

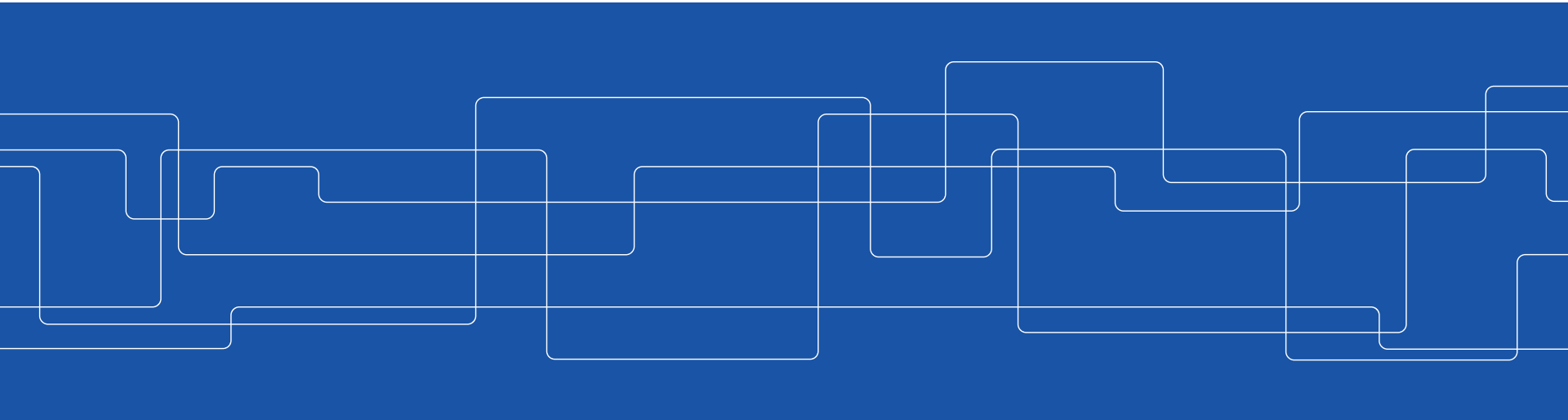


RPL CV / DL Reading Group

26 March 2019

Multitask Learning as Multiobjective Optimization

Ozan Sener and Vladlen Koltun





Outline

- Motivation
- Related Work
- M
- Model definition
- Training details
- Results
- MTL challenges analysis
- Conclusions & Future Work



Multi-task learning

Definition: jointly learn T tasks, sharing inductive bias across them designing a parameterized hypothesis that shares some parameters across tasks.

Strategies:

- Soft-sharing: All parameters specific to each task, jointly constrained.
- Hard-sharing: Part of the parameters fully-shared between tasks.



Multi-task learning

Definition: jointly learn T tasks, sharing inductive bias across them designing a parameterized hypothesis that shares some parameters across tasks.

Strategies:

- Soft-sharing: All parameters specific to each task, jointly constrained.
- **Hard-sharing: Part of the parameters fully-shared between tasks.**



Multi-task learning

Definition: jointly learn T tasks, sharing inductive bias across them designing a parameterized hypothesis that shares some parameters across tasks.

Strategies:

- Soft-sharing: All parameters specific to each task, jointly constrained.
- **Hard-sharing: Part of the parameters fully-shared between tasks.**
 - *Using deep neural networks as model*

Multi-task learning

Usually,

weighted sum of empirical risk for each task

$$\min_{\theta^{sh}, \theta^1, \dots, \theta^T} \sum_{t=1}^T c^t \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$$

- ... cannot handle competing tasks
- ... c^t can be static or dynamic
- ... uniform weights, found as hyperparameter, heuristics



Multi-task learning

Uncertainty weighting (Kendall et al. 2018)

1. Predict heteroscedastic uncertainty model as mean \hat{y}_t and variance σ_t^2 for each task t as new model output
2. Weight loss \mathcal{L}_t with $1 / 2\sigma_t^2$

GradNorm (Chen et al. 2018)

$$G_W^{(i)}(t) = \|\nabla_W w_i(t) L_i(t)\|_2$$

1. For each task, compute gradient wrt selected layer and its norm.
2. Compute average gradient norm $\bar{G}_W(t)$
3. Compute rel. training speed as loss / avg. loss. $r_i(t) = \tilde{L}_i(t) / E_{\text{task}}[\tilde{L}_i(t)]$
4. Compute loss to learn loss weights

$$L_{\text{grad}}(t; w_i(t)) = \sum_i \left| G_W^{(i)}(t) - \bar{G}_W(t) \times [r_i(t)]^\alpha \right|_1$$

5. Update loss weights, then model parameters

Multi-task learning as multi-objective opt.

Instead,

MTL as multi-objective optimization, optimizing set of possibly contrasting objectives.....

$$\min_{\theta^{sh}, \theta^1, \dots, \theta^T} \mathbf{L}(\theta^{sh}, \theta^1, \dots, \theta^T)$$

New goal: Find Pareto optimal solution, not dominated by any other solution.

A solution θ dominates a solution $\bar{\theta}$ if $\hat{\mathcal{L}}^t(\theta^{sh}, \theta^t) \leq \hat{\mathcal{L}}^t(\bar{\theta}^{sh}, \bar{\theta}^t)$ for all tasks t and $\mathbf{L}(\theta^{sh}, \theta^1, \dots, \theta^T) \neq \mathbf{L}(\bar{\theta}^{sh}, \bar{\theta}^1, \dots, \bar{\theta}^T)$.



Multi-task learning as multi-objective opt.

Focus on gradient-based multi-objective optimization...

MGDA - Multiple Gradient Descent Algorithm

...well-suited for multitask deep networks trained with stochastic gradient descent. But two issues need solving:

1. **Does not scale** to high-dimensional gradients
2. Requires separate computation of gradients for each task, i.e. one **backward pass per task**



MGDA - Multiple Gradient Descent Algorithm

Use *Karush-Kuhn-Tucker (KKT)* conditions to find common descent direction of shared parameters for all objectives, necessary for optimality.

- There exist $\alpha^1, \dots, \alpha^T \geq 0$ such that $\sum_{t=1}^T \alpha^t = 1$ and $\sum_{t=1}^T \alpha^t \nabla_{\theta^{sh}} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t) = 0$
- For all tasks t , $\nabla_{\theta^t} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t) = 0$

Solution satisfying these is Pareto stationary but not necessarily Pareto optimal.

MGDA - Multiple Gradient Descent Algorithm

$$\min_{\alpha^1, \dots, \alpha^T} \left\{ \left\| \sum_{t=1}^T \alpha^t \nabla_{\theta^{sh}} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t) \right\|_2^2 \mid \sum_{t=1}^T \alpha^t = 1, \alpha^t \geq 0 \quad \forall t \right\}$$

- Solution is zero:
 - Pareto stationary
- Non-zero:
 - Gives a descent direction improving all objectives

Equivalent to finding a minimum-norm solution in the convex hull of the set of solutions.

MGDA - Multiple Gradient Descent Algorithm

Case for 2 tasks has analytical solution:

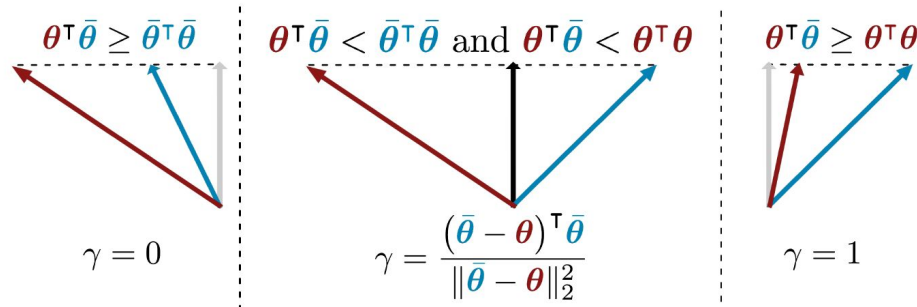


Figure 1: Visualisation of the min-norm point in the convex hull of two points ($\min_{\gamma \in [0,1]} \|\gamma\theta + (1-\gamma)\bar{\theta}\|_2^2$). As the geometry suggests, the solution is either an edge case or a perpendicular vector.

Algorithm 1

$$\min_{\gamma \in [0,1]} \|\gamma\theta + (1-\gamma)\bar{\theta}\|_2^2$$

- 1: **if** $\theta^\top \bar{\theta} \geq \bar{\theta}^\top \bar{\theta}$ **then**
 - 2: $\gamma = 1$
 - 3: **else if** $\theta^\top \bar{\theta} \geq \theta^\top \theta$ **then**
 - 4: $\gamma = 0$
 - 5: **else**
 - 6: $\gamma = \frac{(\bar{\theta} - \theta)^\top \bar{\theta}}{\|\bar{\theta} - \theta\|_2^2}$
 - 7: **end if**
-

MGDA - Multiple Gradient Descent Algorithm

Use 2D case as subroutine for line search in Frank-Wolfe optimizer.

Algorithm 2 Update Equations for MTL

```

1: for  $t = 1$  to  $T$  do
2:    $\theta^t = \theta^t - \eta \nabla_{\theta^t} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$ 
3: end for
4:  $\alpha^1, \dots, \alpha^T = \text{FRANKWOLFESOLVER}(\theta)$ 
5:  $\theta^{sh} = \theta^{sh} - \eta \sum_{t=1}^T \alpha^t \nabla_{\theta^{sh}} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$ 

6: procedure  $\text{FRANKWOLFESOLVER}(\theta)$ 
7:   Initialize  $\alpha = (\alpha^1, \dots, \alpha^T) = (\frac{1}{T}, \dots, \frac{1}{T})$ 
8:   Precompute  $\mathbf{M}$  st.  $\mathbf{M}_{i,j} = (\nabla_{\theta^{sh}} \hat{\mathcal{L}}^i(\theta^{sh}, \theta^i))^\top (\nabla_{\theta^{sh}} \hat{\mathcal{L}}^j(\theta^{sh}, \theta^j))$ 
9:   repeat
10:     $\hat{t} = \arg \max_r \sum_t \alpha^t \mathbf{M}_{rt}$ 
11:     $\hat{\gamma} = \arg \min_{\gamma} ((1 - \gamma)\alpha + \gamma e_{\hat{t}})^\top \mathbf{M} ((1 - \gamma)\alpha + \gamma e_{\hat{t}})$ 
12:     $\alpha = (1 - \hat{\gamma})\alpha + \hat{\gamma} e_{\hat{t}}$ 
13:  until  $\hat{\gamma} \sim 0$  or Number of Iterations Limit
14:  return  $\alpha^1, \dots, \alpha^T$ 
15: end procedure

```

▷ Gradient descent on task-specific parameters

▷ Solve (3) to find a common descent direction

▷ Gradient descent on shared parameters

▷ Using Algorithm 1

MGDA - Multiple Gradient Descent Algorithm

Use 2D case as subroutine for line search in Frank-Wolfe optimizer.

Algorithm 2 Update Equations for MTL

```

1: for  $t = 1$  to  $T$  do
2:    $\theta^t = \theta^t - \eta \nabla_{\theta^t} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$            ▷ Gradient descent on task-specific parameters
3: end for
4:  $\alpha^1, \dots, \alpha^T = \text{FRANKWOLFESOLVER}(\theta)$            ▷ Solve (3) to find a common descent direction
5:  $\theta^{sh} = \theta^{sh} - \eta \sum_{t=1}^T \alpha^t \nabla_{\theta^{sh}} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$    ▷ Gradient descent on shared parameters

6: procedure FRANKWOLFESOLVER( $\theta$ )
7:   Initialize  $\alpha = (\alpha^1, \dots, \alpha^T) = (\frac{1}{T}, \dots, \frac{1}{T})$ 
8:   Precompute  $\mathbf{M}$  st.  $\mathbf{M}_{i,j} = (\nabla_{\theta^{sh}} \hat{\mathcal{L}}^i(\theta^{sh}, \theta^i))^\top (\nabla_{\theta^{sh}} \hat{\mathcal{L}}^j(\theta^{sh}, \theta^j))$ 
9:   repeat
10:     $\hat{t} = \arg \max_r \sum_t \alpha^t \mathbf{M}_{rt}$ 
11:     $\hat{\gamma} = \arg \min_{\gamma} ((1 - \gamma)\alpha + \gamma e_{\hat{t}})^\top \mathbf{M} ((1 - \gamma)\alpha + \gamma e_{\hat{t}})$            ▷ Using Algorithm 1
12:     $\alpha = (1 - \hat{\gamma})\alpha + \hat{\gamma} e_{\hat{t}}$ 
13:  until  $\hat{\gamma} \sim 0$  or Number of Iterations Limit
14:  return  $\alpha^1, \dots, \alpha^T$ 
15: end procedure

```

MGDA - Multiple Gradient Descent Algorithm

Use 2D case as subroutine for line search in Frank-Wolfe optimizer.

Algorithm 2 Update Equations for MTL

```

1: for  $t = 1$  to  $T$  do
2:    $\theta^t = \theta^t - \eta \nabla_{\theta^t} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$            ▷ Gradient descent on task-specific parameters
3: end for
4:  $\alpha^1, \dots, \alpha^T = \text{FRANKWOLFESOLVER}(\theta)$            ▷ Solve (3) to find a common descent direction
5:  $\theta^{sh} = \theta^{sh} - \eta \sum_{t=1}^T \alpha^t \nabla_{\theta^{sh}} \hat{\mathcal{L}}^t(\theta^{sh}, \theta^t)$    ▷ Gradient descent on shared parameters

6: procedure  $\text{FRANKWOLFESOLVER}(\theta)$ 
7:   Initialize  $\alpha = (\alpha^1, \dots, \alpha^T) = (\frac{1}{T}, \dots, \frac{1}{T})$ 
8:   Precompute  $\mathbf{M}$  st.  $\mathbf{M}_{i,j} = (\nabla_{\theta^{sh}} \hat{\mathcal{L}}^i(\theta^{sh}, \theta^i))^\top (\nabla_{\theta^{sh}} \hat{\mathcal{L}}^j(\theta^{sh}, \theta^j))$ 
9:   repeat
10:     $\hat{t} = \arg \max_r \sum_t \alpha^t \mathbf{M}_{rt}$ 
11:     $\hat{\gamma} = \arg \min_{\gamma} ((1 - \gamma)\alpha + \gamma e_{\hat{t}})^\top \mathbf{M} ((1 - \gamma)\alpha + \gamma e_{\hat{t}})$            ▷ Using Algorithm 1
12:     $\alpha = (1 - \hat{\gamma})\alpha + \hat{\gamma} e_{\hat{t}}$ 
13:  until  $\hat{\gamma} \sim 0$  or Number of Iterations Limit
14:  return  $\alpha^1, \dots, \alpha^T$ 
15: end procedure

```



MGDA - Multiple Gradient Descent Algorithm

Frank-Wolfe solver typically converges within a few iterations, negligible addition to training time.

But we still need T backward passes...

MGDA - Upper Bound

Upper bound of objective of min-norm point problem...

$$\left\| \sum_{t=1}^T \alpha^t \nabla_{\boldsymbol{\theta}^{sh}} \hat{\mathcal{L}}^t(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^t) \right\|_2^2 \leq \left\| \frac{\partial \mathbf{Z}}{\partial \boldsymbol{\theta}^{sh}} \right\|_2^2 \left\| \sum_{t=1}^T \alpha^t \nabla_{\mathbf{Z}} \hat{\mathcal{L}}^t(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^t) \right\|_2^2$$

Does not depend on the alphas

Can be computed with a single backward pass

*If $\frac{\partial \mathbf{Z}}{\partial \boldsymbol{\theta}^{sh}}$ is full-rank (tasks not linearly related),
optimizing UB is equivalent.*



MGDA - Upper Bound

The rest of the algorithm is exactly the same

(:



Experiments

Baselines for all experiments

1. Single-task
2. Uniform weights
3. Weights found through grid-search
4. Uncertainty weighting
5. GradNorm

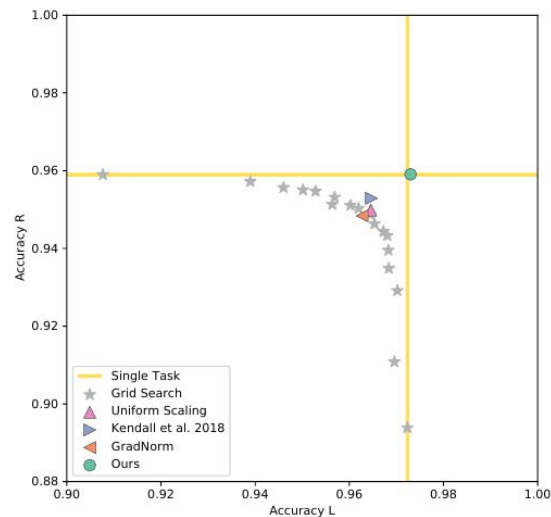
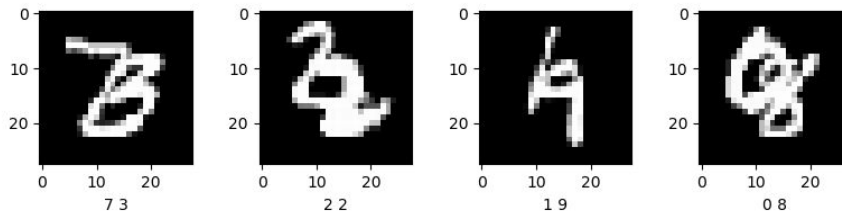
Experiments

MultiMNIST

2 tasks:

- Classify top-left digit “L”
- Classify bottom-right digit “R”

LeNet-based multi-task network



Experiments

CelebA

Eyeglasses



Wearing Hat

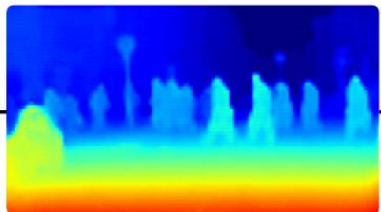
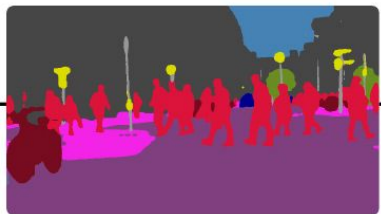


Multi-label classification,
each label is a binary class. task
40 tasks

ResNet-18 encoder, linear decoders.

	Average error
Single task	8.77
Uniform scaling	9.62
Kendall et al. 2018	9.53
GradNorm	8.44
Ours	8.25

Experiments

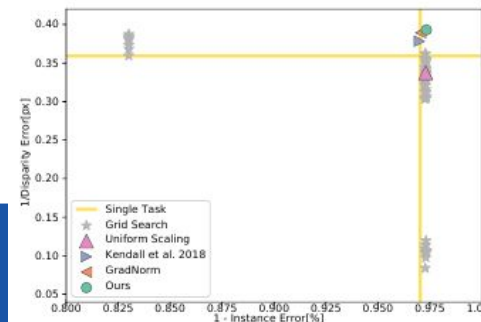
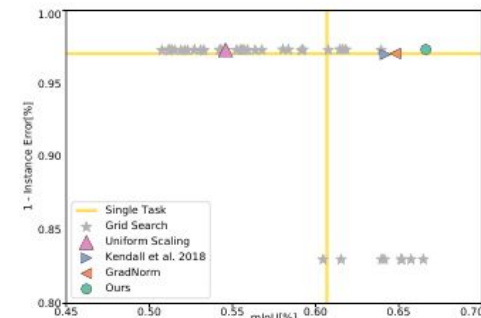
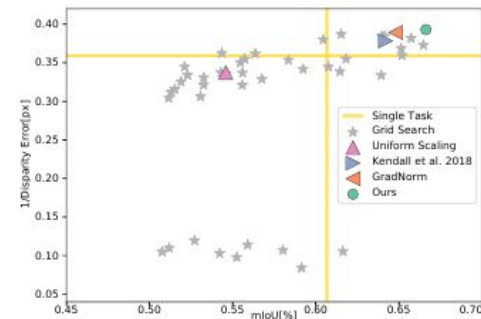


Cityscapes

3 scene-understanding tasks:

- Semantic segmentation
- Instance segmentation
- Depth estimation

ResNet-50 encoder,
pyramid pooling decoders



Experiments

Effect of upper bound approximation

Table 2: Effect of the MGDA-UB approximation. We report the final accuracies as well as training times for our method with and without the approximation.

	Scene understanding (3 tasks)				Multi-label (40 tasks)	
	Training time	Segmentation mIoU [%]	Instance error [px]	Disparity error [px]	Training time (hour)	Average error
Ours (w/o approx.)	38.6	66.13	10.28	2.59	429.9	8.33
Ours	23.3	66.63	10.25	2.54	16.1	8.25
	60 %				3.7%	



Experiments

Effect of upper bound approximation

And *surprisingly* also better accuracies...

...possibly due to solving problem (min norm point)
in lower-dimensional space
(shared representation instead of parameters)

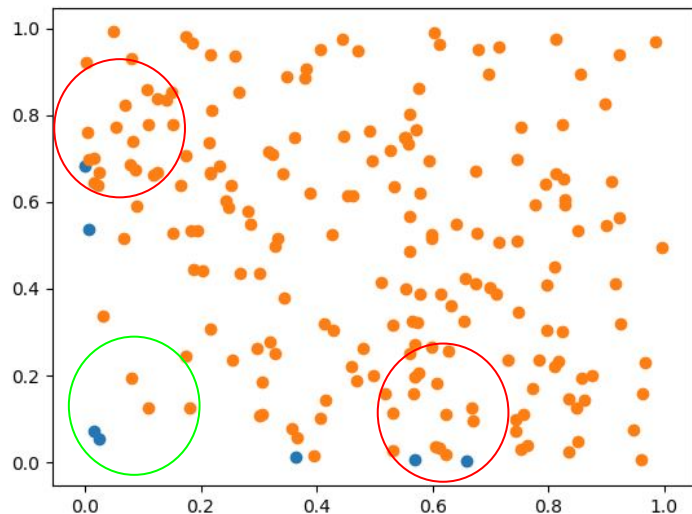


Conclusions

Applying multi-objective optimization to multi-task learning achieves better results than traditional approaches based on weighted sum of losses.

A method and an approximation with negligible computational overhead are proposed and evaluated on 3 different multi-task problems, showing it is effective on a wide range of scenarios.

Extra!



- MGDA tends to give a shortest path to the Pareto front
 - Not necessarily a balanced Pareto optimal solution

In practice, for cases in which gradient magnitudes differ a lot between tasks, this is important. Need to scale gradients:

- **By the loss**
- **By the L2 norm**
- **Other...**

