

# Multi-objective machine learning to predict Pareto fronts

Timo M. Deist

Centrum Wiskunde & Informatica (CWI), Amsterdam (until 2021)

May 6, 2022



Centrum Wiskunde & Informatica



# Acknowledgments

## Co-authors:

Stef  
Maree



Monika  
Grewal



Frank  
Dankers



Tanja  
Alderliesten



Peter  
Bosman




## Feedback:


Marco Virgolin (CWI)




# The Beach-Conference conundrum




# The Beach-Conference conundrum



# The Beach-Conference conundrum



# The Beach-Conference conundrum



# Multi-objective (a posteriori) decision-making<sup>1</sup>

1. Optimize two well-defined, competing objectives
  - ↓ Distance to conference
  - ↓ Distance to beach
2. Present decision-maker with  $p$  solutions on the Pareto front
3. Decision-maker chooses solution


**What is needed:** The Pareto front.


## Problem statement


Find  $p$  solutions that span the Pareto front.

---


<sup>1</sup>Thiele et al. (2009)









Parameter space




Objective space



Parameter space




Objective space




**Goal:** find a set of solutions evenly spread across the Pareto front

Parameter space




Objective space




$$d = w_0 \frac{\partial f_0(x)}{\partial x} + w_1 \frac{\partial f_1(x)}{\partial x}$$

Parameter space



Objective space




$$d_i = a \frac{\partial f_0(x)}{\partial x} + (1 - a) \frac{\partial f_1(x)}{\partial x} \quad \forall a_i \in \{0.1, 0.2, \dots, 1\}$$

What can happen for some MO problems...<sup>2</sup>


---

<sup>1</sup>Das and Dennis (1997)


## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## MO optimization example




## Objective space







Dominated **hypervolume (HV)**




Dominated **hypervolume (HV)**




The dominated **hypervolume (HV)** is maximal in solutions that are spread over the front.




The dominated **hypervolume (HV)** is maximal in solutions that are spread over the front.





The dominated **hypervolume (HV)** is maximal in solutions that are spread over the front.



HV selects good sets of solutions. Can we compute HV gradients to help finding those sets?








$$\frac{\partial \text{HV}}{\partial \mathbf{x}_i} = \frac{\partial \text{HV}}{\partial f_0(\mathbf{x}_i)} \frac{\partial f_0(\mathbf{x}_i)}{\partial \mathbf{x}_i} + \frac{\partial \text{HV}}{\partial f_1(\mathbf{x}_i)} \frac{\partial f_1(\mathbf{x}_i)}{\partial \mathbf{x}_i}$$


---

<sup>2</sup>Emmerich and Deutz (2014)

Parameter space




Objective space




$$\frac{\partial \text{HV}}{\partial \mathbf{x}_i} = \frac{\partial \text{HV}}{\partial f_0(\mathbf{x}_i)} \frac{\partial f_0(\mathbf{x}_i)}{\partial \mathbf{x}_i} + \frac{\partial \text{HV}}{\partial f_1(\mathbf{x}_i)} \frac{\partial f_1(\mathbf{x}_i)}{\partial \mathbf{x}_i}$$


Parameter space



Objective space




$$\frac{\partial \text{HV}}{\partial \mathbf{x}_i} = \frac{\partial \text{HV}}{\partial f_0(\mathbf{x}_i)} \frac{\partial f_0(\mathbf{x}_i)}{\partial \mathbf{x}_i} + \frac{\partial \text{HV}}{\partial f_1(\mathbf{x}_i)} \frac{\partial f_1(\mathbf{x}_i)}{\partial \mathbf{x}_i}$$



**Parameter space**

1) Initialize  $p$  solutions



**Objective space**

2) Compute  $\frac{\partial HV}{\partial f}$   
for non-dominated solutions

3) Compute  $\frac{\partial UD}{\partial f}$   
for dominated solutions

4) Compute update  
direction  $\frac{\partial UHV}{\partial x}$


5) Move solutions

6) Repeat until convergence

## Notation

MO optimization		MO learning
Optimization instance	→	Samples $s_k$
Variables $x_i$	→	Model parameters $\theta_i$
Objectives $f_j$	→	Losses $L_j$

## MO learning problem




- ▶ Each sample has its own Pareto front

$$\theta_1$$


$$\theta_2$$

$$\theta_3$$


# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample




# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample




# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample




# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample




# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample
- ▶ How? Train learners so that the average HV is maximal for all samples




# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample
- ▶ How? Train learners so that the average HV is maximal for all samples




# MO learning problem




- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample
- ▶ How? Train learners so that the average HV is maximal for all samples





# MO learning problem





- ▶ Each sample has its own Pareto front
- ▶ What we want: learners generating Pareto optimal points for each sample
- ▶ How? Train learners so that the average HV is maximal for all samples





## Forward pass



## Forward pass



## Forward pass      Compute HV gradients




## Forward pass      Compute HV gradients      Compute dynamic losses

Forward pass:

$$\begin{array}{l}
 \theta_1 \rightarrow L_1(\theta_1, s_k) \\
 s_k \rightarrow L_2(\theta_1, s_k) \\
 \theta_2 \rightarrow L_1(\theta_2, s_k) \\
 s_k \rightarrow L_2(\theta_2, s_k) \\
 \theta_3 \rightarrow L_1(\theta_3, s_k) \\
 s_k \rightarrow L_2(\theta_3, s_k)
 \end{array}$$

Compute HV gradients:


$$\left. \begin{array}{l}
 \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial \mathcal{L}(\theta_1, s_k)} \rightarrow \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_1(\theta_1, s_k)} L_1(\theta_1, s_k) + \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_2(\theta_1, s_k)} L_2(\theta_1, s_k) \\
 \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial \mathcal{L}(\theta_2, s_k)} \rightarrow \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_1(\theta_2, s_k)} L_1(\theta_2, s_k) + \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_2(\theta_2, s_k)} L_2(\theta_2, s_k) \\
 \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial \mathcal{L}(\theta_3, s_k)} \rightarrow \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_1(\theta_3, s_k)} L_1(\theta_3, s_k) + \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_2(\theta_3, s_k)} L_2(\theta_3, s_k)
 \end{array} \right\}$$



## Forward pass      Compute HV gradients      Compute dynamic losses

$$\begin{array}{c}
 \begin{array}{ccc}
 \text{Forward pass} & \text{Compute HV gradients} & \text{Compute dynamic losses}
 \end{array} \\
 \begin{array}{ccc}
 \begin{array}{c}
 \theta_1 \xrightarrow{s_k} L_1(\theta_1, s_k) \\
 \theta_2 \xrightarrow{s_k} L_2(\theta_1, s_k) \\
 \theta_3 \xrightarrow{s_k} L_1(\theta_2, s_k) \\
 \theta_2 \xrightarrow{s_k} L_2(\theta_2, s_k) \\
 \theta_3 \xrightarrow{s_k} L_1(\theta_3, s_k) \\
 \theta_1 \xrightarrow{s_k} L_2(\theta_3, s_k)
 \end{array} & \left\{ \begin{array}{l}
 \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial \mathcal{L}(\theta_1, s_k)} \rightarrow \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_1(\theta_1, s_k)} L_1(\theta_1, s_k) + \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_2(\theta_1, s_k)} L_2(\theta_1, s_k) \\
 \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial \mathcal{L}(\theta_2, s_k)} \rightarrow \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_1(\theta_2, s_k)} L_1(\theta_2, s_k) + \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_2(\theta_2, s_k)} L_2(\theta_2, s_k) \\
 \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial \mathcal{L}(\theta_3, s_k)} \rightarrow \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_1(\theta_3, s_k)} L_1(\theta_3, s_k) + \frac{\partial \text{HV}(\mathcal{L}(\Theta, s_k))}{\partial L_2(\theta_3, s_k)} L_2(\theta_3, s_k)
 \end{array} \right. & 
 \end{array}
 \end{array}$$

Backpropagate



## MO regression


Given:  $x_k, X \in [0, 2\pi]$



Predict:  $z_k$  that matches  $y_k^{(j)}$

$$Y_1 = \cos(X), \quad Y_2 = \sin(X)$$

Use Mean Square Error (MSE) for learning:


$$L_j = \text{MSE}_j = \frac{1}{|S|} \sum_{k=1}^{|S|} (y_k^{(j)} - z_k)^2$$






# Organ segmentation

## Prostate MRI



2 segmentations



Given: MRI scan<sup>a</sup>  
Predict: Prostate segmentation  
trading off **both**  
segmentations.

$L_1$  = Cross Entropy w.r.t. segmentation 1


$L_2$  = Cross Entropy w.r.t. segmentation 2

---


<sup>a</sup>Data described in Dushatskiy et al. (2020)

# Organ segmentation

Delineations




Predictions



- Observer 1
- Observer 2
- Net 1
- Net 2
- Net 3
- Net 4
- Net 5

## Organ segmentation



# Neural style transfer

Photo



Popularized by Gatys et al. (2016)

Style



Given: Photo & style image  
Optimize: Image trading off photo content  
and style match

$L_1$  = Content loss

$L_2$  = Style loss

Photo by J.C.M. Dankers; Style by R. Lichtenstein, *Drowning Girl*

# Neural style transfer (Gatys et al., 2016)



# Neural style transfer

Photo



Style




How much style do you want?  
10%?  
20%?  
100%?

# Neural style transfer



# Neural style transfer



# Neural multi-style transfer



2



Fanny Tellier



Content



5



1



3



6



4





Frenchman's Bay



Kai Province

# Neural multi-style transfer



# Conclusions

MO optimization to guide decision-making with conflicting goals

HV gradient ascent offers

- ▶ single objective & gradient-based search
- ▶ finding diverse sets of solutions close to Pareto front

We proposed MO learning based on HV-maximization

- ▶ generates Pareto front approximations per sample
- ▶ implemented in PyTorch with prototypes for
  - MO regression
  - medical imaging
  - neural style transfer
- ▶ also works for asymmetric Pareto fronts

# MO learning literature

Cornell University

arXiv.org > cs > arXiv:2102.04523

Computer Science > Machine Learning

[Submitted on 8 Feb 2021]

## Multi-Objective Learning to Predict Pareto Fronts Using Hypervolume Maximization

Timo M. Deist, Monika Grewal, Frank J.W.M. Dankers, Tanja Alderliesten, Peter A.N. Bosman

Real-world problems are often multi-objective with decision-makers unable to specify a prior which trade-off between the conflicting objectives is preferable. A learning approach to estimate the Pareto front by maximizing the dominated hypervolume (HV) of the average loss vectors corresponding to a set of learners' their HV maximizing gradients. Consequently, the learners get trained according to different trade-offs on the Pareto front, which otherwise is not guaranteed that the set of learners are indeed well-spread on the Pareto front. Further, the outputs corresponding to validation samples are also found to closely follow the true Pareto front.

Comments: T.M.D. and M.G. contributed equally

Subjects: Machine Learning (cs.LG)

Cite as: arXiv:2102.04523 [cs.LG]



(or arXiv:2102.04523v1 [cs.LG] for this version)




- ▶ The full method
- ▶ Link to PyTorch code
- ▶ Comparison to existing methods
- ▶ Why one should not learn on average losses.

<https://arxiv.org/abs/2102.04523>

# Questions?



UHV optimization



MO learning

# Appendix


## Existing methods

We compared to:

- ▶ Linear scalarization

$$\text{minimize} \quad a_i L_1(\theta_i, s_k) + (1-a_i) L_2(\theta_i, s_k)$$


- ▶ Pareto MTL (Lin et al., 2019)
- ▶ EPO (Mahapatra and Rajan, 2020)




All these methods require knowing the desired trade-offs **before** training.

# MO regression


Symmetric losses




Different losses




Different scales




## 3D MO regression




## 3D MO regression




# Organ segmentation




(a) Linear scalarization



(b) Pareto MTL




(c) EPO




(d) HV maximization



# Neural style transfer




(a)




(b)




# Figure 2



(a) Dominated subspaces



(b) HV gradients



(c) Domination-ranked fronts

# MO learning algorithm

---

**Algorithm 1** Training networks  $\Theta$  for Pareto front prediction by HV maximization of domination-ranked fronts


---

```
Initialize  $p$  networks  $\Theta = \{\theta_1, \dots, \theta_p\}$ 
for each batch  $\tilde{S}$  do
    for each network  $\theta_i$  do
        for each sample  $s_k \in \tilde{S}$  do
            Compute loss vector  $\mathcal{L}(\theta_i, s_k)$ 
        end for
    end for
    for each sample  $s_k \in \tilde{S}$  do
        Stack loss vectors  $\mathcal{L}(\theta_i, s_k)$  into  $\mathfrak{L}(\Theta, s_k)$ 
        Sort  $\mathfrak{L}(\Theta, s_k)$  into multiple fronts  $\mathfrak{L}(\Theta_l, s_k)$  by domination ranking (Section 3.2)
        for each front  $l$  do
            Compute loss weights  $\frac{\partial \text{HV}(\mathfrak{L}(\Theta_{q(i)}, s_k))}{\partial L_j(\theta_i, s_k)} \forall i, j$  using algorithm by
            Emmerich and Deutz (2014)
        end for
    end for
    for each network  $\theta_i$  do
        Backpropagate on joint loss from Equation (6)
    end for
    Update  $\Theta$  by stepping into gradient direction
end for
```

---


# Comparison on asymmetric fronts

Linear scalarization




(a)

Pareto MTL




(b)

EPO




(c)

HV maximization




(d)


MSE & L1-Norm




(e)



(f)




(g)




(h)


MSE & scaled MSE




(i)



(j)



(k)




(l)

# Learning per sample vs. on mean loss

Training per sample

$$\max \frac{1}{|S|} \sum_{k=1}^{|S|} \text{HV}(\mathcal{L}(\Theta, s_k))$$

(Dynamic loss (6))




(a)


Training on average losses

$$\max \text{HV}(\bar{\mathcal{L}}(\Theta, S))$$


(Dynamic loss (7))




(b)




(c)



(d)



Non-convex



(f)

# Counterexample: asymmetric front when learning on mean loss






Figure 4: An example of a learning problem with strictly convex Pareto fronts in which HV maximization of average losses does not result in well-distributed outputs on both samples' fronts. HV maximization of each sample's losses (left) and of average losses (right).


# Comparison: rescaling losses




(a)



(b)



(c)



(d)