

# Task Adaptive Parameter Sharing for Multi-Task Learning

## [Methodology]

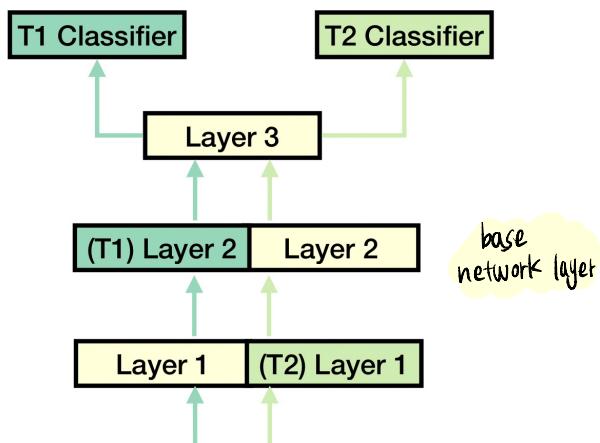
Given:

- Neural Networks w/  $L$  layers of weights
- $K$  target tasks:  $\{T_1, T_2, \dots, T_k\}$

Goal: select the minimal necessary subset of layers that needs to be tuned to achieve the best (or close to the best) performance.

Idea: To relax the combinatorial problem into a continuous one, which will ultimately give us a simple joint loss function to find both the optimal task-specific layers to tune and optimize parameters of those layers.

$\langle$  (Weight parametrization)  $\rangle$



(b) Our Method

$S_i$ : scoring parameter for layer  $i$

reparametrization of the weights for each layer

$$W_i = \bar{W}_i + I_T(S_i) \delta W_i$$

$\bar{W}_i$ : the shared weights of the pretrained network

$\delta W_i$ : a trainable parameter which describes a task-specific perturbation of the base network

$$I_T(S_i) = \begin{cases} 1 & \text{if } S_i \geq \tau \quad \leftarrow \text{task-specific layer} \\ & \quad (\text{i.e., make its parameters task-specific}) \\ 0 & \text{otherwise} \quad \leftarrow \text{the layer is the same as the base network and no new parameters are introduced.} \end{cases}$$

## <Joint Optimization>

Recast the initial combinatorial problem as a joint optimization problem over weight deltas and scores.

tries to optimize  
the task-specific parameters  
to achieve the best performance  
on the task

$$(\delta_W, S) = \underset{W, S}{\operatorname{argmin}} \mathcal{L}_D(W) + \frac{\lambda}{L} \sum_{i=1}^L |S_i|$$

$$S = (S_1, \dots, S_L)$$

$$\delta_W = (\delta_{W_1}, \dots, \delta_{W_L})$$

$$W = (W_1, \dots, W_L)$$

$\mathcal{L}_D(W)$ : the loss of model on dataset D

Sparsity inducing regularizer on S  
which penalizes large values of  $S_i$ ,  
encouraging layer sharing rather  
than introducing task-specific  
parameters.

## <Straight-through gradient estimation>

To make the problem learnable, we utilize a straight-through gradient estimator.  
i.e., modify the backward pass and redefine the gradient update as:

$$\nabla_{S_i} W_i = \delta_{W_i}$$

$$(i.e., \nabla_{S_i} W_i = \nabla_{S_i} \{ \bar{W}_i + S_i \delta_{W_i} \})$$

## <Joint MTL>

How to learn a pretrained representation that reduces the number of task-specific layers that need to be learned to obtain optimal performance:

### ■ Loss function:

$$(\bar{w}, \delta w_1, \dots, \delta w_K, s_1, \dots, s_K) =$$

$$\arg \min_{\bar{w}, \delta w^k, s} \sum_{k=1}^K (\mathcal{L}_{D_k}(\bar{w}, \delta w^k) + \frac{1}{L} \sum_{i=1}^L |s_i^k|)$$

- $k$ : # of tasks
- $\bar{w}$ : common weights
- $\delta w^k, s^k$ : task specific parameters

→ encourages learning common weights  $\bar{w}$  in such a way that the number of task specific parameters is minimized, due to the  $L1$  penalty on  $s$

## ⊕ Batch Normalization

## [Experiments]

