\| \leq \varepsilon, \mathbf{0} \leq J\mathbf{d}\}$. With $\mathbf{e}_i^\top J$ denoting the i 'th row of J , $\mathcal{A}_{\text{Nash-MTL}}$ is then defined as

$$\mathcal{A}_{\text{Nash-MTL}}(J) = \arg \max_{\mathbf{d} \in B_\varepsilon} \sum_{i \in [m]} \log(\mathbf{e}_i^\top J\mathbf{d}) \quad (37)$$

✓ **Non-conflicting.** By the constraint, $\mathcal{A}_{\text{Nash-MTL}}$ is non-conflicting.

✗ **Linear under scaling.** If an aggregator \mathcal{A} is linear under scaling, it should be the case that $\mathcal{A}(aJ) = a\mathcal{A}(J)$ for any scalar $a > 0$ and any $J \in \mathbb{R}^{m \times n}$. However, $\log(a\mathbf{e}_i^\top J\mathbf{d}) = \log(\mathbf{e}_i^\top J\mathbf{d}) + \log(a)$. This means that scaling by a scalar does not impact aggregation. Since this is not the trivial $\mathbf{0}$ aggregator, $\mathcal{A}_{\text{Nash-MTL}}$ is not linear under scaling.

✓ **Weighted.** Suppose towards contradiction that \mathbf{d} is both optimal for (37) and not in the span of J^\top . Let \mathbf{d}' be the projection of \mathbf{d} onto the span of J^\top . Since $\|\mathbf{d}'\| < \|\mathbf{d}\| < \varepsilon$ and $J\mathbf{d} = J\mathbf{d}'$, we have $J\mathbf{d} < J\left(\frac{\|\mathbf{d}\|}{\|\mathbf{d}'\|}\mathbf{d}'\right)$, contradicting the optimality of \mathbf{d} . Therefore, $\mathcal{A}_{\text{Nash-MTL}}$ is weighted.

B.10 ALIGNED-MTL

The Aligned-MTL method for balancing the Jacobian is described in Senushkin et al. (2023). For simplicity, we fix the vector of preferences to $\frac{1}{m}\mathbf{1}$, but the proofs can be adapted for any non-trivial vector. Given $J \in \mathbb{R}^{m \times n}$, let $V\Sigma^2V^\top$ be the eigen-decomposition of JJ^\top , let Σ^\dagger be the diagonal matrix whose non-zero elements are the inverse of corresponding non-zero diagonal elements of Σ and let $\sigma_{\min} = \min_{i \in [m], \Sigma_{ii} \neq 0} \Sigma_{ii}$. The aggregation is then defined as

$$\mathcal{A}_{\text{Aligned-MTL}}(J) = \frac{1}{m} J^\top \cdot \mathbf{w} \quad (38)$$

$$\text{with } \mathbf{w} = \sigma_{\min} \cdot V\Sigma^\dagger V^\top \cdot \mathbf{1} \quad (39)$$

✗ **Non-conflicting.** If the SVD of J is $V\Sigma U^\top$, then $J^\top \mathbf{w} = \sigma_{\min} U P V^\top \mathbf{1}$ with $P = \Sigma^\dagger \Sigma$ a diagonal projection matrix with 1s corresponding to non zero elements of Σ and 0s everywhere else.

Further, $J \cdot \mathcal{A}_{\text{Aligned-MTL}}(J) = \frac{\sigma_{\min}}{m} \cdot V\Sigma V^\top \mathbf{1}$. If $V = \frac{1}{2} \begin{bmatrix} \sqrt{3} & 1 \\ -1 & \sqrt{3} \end{bmatrix}$ and $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$, we have

$J \cdot \mathcal{A}_{\text{Aligned-MTL}}(J) = \frac{1}{2} \cdot V\Sigma V^\top \mathbf{1} = \frac{1}{8} [3 - \sqrt{3} \quad 1 - \sqrt{3}]^\top$ which is not non-negative. $\mathcal{A}_{\text{Aligned-MTL}}$ is thus not non-conflicting.

✗ **Linear under scaling.** If $J = \begin{bmatrix} 1 & 0 \\ 0 & a \end{bmatrix}$, then $U = V = I$, $\Sigma = J$. For $0 < a \leq 1$, $\sigma_{\min} = a$, therefore $\mathcal{A}_{\text{Aligned-MTL}}(J) = \frac{a}{2} \cdot \mathbf{1}$. For $1 < a$, $\sigma_{\min} = 1$ and therefore $\mathcal{A}_{\text{Aligned-MTL}}(J) = \frac{1}{2} \cdot \mathbf{1}$, which makes $\mathcal{A}_{\text{Aligned-MTL}}$ not linear under scaling.

✓ **Weighted.** By (38), $\mathcal{A}_{\text{Aligned-MTL}}$ is weighted.

C EXPERIMENTAL SETTINGS

For all of our experiments, we used `PyTorch` (Paszke et al., 2019). We have developed an open-source library⁴ on top of it to enable Jacobian descent easily. This library is designed to be reusable for many other use cases than the experiments presented in our work. To separate them from the library, the experiments have been conducted with a different code repository⁵ mainly using `PyTorch` and our library.

C.1 LEARNING RATE SELECTION

The learning rate has a very important impact on the speed of optimization. To make the comparisons as fair as possible, we always show the results corresponding to the best learning rate. We have selected the area under the loss curve as the criterion to compare the learning rates. This choice is arbitrary but seems to work well in practice: a lower area under the loss curve means the optimization is fast (quick loss decrease) and stable (few bumps in the loss curve). Concretely, for each random rerun and each aggregator, we first try 22 learning rates from 10^{-5} to 10^2 , increasing by a factor $10^{\frac{1}{3}}$ every time. The two best learning rates from this range then define a refined range of plausible good learning rates, going from the smallest of those two multiplied by $10^{-\frac{1}{3}}$ to the largest of those two multiplied by $10^{\frac{1}{3}}$. This margin makes it unlikely for the best learning rate to lie out of the refined range. After this, 50 learning rates from the refined range are tried. These learning rates are evenly spaced in the exponent domain. The one with the best area under the loss curve is then selected and presented in the plots. For simplicity, we have always used a constant learning rate, i.e. no learning rate scheduler was used.

This approach has the advantage of being simple and precise, thus giving trustworthy results. However, it requires 72 trainings for each aggregator, random rerun, and dataset, i.e. a total of 43776 trainings for all of our experiments. For this reason, we have opted to work on small subsets of the original datasets.

C.2 RANDOM RERUNS AND STANDARD ERROR OF THE MEAN

To get an idea of confidence in our results, every experiment is performed 8 times on a different seed and a different subset, of size 1024, of the training dataset. The seed used for run $i \in [8]$ is always simply set to i . Because each random rerun includes the full learning rate selection method described in Appendix C.1, it is sensible to consider the 8 sets of results as i.i.d. For each point of both the loss curves and the cosine similarity curves, we thus compute the estimated standard error of the mean with the usual formula $\frac{1}{\sqrt{8}} \sqrt{\frac{\sum_{i \in [8]} (v_i - \bar{v})^2}{8-1}}$, where v_i is the value of a point of the curve for random rerun i , and \bar{v} is the average value of this point over the 8 runs.

C.3 MODEL ARCHITECTURES

In all experiments, the models are simple convolutional neural networks. All convolutions always have a stride of 1×1 , a kernel size of 3×3 , a learnable bias, and no padding. All linear layers always have a learnable bias. The activation function is the exponential linear unit (Clevert et al., 2015). The full architectures are given in Tables 2, 3, 4 and 5. Note that these architectures have been fixed arbitrarily, i.e. they were not optimized through some hyper-parameter selection. The weights of the model have been initialized with the default initialization scheme of `PyTorch`.

C.4 OPTIMIZER

For all experiments except those described in Appendix D.3, we always use the basic `SGD` optimizer of `PyTorch`, without any regularization or momentum. Here, `SGD` refers to the `PyTorch` optimizer that updates the parameters of the model in the opposite direction of the gradient, which, in our case, is replaced by the aggregation of the Jacobian matrix. In the rest of this paper, `SGD` refers to the

⁴Available at https://github.com/***/***

⁵Available at https://github.com/***/***

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

Table 2: Architecture used for SVHN

Conv2d (3 input channels, 16 output channels, 1 group), ELU
Conv2d (16 input channels, 32 output channels, 16 groups)
MaxPool2d (stride of 2×2, kernel size of 2×2), ELU
Conv2d (32 input channels, 32 output channels, 32 groups)
MaxPool2d (stride of 3×3, kernel size of 3×3), ELU, Flatten
Linear (512 input features, 64 output features), ELU
Linear (64 input features, 10 outputs)

Table 3: Architecture used for CIFAR-10

Conv2d (3 input channels, 32 output channels, 1 group), ELU
Conv2d (32 input channels, 64 output channels, 32 groups)
MaxPool2d (stride of 2×2, kernel size of 2×2), ELU
Conv2d (64 input channels, 64 output channels, 64 groups)
MaxPool2d (stride of 3×3, kernel size of 3×3), ELU, Flatten
Linear (1024 input features, 128 output features), ELU
Linear (128 input features, 10 outputs)

Table 4: Architecture used for EuroSAT

Conv2d (3 input channels, 32 output channels, 1 group)
MaxPool2d (stride of 2×2, kernel size of 2×2), ELU
Conv2d (32 input channels, 64 output channels, 32 groups)
MaxPool2d (stride of 2×2, kernel size of 2×2), ELU
Conv2d (64 input channels, 64 output channels, 64 groups)
MaxPool2d (stride of 3×3, kernel size of 3×3), ELU, Flatten
Linear (1024 input features, 128 output features), ELU
Linear (128 input features, 10 outputs)

Table 5: Architecture used for MNIST, Fashion-MNIST and Kuzushiji-MNIST

Conv2d (1 input channel, 32 output channels, 1 group), ELU
Conv2d (32 input channels, 64 output channels, 1 group)
MaxPool2d (stride of 2×2, kernel size of 2×2), ELU
Conv2d (64 input channels, 64 output channels, 1 group)
MaxPool2d (stride of 3×3, kernel size of 3×3), ELU, Flatten
Linear (576 input features, 128 output features), ELU
Linear (128 input features, 10 outputs)

1458 whole stochastic gradient descent algorithm. In the experiments of Appendix D.3, we instead use
 1459 Adam to study its interactions with JD.

1460

1461 C.5 LOSS FUNCTION

1462

1463 The loss function is always the usual cross-entropy, with the default parameters of PyTorch.

1464

1465 C.6 PREPROCESSING

1466

1467 The inputs are always normalized per channel based on the mean and standard deviation computed
 1468 on the entire training split of the dataset.

1469

1470 C.7 ITERATIONS AND COMPUTATIONAL BUDGET

1471

1472 The numbers of epochs and the corresponding numbers of iterations for all datasets are provided in
 1473 Table 6, along with the required number of NVIDIA L4 GPU-hours, to run all 72 learning rates for
 1474 the 11 aggregators on a single seed. The total computational budget to run the main experiments on 8
 1475 seeds was thus around 760 GPU-hours. Additionally, we used a total of about 100 GPU hours for the
 1476 experiments varying the batch size and using Adam, and about 200 more GPU hours were used for
 1477 early investigations.

1477

1478 Table 6: Numbers of epochs, iterations, and GPU-hours for each dataset

1479

1480 Dataset	1481 Epochs	1482 Iterations	1483 GPU-Hours
1484 SVHN	25	800	17
1485 CIFAR-10	20	640	15
1486 EuroSAT	30	960	32
1487 MNIST	8	256	6
1488 Fashion-MNIST	25	800	17
1489 Kuzushiji-MNIST	10	320	8

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512 D ADDITIONAL EXPERIMENTAL RESULTS

1513
1514 In this appendix, we provide additional experimental results about IWRM.

1516 D.1 ALL DATASETS AND ALL AGGREGATORS

1517
1518 Figures 3, 4, 5, 6, 7 and 8 show the full results of the experiments described in Section 5 on SVHN,
1519 CIFAR-10, EuroSAT, MNIST, Fashion-MNIST and Kuzushiji-MNIST, respectively. For readability,
1520 the results are displayed on three different plots for each dataset. We always show $\mathcal{A}_{\text{UPGrad}}$ and $\mathcal{A}_{\text{Mean}}$
1521 for reference. The exact experimental settings are described in Appendix C.

1522 It should be noted that some of these aggregators were not developed as general-purpose aggregators,
1523 but mainly for the use case of multi-task learning, with one gradient per task. Our experiments present
1524 a more challenging setting than multi-task learning optimization because conflict between rows of
1525 the Jacobian is typically higher. Besides, for some aggregators, e.g. $\mathcal{A}_{\text{GradDrop}}$ and $\mathcal{A}_{\text{IMTL-G}}$, it was
1526 advised to make the aggregation of gradients w.r.t. an internal activation (such as the last shared
1527 representation), rather than w.r.t. the parameters of the model (Chen et al., 2020; Liu et al., 2021b).
1528 To enable comparison, we instead always aggregated the Jacobian w.r.t. all parameters.

1529 We can see that $\mathcal{A}_{\text{UPGrad}}$ provides a significant improvement over $\mathcal{A}_{\text{Mean}}$ on all datasets. Moreover, the
1530 performance gaps seem to be linked to the difficulty of the dataset, which suggests that experimenting
1531 with harder tasks is a promising future direction. The intrinsic randomness of \mathcal{A}_{RGW} and $\mathcal{A}_{\text{GradDrop}}$
1532 reduces the train set performance, but it could positively impact the generalization, which we do not
1533 study here. We suspect the disappointing results of $\mathcal{A}_{\text{Nash-MTL}}$ to be caused by issues in the official
1534 implementation that we used, leading to instability.

1536 D.2 VARYING THE BATCH SIZE

1537
1538 Figure 9 shows the results on CIFAR-10 with $\mathcal{A}_{\text{UPGrad}}$ when varying the batch size from 4 to 64.
1539 Concretely, because we are using SSJD, this makes the number of rows of the sub-Jacobian aggregated
1540 at each step vary from 4 to 64. Recall that IWRM with SSJD and $\mathcal{A}_{\text{Mean}}$ is equivalent to ERM with
1541 SGD. We observe that with a small batch size, $\mathcal{A}_{\text{UPGrad}}$ becomes very similar to $\mathcal{A}_{\text{Mean}}$. This is not
1542 surprising since both would be equivalent with a batch size of 1. Conversely, a larger batch size
1543 increases the gap between $\mathcal{A}_{\text{UPGrad}}$ and $\mathcal{A}_{\text{Mean}}$. Since the projections of $\mathcal{A}_{\text{UPGrad}}$ are onto the dual cone
1544 of more rows, each step becomes non-conflicting with respect to more of the original 1024 objectives,
1545 pushing even further the benefits of the non-conflicting property. In other words, increasing the batch
1546 size refines the dual cone, thereby improving the quality of the projections. It would be interesting to
1547 theoretically analyze the impact of the batch size in this setting.

1548 D.3 COMPATIBILITY WITH ADAM

1549
1550 Figure 10 gives the results on CIFAR-10 and SVHN when using Adam rather than the SGD optimizer.
1551 Concretely, this corresponds to the Adam algorithm in which the gradient is replaced by the aggre-
1552 gation of the Jacobian. The learning rate is still tuned as described in Appendix C.1, but the other
1553 hyperparameters of Adam are fixed to the default values of PyTorch, i.e. $\beta_1 = 0.9$, $\beta_2 = 0.999$ and
1554 $\epsilon = 10^{-8}$. Because optimization with Adam is faster, the number of epochs for SVHN and CIFAR-10
1555 is reduced to 20 and 15, respectively. While the performance gap is smaller with this optimizer, it is
1556 still significant and suggests that our methods are beneficial with other optimizers than the simple
1557 SGD. Note that this analysis is fairly superficial. The thorough investigation of the interplay between
1558 aggregators and momentum-based optimizers is a compelling future research direction.

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

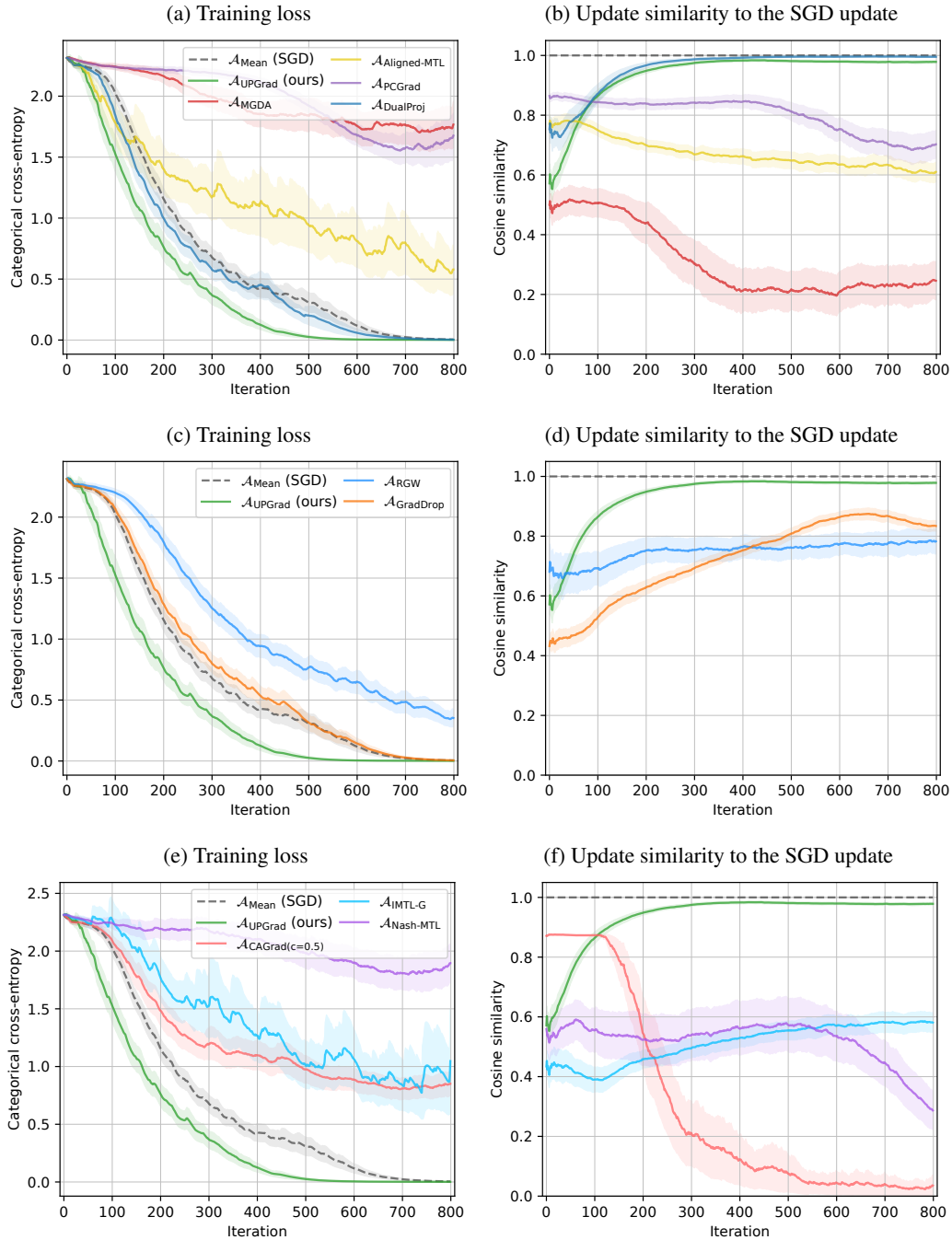


Figure 3: SVHN results.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

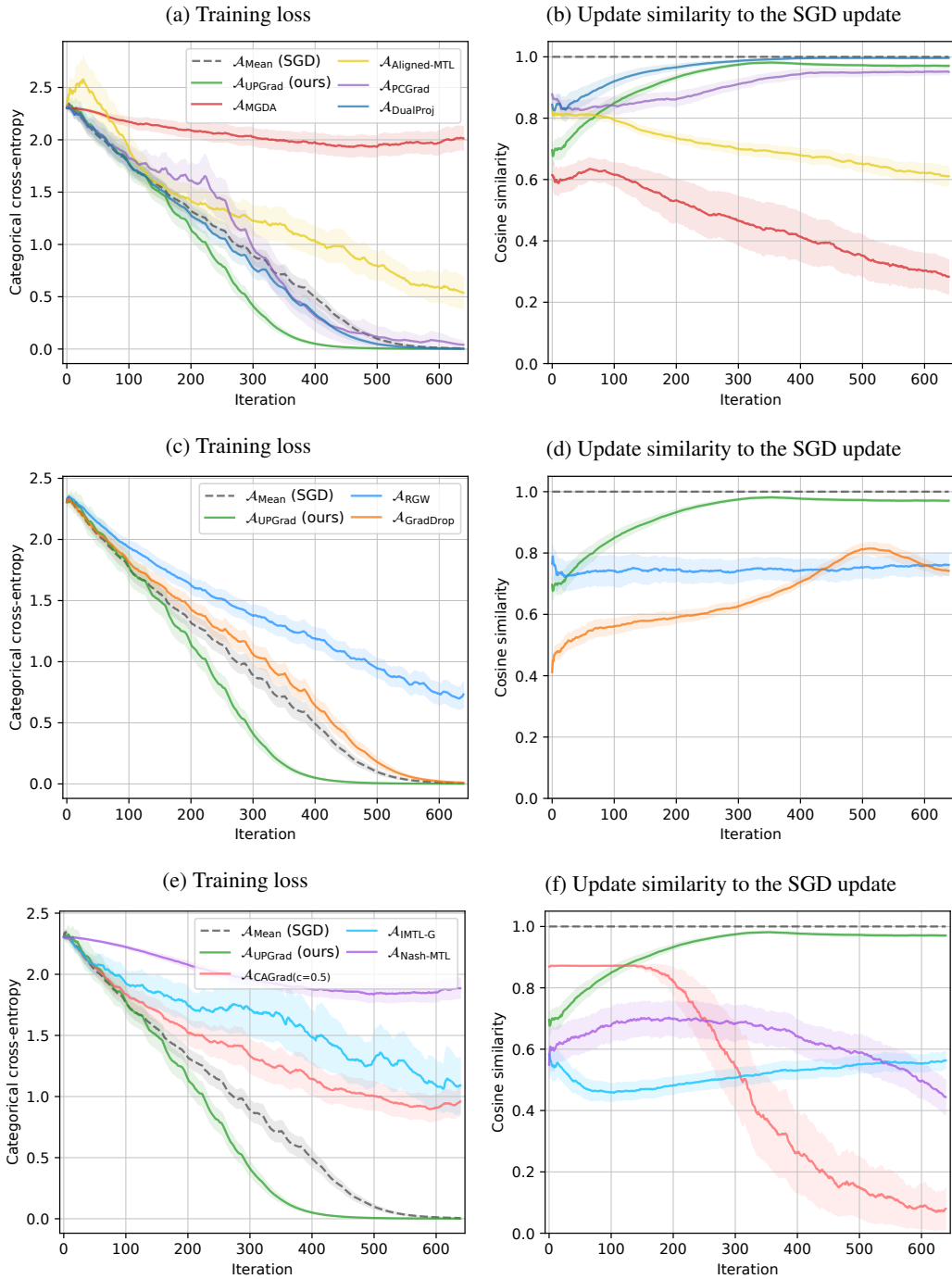


Figure 4: CIFAR-10 results.

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

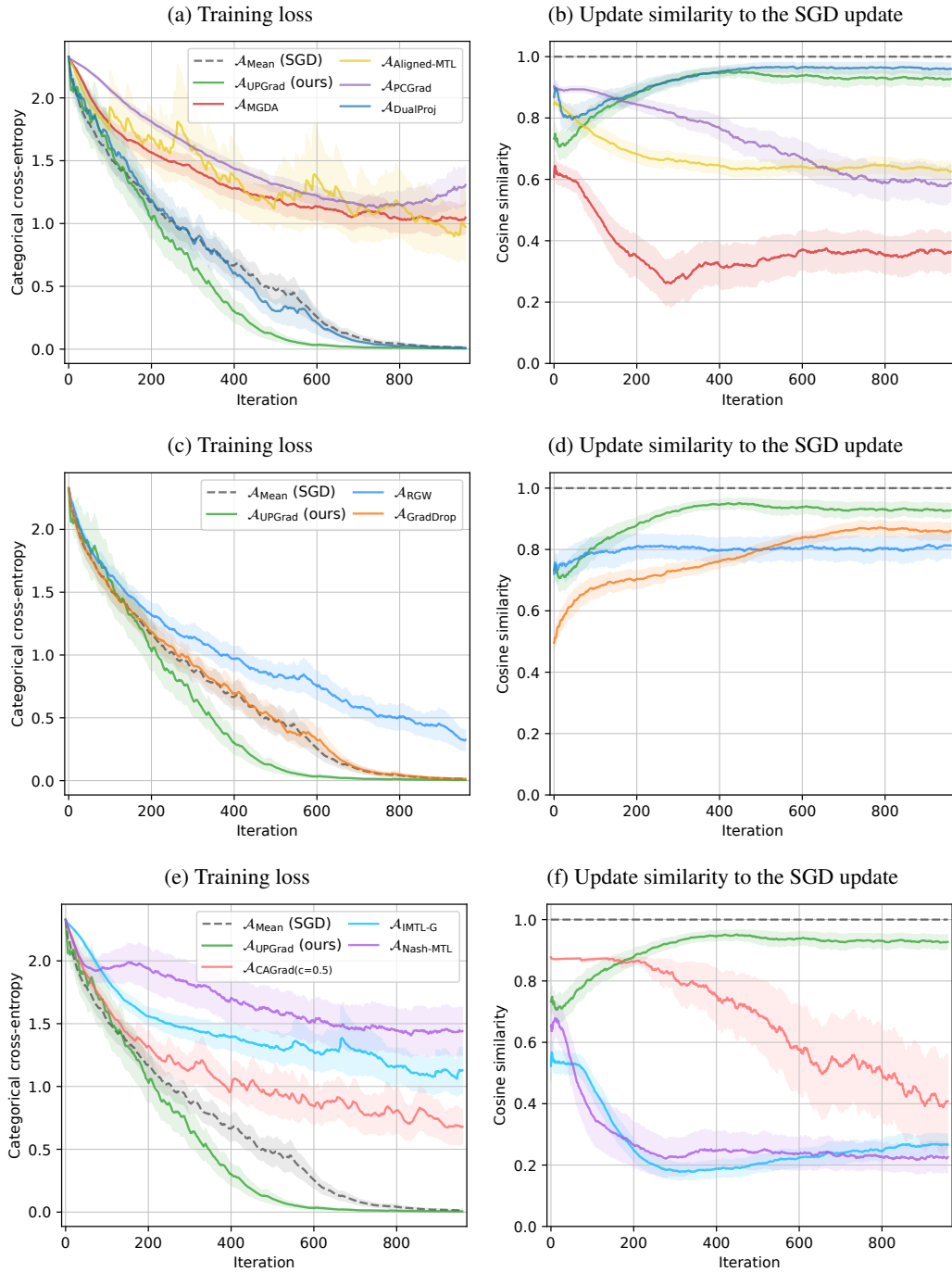


Figure 5: EuroSAT results.

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

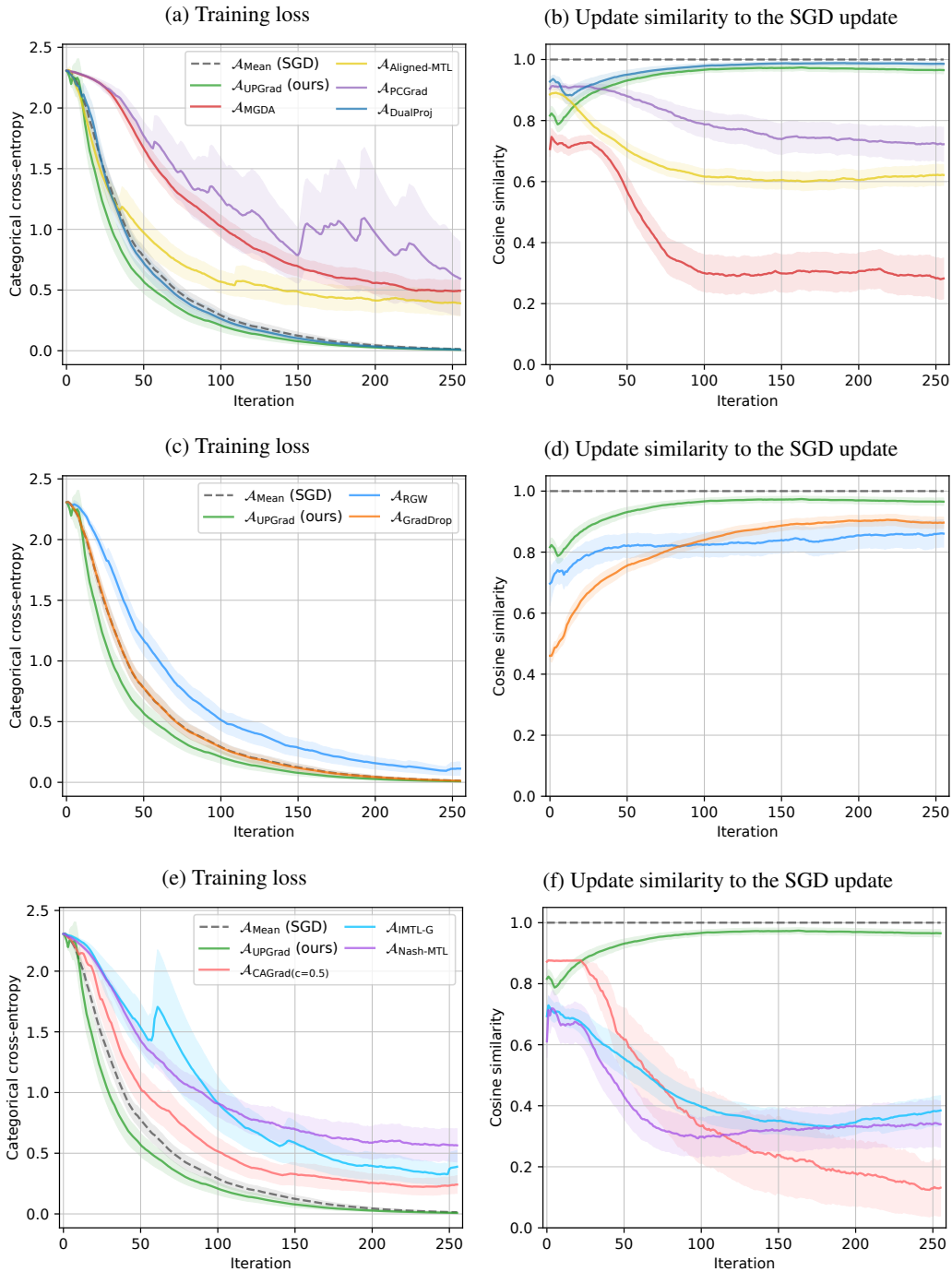


Figure 6: MNIST results.

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

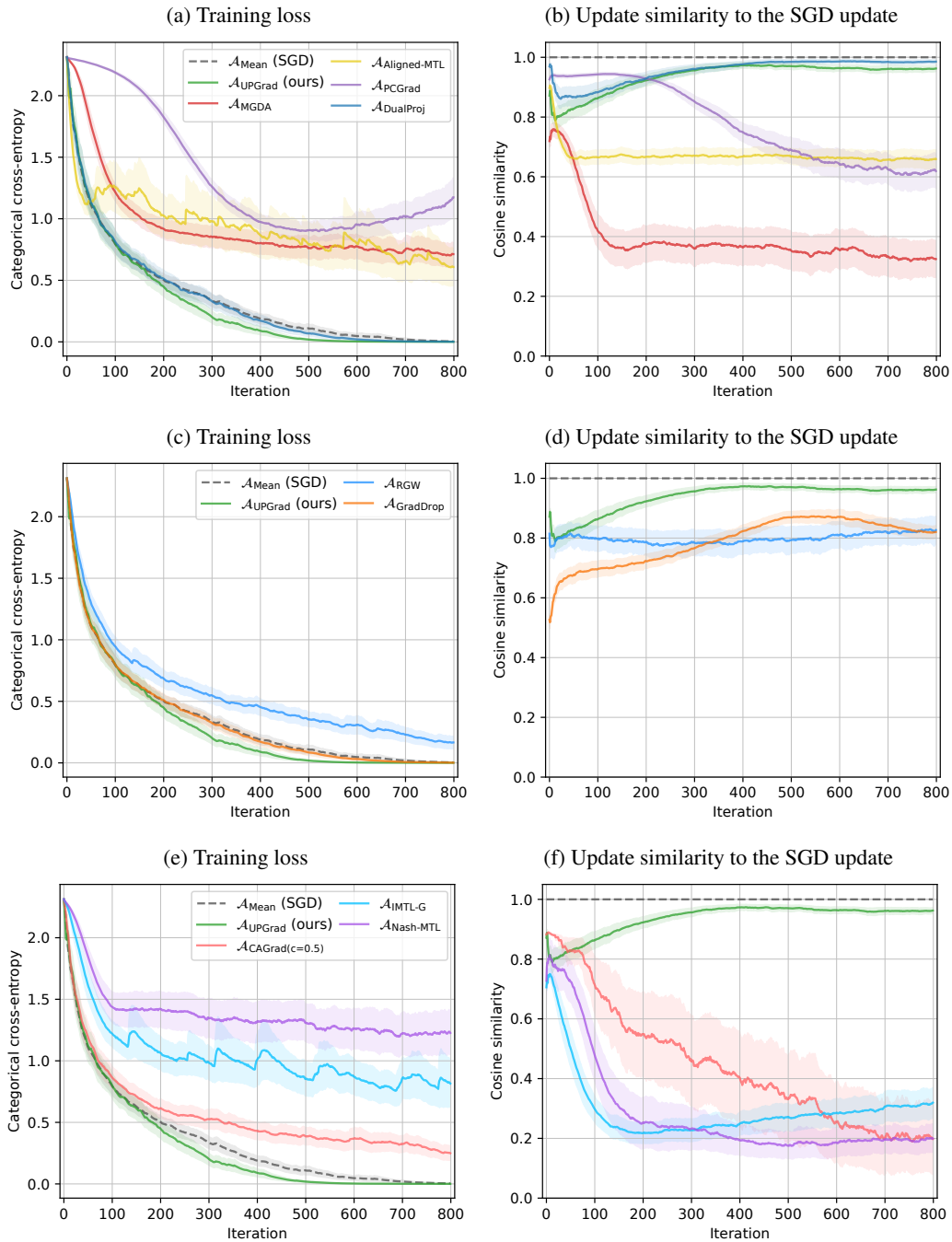


Figure 7: Fashion-MNIST results.

1836
 1837
 1838
 1839
 1840
 1841
 1842
 1843
 1844
 1845
 1846
 1847
 1848
 1849
 1850
 1851
 1852
 1853
 1854
 1855
 1856
 1857
 1858
 1859
 1860
 1861
 1862
 1863
 1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889

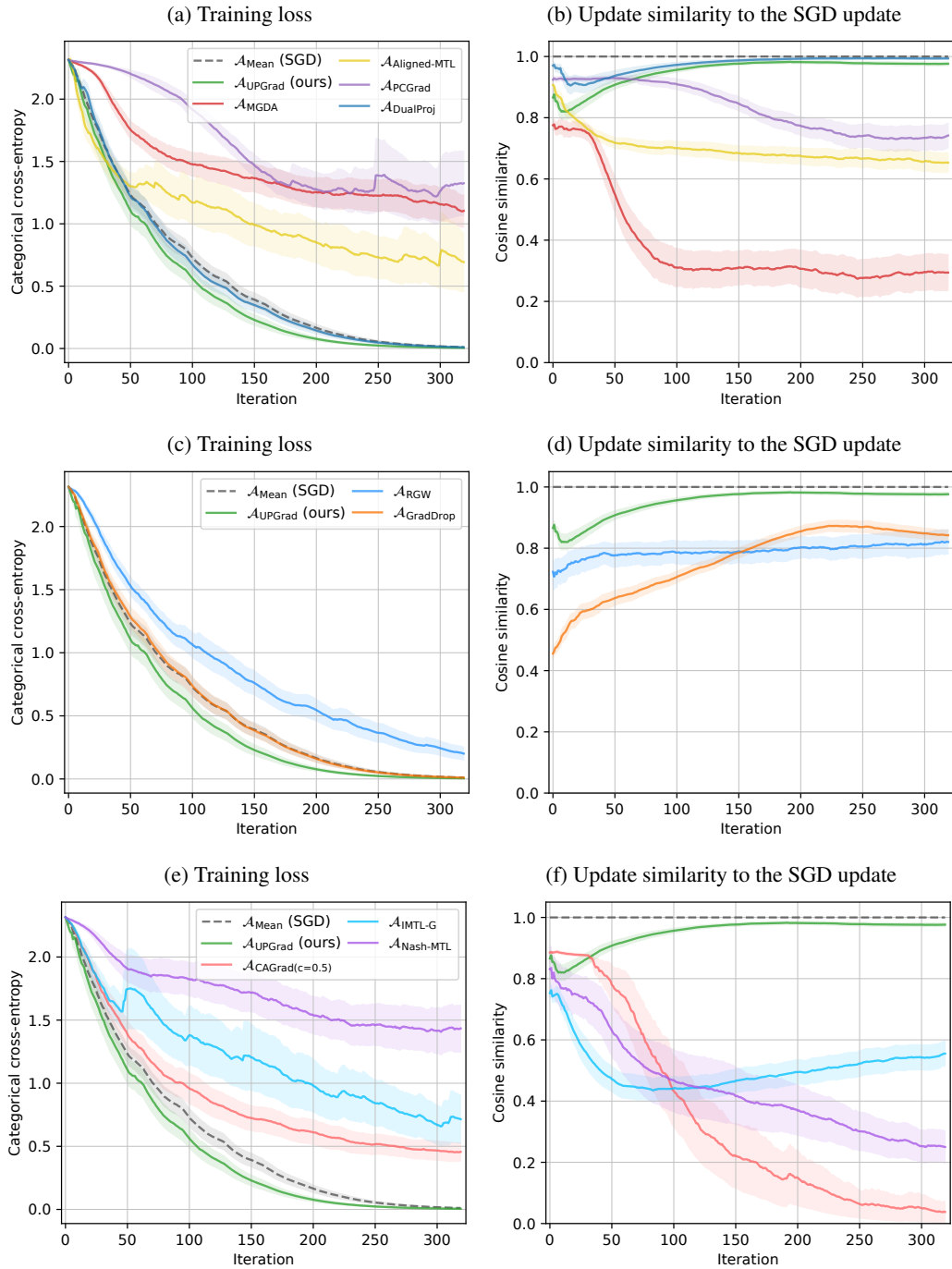


Figure 8: Kuzushiji-MNIST results.

1890
 1891
 1892
 1893
 1894
 1895
 1896
 1897
 1898
 1899
 1900
 1901
 1902
 1903
 1904
 1905
 1906
 1907
 1908
 1909
 1910
 1911
 1912
 1913
 1914
 1915
 1916
 1917
 1918
 1919
 1920
 1921
 1922
 1923
 1924
 1925
 1926
 1927
 1928
 1929
 1930
 1931
 1932
 1933
 1934
 1935
 1936
 1937
 1938
 1939
 1940
 1941
 1942
 1943

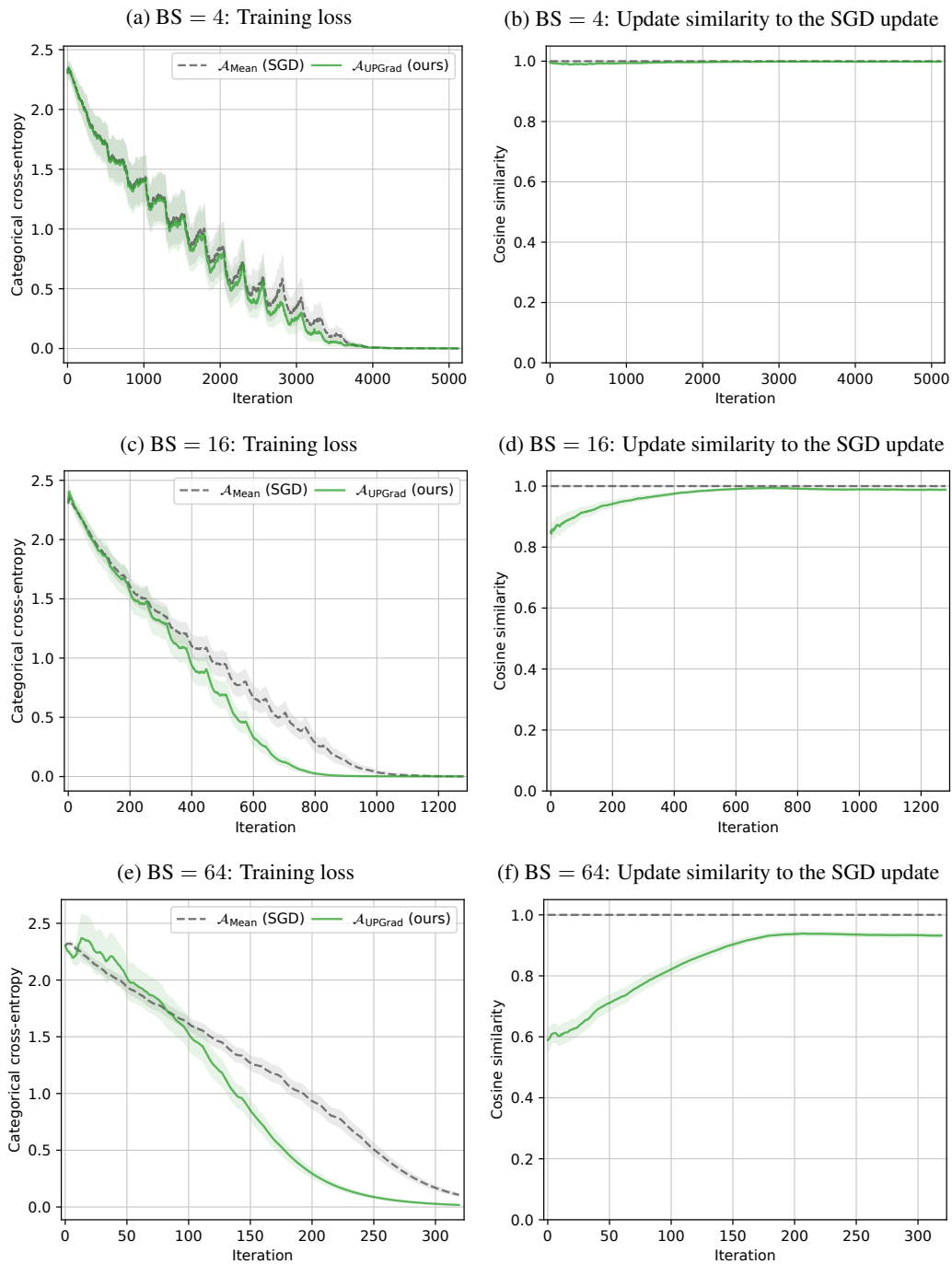


Figure 9: CIFAR-10 results with different batch sizes (BS). The number of epochs is always 20, so the number of iterations varies.

1944
 1945
 1946
 1947
 1948
 1949
 1950
 1951
 1952
 1953
 1954
 1955
 1956
 1957
 1958
 1959
 1960
 1961
 1962
 1963
 1964
 1965
 1966
 1967
 1968
 1969
 1970
 1971
 1972
 1973
 1974
 1975
 1976
 1977
 1978
 1979
 1980
 1981
 1982
 1983
 1984
 1985
 1986
 1987
 1988
 1989
 1990
 1991
 1992
 1993
 1994
 1995
 1996
 1997

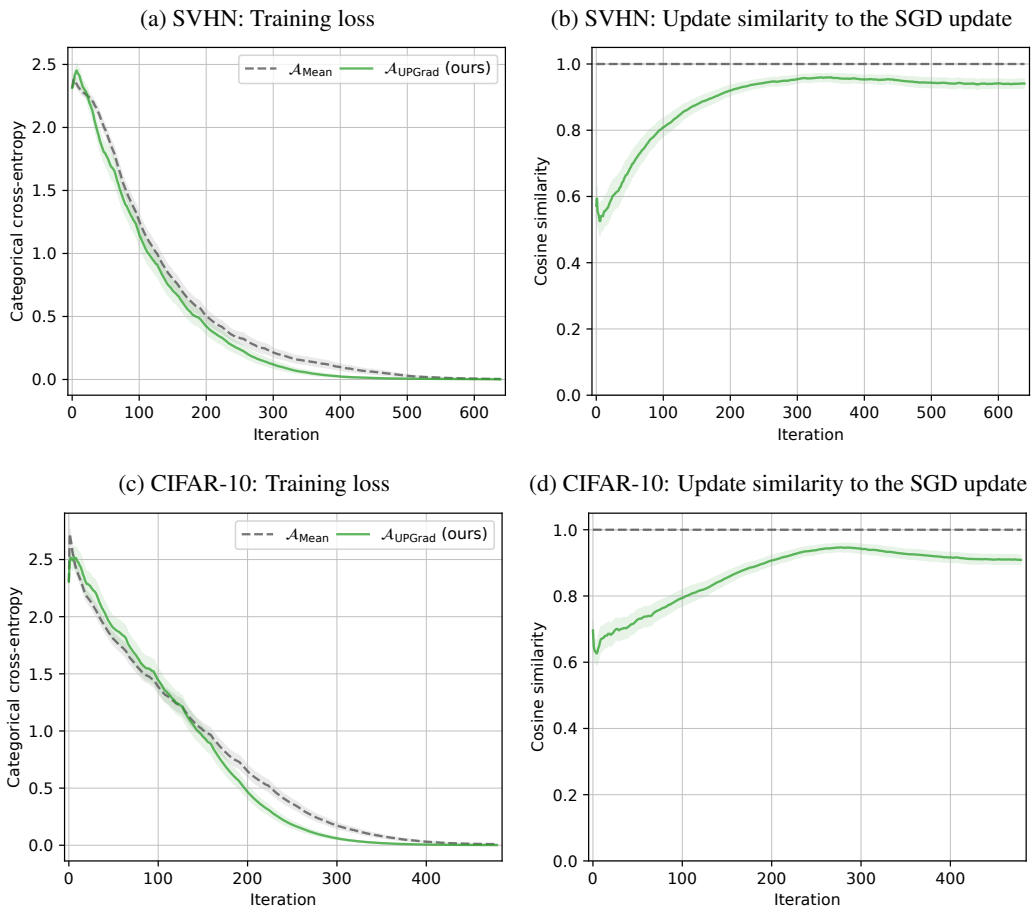


Figure 10: Results with the Adam optimizer.

1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

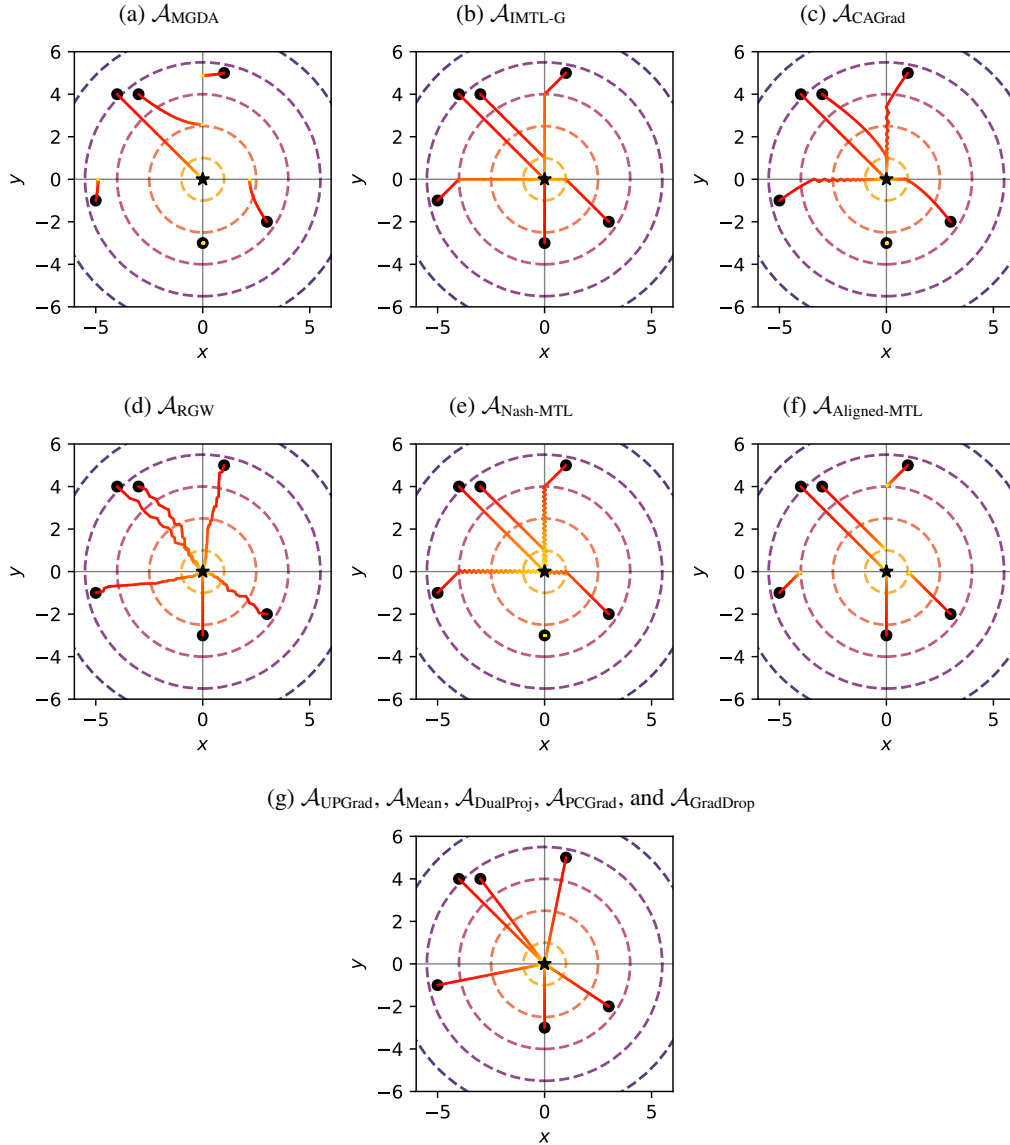


Figure 11: Optimization trajectories in the parameter space for various aggregators when optimizing $\mathbf{f} : [x \ y]^\top \mapsto [x^2 \ y^2]^\top$ with JD. Black dots represent initial parameter values, the black star is the Pareto set, the trajectories start in red and evolve towards yellow. The dashed lines are contour lines of the mean objective.

D.4 OPTIMIZATION TRAJECTORIES

Figure 11 illustrates the optimization trajectories of various aggregators from Table 1 for the function $\mathbf{f} : [x \ y]^\top \mapsto [x^2 \ y^2]^\top$, with several initializations. Notably, the Pareto set only contains the origin, while the set of Pareto stationary points is the union of the two axes. The Jacobian of \mathbf{f} at $[x \ y]^\top$ is $\begin{bmatrix} 2x & 0 \\ 0 & 2y \end{bmatrix}$, indicating that the rows do not conflict, which makes this function relatively simple to optimize. Nevertheless, several aggregators, including $\mathcal{A}_{\text{MGDA}}$, $\mathcal{A}_{\text{CAGrad}}$, $\mathcal{A}_{\text{Nash-MTL}}$ and $\mathcal{A}_{\text{Aligned-MTL}}$, fail to converge to the Pareto front for some initializations.

E COMPUTATION TIME AND MEMORY USAGE

E.1 MEMORY CONSIDERATIONS OF SSJD FOR IWRM

The main overhead of SSJD on the IWRM objective comes from having to store the full Jacobian in memory. Remarkably, when we use SGD with ERM, every activation is a tensor whose first dimension is the batch size. Automatic differentiation engines thus have to compute the Jacobian anyway. Since the gradients can be averaged at each layer as soon as they are obtained, the full Jacobian does not have to be stored. In the naive implementation of SSJD, however, storing the full Jacobian costs memory, which given the high parallelization ability of GPUs, increases the computational time. With the Gramian-based method proposed in Section 6, only the Gramian, which is typically small, has to be stored: the memory requirement would then be similar to that of SGD.

E.2 TIME COMPLEXITY OF THE UNCONFLICTING PROJECTION OF GRADIENTS

Let $J \in \mathbb{R}^{m \times n}$ be the matrix to aggregate. Apart from computing the Gramian JJ^\top , applying $\mathcal{A}_{\text{UPGrad}}$ to J requires solving m instances of the quadratic program (5) of Proposition 1. Solvers for such problems of dimension m typically have a computational complexity upper bounded by $\mathcal{O}(m^4)$ or less in recent implementations (e.g. $\mathcal{O}(m^{3.67} \log m)$). This induces a $\mathcal{O}(m^5)$ time complexity for extracting the weights of $\mathcal{A}_{\text{UPGrad}}$. Note that solving these m problems in parallel would reduce this complexity to $\mathcal{O}(m^4)$.

E.3 EMPIRICAL COMPUTATIONAL TIMES

In Table 7, we compare the computation time of SGD with that of SSJD for all the aggregators that we experimented with. Since we used the same architecture for MNIST, Fashion-MNIST and Kuzushiji-MNIST, we only report the results for one of them. Several factors affect this computation time. First, the batch size affects the number of rows in the Jacobian to aggregate. Increasing the batch size thus requires more GPU memory and the aggregation of a taller matrix. Then, some aggregators, e.g. $\mathcal{A}_{\text{Nash-MTL}}$ and $\mathcal{A}_{\text{MGDA}}$, seem to greatly increase the run time. When the aggregation is the bottleneck, a faster implementation will be necessary to make them usable in practice. Lastly, the current implementation of JD in our library is still fairly inefficient in terms of memory management, which in turn limits how well the GPU can parallelize. Also, our implementation of $\mathcal{A}_{\text{UPGrad}}$ does not solve the m quadratic programs in parallel. Therefore, these results just give a rough indication of the current computation times.

Table 7: Time required in seconds for one epoch of training on the ERM objective with SGD and on the IWRM objective with SSJD and different aggregators, on an NVIDIA L4 GPU. The batch size is always 32.

Objective	Method	SVHN	CIFAR-10	EuroSAT	MNIST
ERM	SGD	0.79	0.50	0.81	0.47
IWRM	SSJD- $\mathcal{A}_{\text{Mean}}$	1.41	1.76	2.93	1.64
IWRM	SSJD- $\mathcal{A}_{\text{MGDA}}$	5.50	5.22	6.91	5.22
IWRM	SSJD- $\mathcal{A}_{\text{DualProj}}$	1.51	1.88	3.02	1.76
IWRM	SSJD- $\mathcal{A}_{\text{PCGrad}}$	2.78	3.13	4.18	3.01
IWRM	SSJD- $\mathcal{A}_{\text{GradDrop}}$	1.57	1.90	3.06	1.78
IWRM	SSJD- $\mathcal{A}_{\text{IMTL-G}}$	1.48	1.79	2.94	1.69
IWRM	SSJD- $\mathcal{A}_{\text{CAGrad}}$	1.93	2.26	3.42	2.17
IWRM	SSJD- \mathcal{A}_{RGW}	1.42	1.76	2.89	1.73
IWRM	SSJD- $\mathcal{A}_{\text{Nash-MTL}}$	7.88	8.12	9.33	7.91
IWRM	SSJD- $\mathcal{A}_{\text{Aligned-MTL}}$	1.53	1.98	2.97	1.71
IWRM	SSJD- $\mathcal{A}_{\text{UPGrad}}$	1.80	2.01	3.21	1.90