

# Multi-Task Learning Notes (II): Acting on Gradients

Su Jianlin

February 8, 2022

In "[Multi-Task Learning Notes \(I\): In the Name of Loss](#)", we initially explored the problem of multi-task learning from the perspective of loss functions. We eventually found that if we want the results to possess both scaling invariance and translation invariance, using the reciprocal of the gradient norm as the task weight is a relatively simple choice. We further analyzed that this design is equivalent to normalizing the gradient of each task individually before summing them. This means the "battlefield" of multi-task learning has shifted from the loss function to the gradients: it appears we are designing loss functions, but in reality, we are designing better gradients—the so-called "acting on gradients in the name of loss."

So, what are the criteria for a "better" gradient? How can we design better gradients? In this article, we will understand multi-task learning from the perspective of gradients and attempt to construct multi-task learning algorithms directly from the idea of designing gradients.

## 1 Overall Idea

We know that for single-task learning, the commonly used optimization method is gradient descent. How is it derived? Can the same logic be directly applied to multi-task learning? This is the question this section aims to answer.

### 1.1 Descent Direction

In fact, we answered the first question in "[Optimization Algorithms from a Dynamics Perspective \(III\): A More Holistic View](#)". Suppose the loss function is  $\mathcal{L}$  and the current parameter is  $\theta$ . We want to design a parameter increment  $\Delta\theta$  that makes the loss function smaller, i.e.,  $\mathcal{L}(\theta + \Delta\theta) < \mathcal{L}(\theta)$ . To this end, we consider the first-order Taylor expansion:

$$\mathcal{L}(\theta + \Delta\theta) \approx \mathcal{L}(\theta) + \langle \nabla_{\theta} \mathcal{L}, \Delta\theta \rangle \quad (1)$$

Assuming the accuracy of this approximation is sufficient,  $\mathcal{L}(\theta + \Delta\theta) < \mathcal{L}(\theta)$  implies  $\langle \nabla_{\theta} \mathcal{L}, \Delta\theta \rangle < 0$ . This means the angle between the update amount and the gradient must be at least greater than 90 degrees. The most natural choice is:

$$\Delta\theta = -\eta \nabla_{\theta} \mathcal{L} \quad (2)$$

This is gradient descent, where the update is taken in the opposite direction of the gradient, and  $\eta > 0$  is the learning rate.

### 1.2 No Exceptions

Returning to multi-task learning, if we assume each task is equally important, we can interpret this assumption as requiring  $\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n$  to all decrease or remain unchanged at each update step. If, after the parameters reach  $\theta^*$ , any further change leads to an increase in at least one  $\mathcal{L}_i$ , then  $\theta^*$  is said to be a Pareto Optimal solution. Simply put, Pareto optimality means

we cannot improve one task by sacrificing another; it means there is no "internal competition" (involution) between tasks.

Assuming the approximation (1) still holds, seeking Pareto optimality means we need to find  $\Delta\theta$  such that:

$$\begin{cases} \langle \nabla_{\theta} \mathcal{L}_1, \Delta\theta \rangle \leq 0 \\ \langle \nabla_{\theta} \mathcal{L}_2, \Delta\theta \rangle \leq 0 \\ \vdots \\ \langle \nabla_{\theta} \mathcal{L}_n, \Delta\theta \rangle \leq 0 \end{cases} \quad (3)$$

Note that a trivial solution  $\Delta\theta = \mathbf{0}$  exists, so the feasible region of the above system of inequalities is certainly non-empty. We are primarily concerned with whether non-zero solutions exist in the feasible region. If they do, we find one to use as the update direction; if not, we may have reached Pareto optimality (necessary but not sufficient), and we call this state a Pareto Stationary point.

## 2 Solving Algorithms

For convenience, we denote  $\mathbf{g}_i = \nabla_{\theta} \mathcal{L}_i$ . We seek a vector  $\mathbf{u}$  such that  $\langle \mathbf{g}_i, \mathbf{u} \rangle \geq 0$  for all  $i$ . Then we can take  $\Delta\theta = -\eta\mathbf{u}$  as the update amount, similar to single-task gradient descent. If there are only two tasks, it can be verified that  $\mathbf{u} = \mathbf{g}_1/\|\mathbf{g}_1\| + \mathbf{g}_2/\|\mathbf{g}_2\|$  automatically satisfies  $\langle \mathbf{g}_1, \mathbf{u} \rangle \geq 0$  and  $\langle \mathbf{g}_2, \mathbf{u} \rangle \geq 0$ . That is to say, in dual-task learning, the previously mentioned gradient normalization can reach a Pareto stationary point.

When the number of tasks is greater than 2, the problem becomes more complex. Here we introduce two solving methods. The first is a derivation provided by the author, and the second is the "standard answer" given in ["Multi-Task Learning as Multi-Objective Optimization"](#).

### 2.1 Problem Transformation

First, we further transform the problem. Note that:

$$\forall i, \langle \mathbf{g}_i, \mathbf{u} \rangle \geq 0 \quad \Leftrightarrow \quad \min_i \langle \mathbf{g}_i, \mathbf{u} \rangle \geq 0 \quad (4)$$

Therefore, we only need to maximize the minimum  $\langle \mathbf{g}_i, \mathbf{u} \rangle$  as much as possible to find the ideal  $\mathbf{u}$ . The problem becomes:

$$\max_{\mathbf{u}} \min_i \langle \mathbf{g}_i, \mathbf{u} \rangle \quad (5)$$

However, this is somewhat dangerous because if a non-zero  $\mathbf{u}$  exists such that  $\min_i \langle \mathbf{g}_i, \mathbf{u} \rangle > 0$ , then letting the norm of  $\mathbf{u}$  tend to infinity would cause the maximum value to tend to infinity. For stability, we add a regularization term:

$$\max_{\mathbf{u}} \min_i \langle \mathbf{g}_i, \mathbf{u} \rangle - \frac{1}{2} \|\mathbf{u}\|^2 \quad (6)$$

In this way, a  $\mathbf{u}$  with an infinite norm cannot be the optimal solution. Note that substituting  $\mathbf{u} = \mathbf{0}$  yields  $\min_i \langle \mathbf{g}_i, \mathbf{u} \rangle - \frac{1}{2} \|\mathbf{u}\|^2 = 0$ . Thus, if the optimal solution for max is  $\mathbf{u}^*$ , it must satisfy:

$$\min_i \langle \mathbf{g}_i, \mathbf{u}^* \rangle - \frac{1}{2} \|\mathbf{u}^*\|^2 \geq 0 \quad \Leftrightarrow \quad \min_i \langle \mathbf{g}_i, \mathbf{u}^* \rangle \geq \frac{1}{2} \|\mathbf{u}^*\|^2 \geq 0 \quad (7)$$

So the solution to problem (6) must satisfy condition (4), and if it is a non-zero solution, its opposite direction must be a direction that decreases the loss of all tasks.

## 2.2 Smooth Approximation

Now we introduce the first solution for problem (6). It assumes the reader, like the author, is unfamiliar with solving **min-max problems**. We can replace the min in the first step with a smooth approximation (refer to "[Seeking a Smooth Maximum Function](#)"):

$$\min_i \langle \mathbf{g}_i, \mathbf{u} \rangle \approx -\frac{1}{\lambda} \log \sum_i e^{-\lambda \langle \mathbf{g}_i, \mathbf{u} \rangle} \quad (\text{for sufficiently large } \lambda) \quad (8)$$

Thus, we can first solve:

$$\max_{\mathbf{u}} -\frac{1}{\lambda} \log \sum_i e^{-\lambda \langle \mathbf{g}_i, \mathbf{u} \rangle} - \frac{1}{2} \|\mathbf{u}\|^2 \quad (9)$$

And then let  $\lambda \rightarrow \infty$ . This transforms the problem into an unconstrained maximization of a single function. Taking the gradient and setting it to zero gives:

$$\frac{\sum_i e^{-\lambda \langle \mathbf{g}_i, \mathbf{u} \rangle} \mathbf{g}_i}{\sum_i e^{-\lambda \langle \mathbf{g}_i, \mathbf{u} \rangle}} = \mathbf{u} \quad (10)$$

Assuming the differences between various  $\langle \mathbf{g}_i, \mathbf{u} \rangle$  are larger than the order of  $\mathcal{O}(1/\lambda)$ , then as  $\lambda \rightarrow \infty$ , the above equation is effectively:

$$\mathbf{u} = \mathbf{g}_\tau, \quad \tau = \underset{i}{\operatorname{argmin}} \langle \mathbf{g}_i, \mathbf{u} \rangle \quad (11)$$

However, if we iterate directly according to  $\mathbf{u}^{(k+1)} = \mathbf{g}_\tau, \tau = \underset{i}{\operatorname{argmin}} \langle \mathbf{g}_i, \mathbf{u}^{(k)} \rangle$ , it will likely oscillate. This is because it asks us to find the  $\mathbf{g}_i$  that minimizes  $\langle \mathbf{g}_i, \mathbf{u}^{(k)} \rangle$  and use it as  $\mathbf{u}^{(k+1)}$ . If  $\mathbf{u}^{(k+1)} = \mathbf{g}_{i^*}$ , then in the next step, the  $\mathbf{g}_i$  that minimizes  $\langle \mathbf{g}_i, \mathbf{u}^{(k+1)} \rangle = \langle \mathbf{g}_i, \mathbf{g}_{i^*} \rangle$  is very likely no longer  $\mathbf{g}_{i^*}$ ; in fact,  $\mathbf{g}_{i^*}$  might be the largest one.

Intuitively, although the algorithm oscillates, it should oscillate around the optimal point  $\mathbf{u}^*$ . Therefore, if we average all results during the oscillation process, we should obtain the optimal point. This means the iteration format converging to the optimal point is:

$$\mathbf{u}^{(k+1)} = \frac{k\mathbf{u}^{(k)} + \mathbf{g}_\tau}{k+1}, \quad \tau = \underset{i}{\operatorname{argmin}} \langle \mathbf{g}_i, \mathbf{u}^{(k)} \rangle \quad (12)$$

Note that since each added term is some  $\mathbf{g}_i$ , the final  $\mathbf{u}^*$  must be a weighted average of the  $\mathbf{g}_i$ 's. That is, there exist  $\alpha_1, \alpha_2, \dots, \alpha_n \geq 0$  and  $\alpha_1 + \alpha_2 + \dots + \alpha_n = 1$  such that:

$$\mathbf{u}^* = \sum_i \alpha_i \mathbf{g}_i \quad (13)$$

We can also interpret  $\alpha_1, \alpha_2, \dots, \alpha_n$  as the current optimal weight distribution for each  $\mathcal{L}_i$ .

## 2.3 Dual Problem

The advantage of the smooth approximation technique is that it is simple and intuitive, requiring little background in optimization algorithms. However, it is ultimately a "non-mainstream" approach with some lack of rigor (though the result is correct). Next, we introduce the "standard answer" based on the idea of duality.

First, define  $\mathbb{P}^n$  as the set of all  $n$ -dimensional discrete distributions:

$$\mathbb{P}^n = \left\{ (\alpha_1, \alpha_2, \dots, \alpha_n) \mid \alpha_1, \alpha_2, \dots, \alpha_n \geq 0, \sum_i \alpha_i = 1 \right\} \quad (14)$$

Then it is easy to verify:

$$\min_i \langle \mathbf{g}_i, \mathbf{u} \rangle = \min_{\alpha \in \mathbb{P}^n} \langle \tilde{\mathbf{g}}(\alpha), \mathbf{u} \rangle, \quad \tilde{\mathbf{g}}(\alpha) = \sum_i \alpha_i \mathbf{g}_i \quad (15)$$

Thus, problem (6) is equivalent to:

$$\max_{\mathbf{u}} \min_{\alpha \in \mathbb{P}^n} \langle \tilde{\mathbf{g}}(\alpha), \mathbf{u} \rangle - \frac{1}{2} \|\mathbf{u}\|^2 \quad (16)$$

The above function is concave with respect to  $\mathbf{u}$  and convex with respect to  $\alpha$ , and the feasible regions for  $\mathbf{u}$  and  $\alpha$  are both convex sets. Therefore, according to von Neumann's [Minimax Theorem](#), the min and max in (16) can be swapped:

$$\min_{\alpha \in \mathbb{P}^n} \max_{\mathbf{u}} \langle \tilde{\mathbf{g}}(\alpha), \mathbf{u} \rangle - \frac{1}{2} \|\mathbf{u}\|^2 = \min_{\alpha \in \mathbb{P}^n} \frac{1}{2} \|\tilde{\mathbf{g}}(\alpha)\|^2 \quad (17)$$

The right side of the equation follows because the max part is just an unconstrained maximum problem for a quadratic function, which yields  $\mathbf{u}^* = \tilde{\mathbf{g}}(\alpha)$ . Finally, only the min remains, and the problem becomes finding a weighted average of  $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_n$  that minimizes its norm.

When  $n = 2$ , the solution is relatively simple, equivalent to constructing the altitude of a triangle, as shown below:

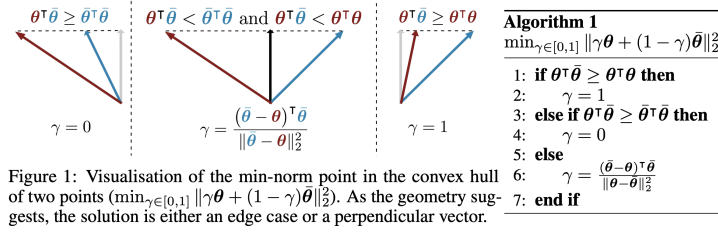


Figure 1: Visualisation of the min-norm point in the convex hull of two points ( $\min_{\gamma \in [0,1]} \|\gamma \theta + (1-\gamma) \bar{\theta}\|_2^2$ ). As the geometry suggests, the solution is either an edge case or a perpendicular vector.

Figure 1: Solving algorithm and geometric meaning when  $n = 2$

When  $n > 2$ , we can use the [Frank-Wolfe algorithm](#) to transform it into multiple  $n = 2$  cases for iteration. The Frank-Wolfe algorithm can be understood as a constrained gradient descent algorithm suitable for cases where the parameter's feasible region is a convex set. Explaining it in detail would take too much space, so I won't elaborate here; readers are encouraged to find materials to study it. Briefly, the Frank-Wolfe algorithm first linearizes the objective to find the next update direction  $e_\tau$ , where  $\tau = \operatorname{argmin}_i \langle \mathbf{g}_i, \tilde{\mathbf{g}}(\alpha) \rangle$  and  $e_\tau$  is a one-hot vector with 1 at position  $\tau$ . Then it performs an interpolation search between  $\alpha$  and  $e_\tau$  to find the optimal result. The iteration process is:

$$\begin{cases} \tau = \operatorname{argmin}_i \langle \mathbf{g}_i, \tilde{\mathbf{g}}(\alpha^{(k)}) \rangle \\ \gamma = \operatorname{argmin}_\gamma \|\tilde{\mathbf{g}}((1-\gamma)\alpha^{(k)} + \gamma e_\tau)\|^2 = \operatorname{argmin}_\gamma \|(1-\gamma)\tilde{\mathbf{g}}(\alpha^{(k)}) + \gamma \mathbf{g}_\tau\|^2 \\ \alpha^{(k+1)} = (1-\gamma)\alpha^{(k)} + \gamma e_\tau \end{cases} \quad (18)$$

The solution for  $\gamma$  is exactly the  $n = 2$  special case, which can be solved using the algorithm in the screenshot above. If  $\gamma$  is not obtained through search but fixed at  $1/(k+1)$ , the result is equivalent to (12), which is a simplified version of the Frank-Wolfe algorithm. In other words, the result obtained through smooth approximation is equivalent to the result of the simplified Frank-Wolfe algorithm.

## 2.4 De-constraint

In fact, for the solution of problem (17), we could theoretically solve it directly via gradient descent by removing constraints. For example, we can set parameters  $\beta_1, \beta_2, \dots, \beta_n \in \mathbb{R}$  and:

$$\alpha_i = \frac{e^{\beta_i}}{Z}, \quad Z = \sum_i e^{\beta_i} \quad (19)$$

Then it can be transformed into:

$$\min_{\beta} \frac{1}{2Z^2} \left\| \sum_i e^{\beta_i} \mathbf{g}_i \right\|^2 \quad (20)$$

This is an unconstrained optimization problem that can be solved with conventional gradient descent. However, for some reason, the author has not seen this approach used (perhaps to avoid tuning the learning rate?).

## 3 Some Techniques

In the previous section, we provided two schemes for finding the update direction of a Pareto stationary point. Both require us to perform additional multiple iterations to determine the weight of each task at every step of training before updating the model parameters. As one can imagine, the computational cost is quite high in practice, so we need some techniques to reduce it.

### 3.1 Gradient Inner Product

As seen, in both schemes, the key step is  $\operatorname{argmin}_i \langle \mathbf{g}_i, \tilde{\mathbf{g}}(\alpha) \rangle$ , which means we must traverse gradients to calculate inner products. However, in deep learning scenarios, model parameters are often very large, so the gradient is a very high-dimensional vector. Calculating an inner product at every iteration is computationally expensive. We can utilize the expansion:

$$\langle \mathbf{g}_i, \tilde{\mathbf{g}}(\alpha) \rangle = \left\langle \mathbf{g}_i, \sum_j \alpha_j \mathbf{g}_j \right\rangle = \sum_j \alpha_j \langle \mathbf{g}_i, \mathbf{g}_j \rangle \quad (21)$$

In each iteration, only  $\alpha$  changes. Therefore,  $\langle \mathbf{g}_i, \mathbf{g}_j \rangle$  only needs to be calculated once and stored for each training step, avoiding repeated high-dimensional vector inner product calculations.

### 3.2 Shared Encoding

However, when the model is large enough, it is difficult to calculate the gradient for each task separately and then perform iterative calculations. If we assume that the various tasks share the same encoder, we can further simplify the algorithm with an approximation.

Specifically, assume the batch size is  $b$ , and the encoder output for the  $j$ -th sample is  $\mathbf{h}_j$ . From the chain rule, we know:

$$\mathbf{g}_i = \nabla_{\theta} \mathcal{L}_i = \sum_j (\nabla_{\mathbf{h}_j} \mathcal{L}_i) (\nabla_{\theta} \mathbf{h}_j) = \underbrace{(\nabla_{\mathbf{h}_1} \mathcal{L}_i, \dots, \nabla_{\mathbf{h}_b} \mathcal{L}_i)}_{\nabla_{\mathbf{H}} \mathcal{L}_i} \underbrace{\begin{pmatrix} \nabla_{\theta} \mathbf{h}_1 \\ \vdots \\ \nabla_{\theta} \mathbf{h}_b \end{pmatrix}}_{\nabla_{\theta} \mathbf{H}} \quad (22)$$

Letting  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_b)$ , we get  $\mathbf{g}_i = (\nabla_{\mathbf{H}} \mathcal{L}_i) (\nabla_{\theta} \mathbf{H})$ . Using the matrix norm inequality:

$$\left\| \sum_i \alpha_i \mathbf{g}_i \right\|^2 = \left\| \sum_i \alpha_i (\nabla_{\mathbf{H}} \mathcal{L}_i) (\nabla_{\theta} \mathbf{H}) \right\|^2 \leq \left\| \sum_i \alpha_i \nabla_{\mathbf{H}} \mathcal{L}_i \right\|^2 \|\nabla_{\theta} \mathbf{H}\|^2 \quad (23)$$

Naturally, if we minimize  $\left\| \sum_i \alpha_i \nabla_{\mathbf{H}} \mathcal{L}_i \right\|^2$ , the computational cost will be significantly reduced because this only requires the gradient of the final output encoding vector, not the gradients of all parameters. The above formula tells us that minimizing  $\left\| \sum_i \alpha_i \nabla_{\mathbf{H}} \mathcal{L}_i \right\|^2$  is actually minimizing the upper bound of equation (17). Like many problems that are difficult to optimize directly, we hope that minimizing the upper bound will yield similar results.

However, while this upper bound is more efficient, it has limitations. It generally only applies to multi-task learning where each sample has multiple types of labels. It is not suitable for scenarios where training data for different tasks have no intersection (i.e., each task is labeled on different samples, and a single sample has only one type of label). In the latter case, the various  $\nabla_{\mathbf{H}} \mathcal{L}_i$  are mutually orthogonal, meaning there is no interaction between tasks, and the upper bound becomes too loose to be meaningful.

### 3.3 Error in Proof

The "standard answer" mentioned earlier and the optimization of the upper bound for shared encoders both come from the paper "[Multi-Task Learning as Multi-Objective Optimization](#)". The original paper then attempts to prove that when  $\nabla_{\boldsymbol{\theta}} \mathbf{H}$  is full rank, optimizing the upper bound can also find a Pareto stationary point. Unfortunately, the proof in the original paper is incorrect.

The proof is located in Appendix A of the original paper, where it uses an incorrect conclusion:

If  $\mathbf{M}$  is a symmetric positive definite matrix, then  $\mathbf{x}^\top \mathbf{y} \geq 0$  if and only if  $\mathbf{x}^\top \mathbf{M} \mathbf{y} \geq 0$ .

It is easy to provide a counterexample to prove this conclusion is wrong. For instance, let  $\mathbf{x} = \begin{pmatrix} 1 \\ -2 \end{pmatrix}$ ,  $\mathbf{y} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ ,  $\mathbf{M} = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}$ . Here  $\mathbf{x}^\top \mathbf{y} < 0$  but  $\mathbf{x}^\top \mathbf{M} \mathbf{y} > 0$ .

After reflection, the author believes the proof in the original paper is difficult to fix, meaning the original paper's conjecture does not hold. In other words, even if  $\nabla_{\boldsymbol{\theta}} \mathbf{H}$  is full rank, the update direction derived from the upper bound may not necessarily be a direction that prevents all task losses from rising, and thus may not find a Pareto stationary point. As for why the experimental results of the upper bound optimization in the original paper are good, it can only be said that the parameter space of deep learning models is so large that there is plenty of "room to maneuver," allowing the upper bound approximation to achieve decent results.

## 4 Summary

In this article, we understood multi-task learning from the perspective of gradients. From this viewpoint, the main task of multi-task learning is to find an update direction that is as opposite as possible to the gradient of each task, so that the loss of each task can decrease as much as possible, without sacrificing one task for the sake of another. This is the ideal state of no "internal competition" between tasks.

*Original Address:* <https://kexue.fm/archives/8896>