

A Multi-objective / Multi-task Learning Framework Induced by Pareto Stationarity

Michinari Momma¹, Chaosheng Dong¹, Jia Liu²

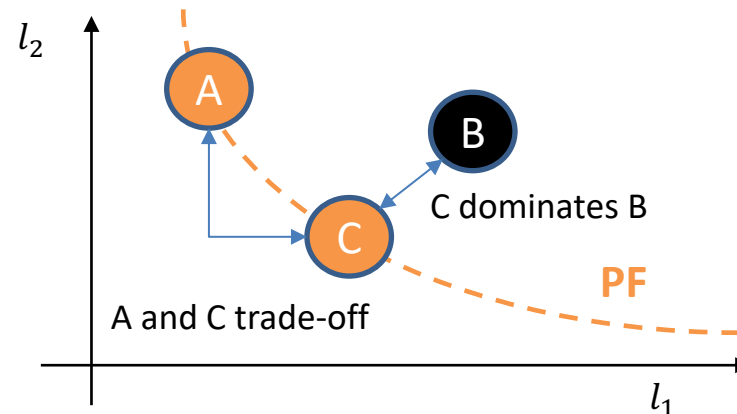
¹Amazon.com, ²The Ohio State University

Introduction

- MOO / MTL has become very popular modeling technique
 - MOO: Product search models are trained with multi-objectives
 - MTL: Popular way to training deep learning models
- Trade-off: common property in MOO / MTL
 - Many models with different trade-offs
 - Need to specify preference to select a model for experiments, deployment, etc.
- Challenges
 1. How to explore (all) trade-offs over multi-objectives / tasks
 2. How to train a (specific) model that **aligns with user preference**
 3. How to improve all objective / tasks when **re-train / fine-tune existing models**

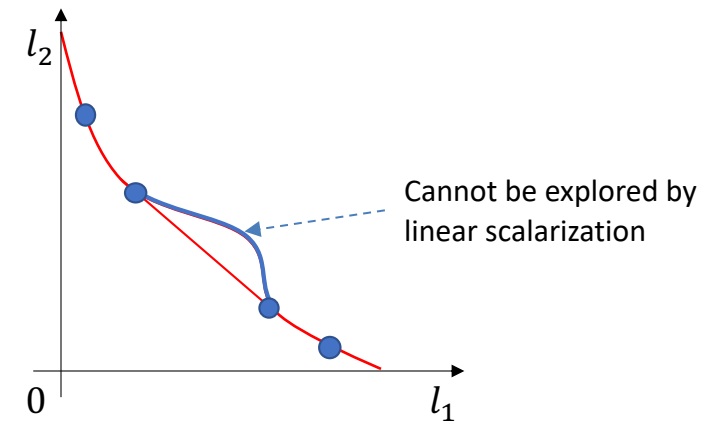
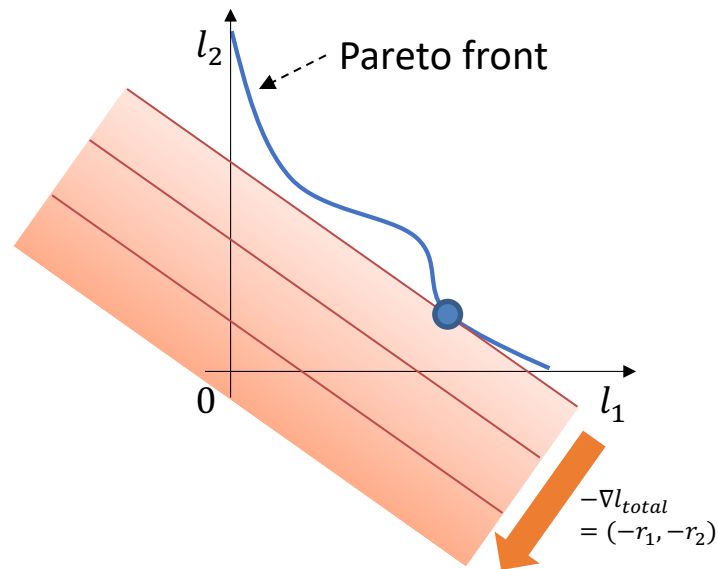
MOO / MTL problem setting

- Given a model parameter vector $\mathbf{x} \in R^n$, loss for m objectives / tasks: $l_i(\mathbf{x}), i = 1, \dots, m$
- Goal: $\min \mathbf{l}(\mathbf{x}) = [l_1(\mathbf{x}), \dots, l_m(\mathbf{x})]$ ($n \gg m$)
- Solutions typically have **trade-offs**, where there is no better solution (i.e., **dominates**) – **Pareto optimal (PO)** solution
- **Pareto Front (PF)** is a set of PO solutions.



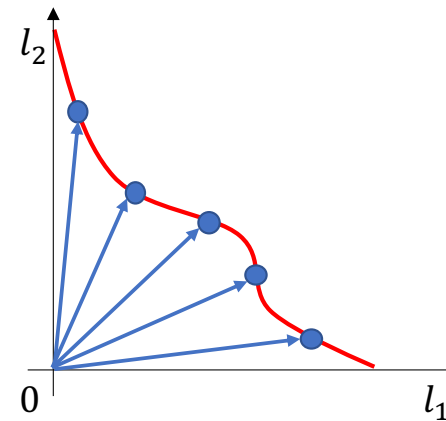
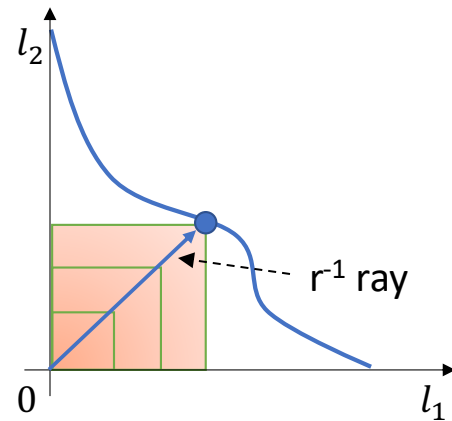
Example of preferences (Linear scalarization)

- Linear scalarization ($l_{LS} = r_1 l_1 + r_2 l_2$)
 - Higher r , the smaller l .
 - “level set” $r_1 l_1 + r_2 l_2 = \text{const}$, on a tangential hyperplane
 - Cannot explore non-convex portion of PF



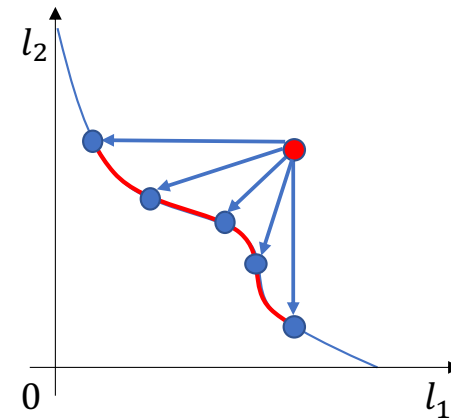
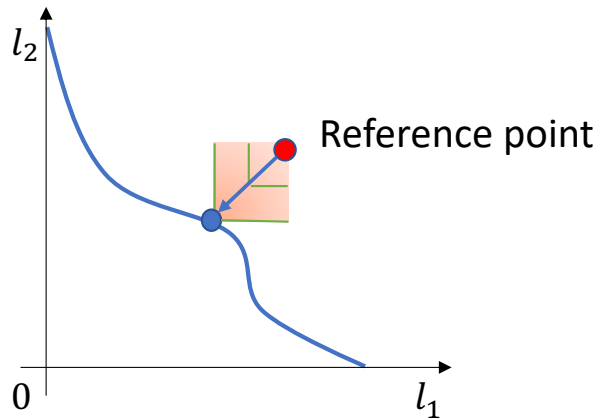
Example of preference (weighted Chebyshev)

- Ratio between loss ($r_1 l_1 = r_2 l_2$) – “weighted Chebyshev (WC)”
 - Higher r , the smaller l .
 - Level set on a hyper-rectangular box:
$$r_1 l_1 = r_2 l_2 = c \Rightarrow (l_1, l_2) = \left(\frac{c}{r_1}, \frac{c}{r_2} \right) : \mathbf{r}^{-1} \text{ray}$$
 - Can explore full PF by various \mathbf{r}^{-1} rays
 - EPO Search (ICML 2020) solves WC



Example of preference (Reference point)

- Preference from **Reference point** – Extended Weighted Chebyshev (XWC)
 - Designed to find solution better than the reference point
 - Can explore a **specific portion** of PF, pivoted on the reference point
 - Useful for model retraining / finetuning on an existing baseline / pretrained model



Gradient based method to solve PO: MGDA

- Multiple Gradient Descent Algorithm (MGDA):

- Gradient matrix: $K = f(\nabla l_1, \dots, \nabla l_m) \in R^{m \times m}$

- Pareto stationarity:

$$K\alpha = 0, e^T \alpha = 1, \alpha \geq 0, e = [1, \dots, 1]$$

- MGDA solves: $\min_{\alpha} \|K\alpha\|_2, \text{ s.t. } e^T \alpha = 1, \alpha \geq 0$

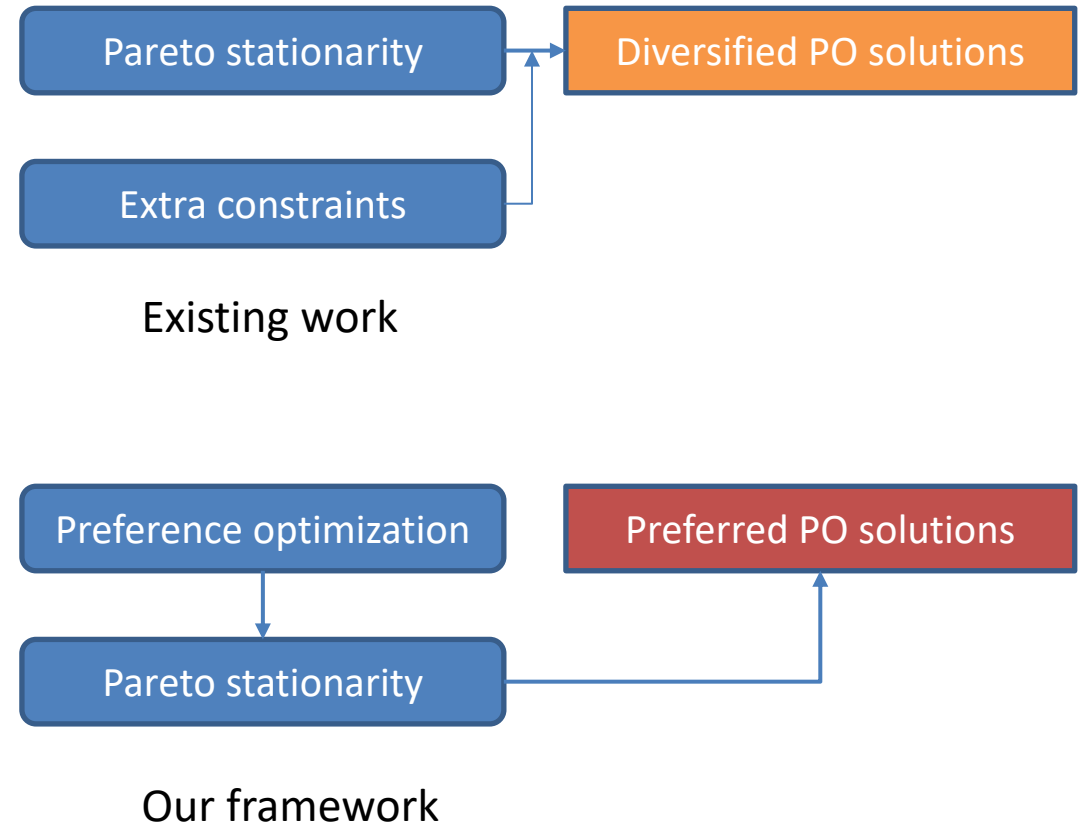
- Once we find α , run gradient descent

$$x^{new} := x^{old} - \eta G\alpha$$

- [challenge] Although MGDA finds any PO, **no control** on what solution to discover.

Improvement over MGDA

- Existing approaches try to modify MGDA (i.e. Pareto stationarity) by adding extra constraints
 - PMTL (Lin et al., NeurIPS 2019), etc.
- Note Pareto stationarity itself is a part of optimality condition in “some” optimization problem
- We formulate “preference” first, then derive Pareto stationarity



Implementation of WC: WC-MGDA

- Weighted Chebyshev problem: ℓ_∞ -norm minimization of weighted loss functions

$$\min \rho \quad s.t. \quad \mathbf{r} \odot \mathbf{l}(\mathbf{x}) \leq \rho \mathbf{e}$$

- By KKT conditions and $\mathbf{K}_r = \text{diag}(\mathbf{r})\mathbf{K}\text{diag}(\mathbf{r})$, Wolfe Dual is

$$\begin{aligned} \max \quad & \boldsymbol{\alpha}^T (\mathbf{r} \odot \mathbf{l}(\mathbf{x})) \\ s.t. \quad & \mathbf{e}^T \boldsymbol{\alpha} = 1, \boldsymbol{\alpha} \geq 0, \mathbf{K}_r \boldsymbol{\alpha} = \mathbf{0} \end{aligned}$$

- Like MGDA, we minimize ℓ_2 -norm of $\mathbf{K}_r \boldsymbol{\alpha}$:

$$\begin{aligned} \max \quad & \boldsymbol{\alpha}^T (\mathbf{r} \odot \mathbf{l}(\mathbf{x})) - u\gamma \\ s.t. \quad & \mathbf{e}^T \boldsymbol{\alpha} = 1, \boldsymbol{\alpha} \geq 0, \|\mathbf{K}_r \boldsymbol{\alpha}\|_2 \leq \gamma \end{aligned}$$

“MGDA” portion to optimize Pareto stationarity

- “trade-off” u is auto-tuned

Formulation of XWC: XWC-MGDA

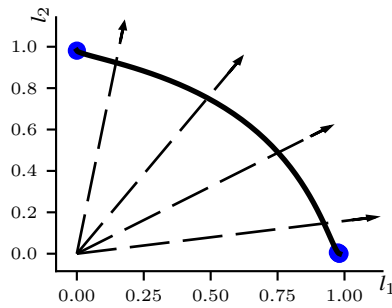
- Extend WC by two modifications:
 - Including reference point (\mathbf{b}) to the formulation
 - Setting lower bound of α to ensure strict optimality ($\alpha \geq \mathbf{w}$)
- The XWC-MGDA problem is;

$$\begin{aligned} \max \quad & \alpha^T (\mathbf{r} \odot (\mathbf{l}(\mathbf{x}) - \mathbf{b})) - u\gamma \\ \text{s.t.} \quad & \mathbf{e}^T \alpha = 1, \alpha \geq \mathbf{w}, \|\mathbf{K}_r \alpha\|_2 \leq \gamma \end{aligned}$$

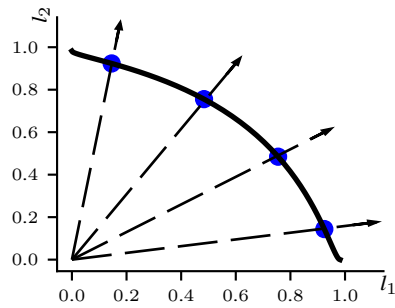
Experimental results

Experiment on non-convex synthetic data:

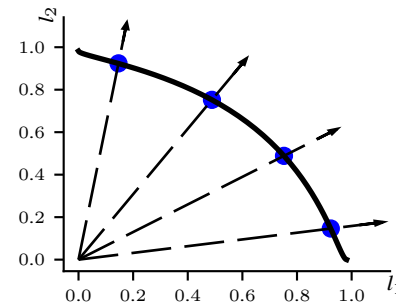
$$l(\mathbf{x}) = \left[1 - \exp\left(-\left\|\mathbf{x} - \frac{1}{\sqrt{n}}\mathbf{e}\right\|^2\right), 1 - \exp\left(-\left\|\mathbf{x} + \frac{1}{\sqrt{n}}\mathbf{e}\right\|^2\right) \right]$$



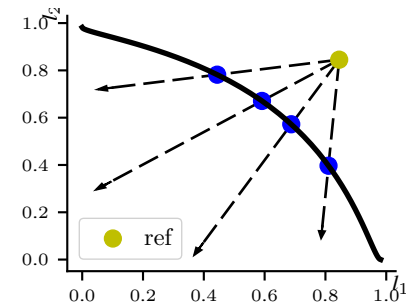
LinScalar



EPO



WC-MGDA



XWC-MGDA
(w/ ref point)

Summary

- Proposed a new framework to formulate preferences in MOO / MTL
- XWC-MGDA allows us to explore PF pivoted on a given reference point
- Experimental results (see main paper) quantify competitive performance of XWC-MGDA