

## Lecture Plan

\* Overview of weakly supervised learning: a bridge between supervised and unsupervised learning

\* Multi-task learning

- Hard parameter sharing
- Soft parameter sharing

\* Transfer learning

- Fine tuning

---

So far in the class, we have been studying supervised learning. This makes several assumptions:

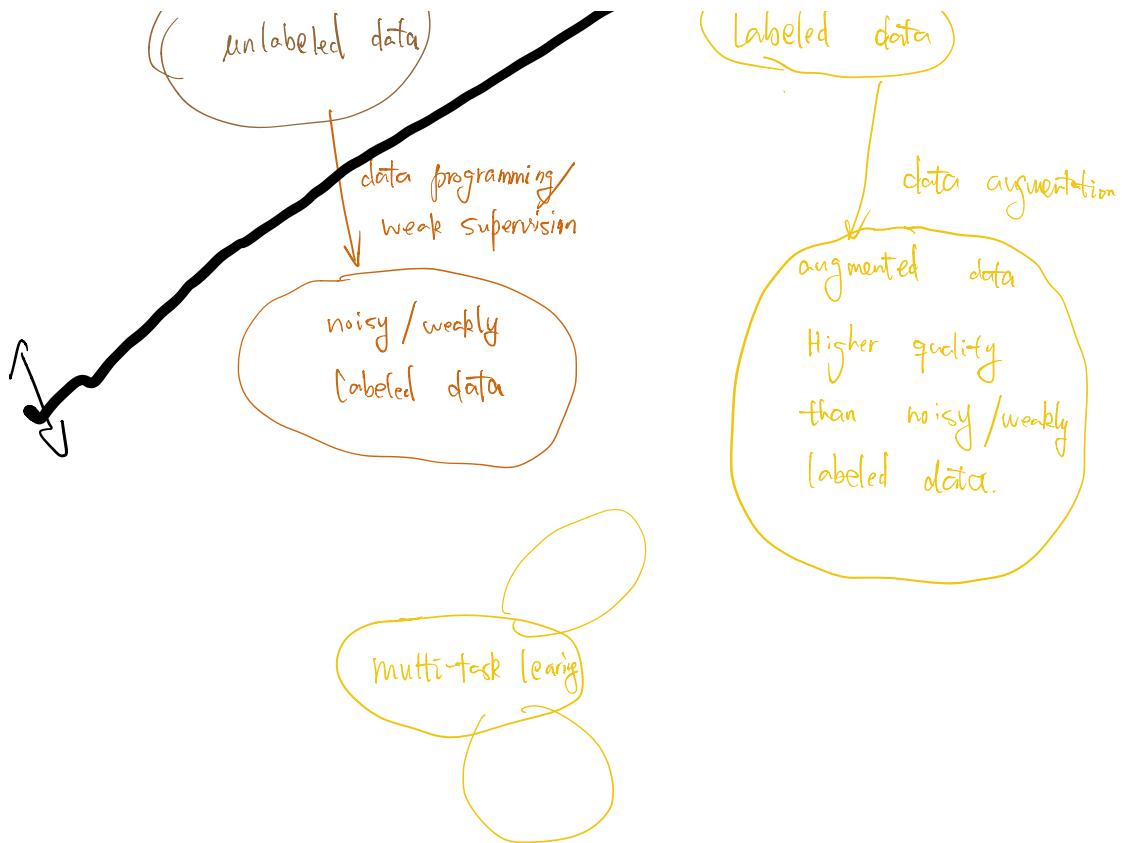
- ① \* We have lots of labeled data already (Neural net training requires lots of labeled data)
- ② \* The training data distribution and the test data distribution are the same.

The goal of weakly supervised learning is to address problem ①.

The goal of transfer learning or domain adaptation is to address problem ②.

A picture of weakly supervised learning





Crossing the boundary from supervised to unsupervised learning is a challenging task.

### — Self-supervised learning

Ex. Noisy student (Xie et al. CVPR 2020)

- ~~ Imagine that we only have a limited amount of labeled data.
- ~~ Bootstrap :
  - (i) we learn a neural net using the labeled data.
  - (ii) use the learned neural net to infer pseudo-labels over unlabeled data
  - (iii) re-train using the labeled data + pseudo-labeled data.

Ex. Masked next word prediction (BERT)

- Take the entire Wikipedia corpus (!?)
- Train a multi-layer (12/24) transformer that predicts the next word.

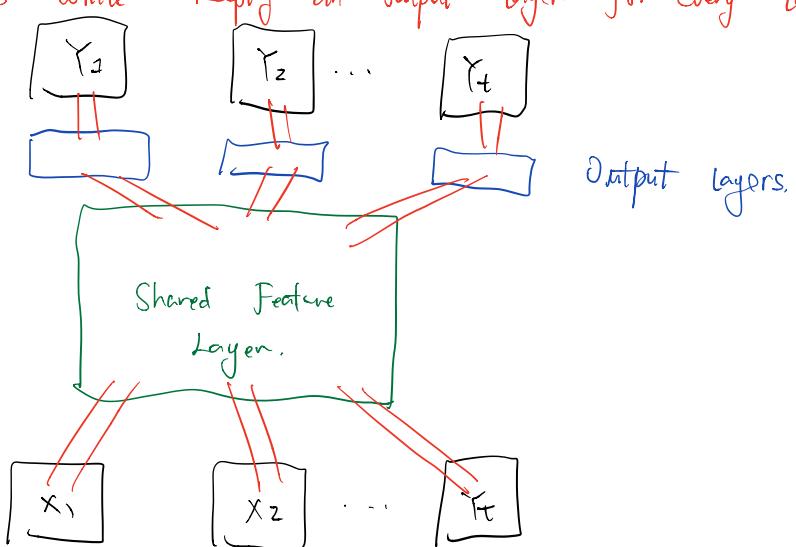
In both of these examples, the model is trained using zero or very little labeled data.

---

Multi-task learning : There are two approaches in neural nets

- Hard parameter sharing
- Soft parameter sharing

\* Hard parameter sharing is generally applied by sharing the feature layers between all tasks while keeping an output layer for every task.



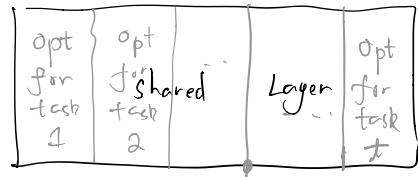
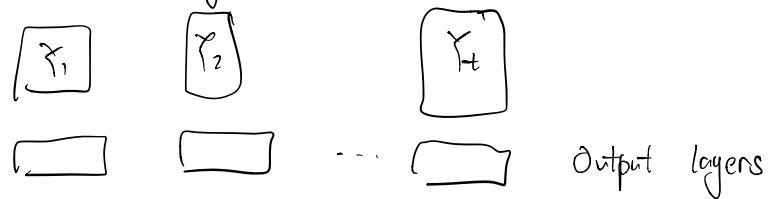
Ex. For text classification tasks, a commonly used shared feature layer is BERT.

Why does hard parameter sharing work?

- The first reason is regularization.

Claim. By restricting the shared layer's size, one can reduce overfitting.

Setting. Consider a two-layer neural network:



Observation. If the shared layer is too large, then it can fit the optimal solution of every single task.

On the other hand, by restricting the shared layer's size, we encourage information sharing!

- The second reason is adding inductive bias from multiple tasks.

Question: the inductive bias depends on datasets' properties such as sample sizes, covariance matrices, and so on. How are they related?

Claim. Hand parameter sharing helps by reducing variance, but hurts by increasing bias.

Ex. (Multi-label regression).

- We have  $t$  linear regression tasks.

$$\text{Task 1 : } x_1, y_1^{(1)} = x_1^T \cdot \beta^{(1)} + \varepsilon_1^{(1)} \rightarrow N(0, \sigma^2)$$

$$x_2, y_2^{(2)} = x_2^T \cdot \beta^{(2)} + \varepsilon_2^{(2)}$$

⋮

$$x_n, y_n^{(t)} = x_n^T \cdot \beta^{(t)} + \varepsilon_n^{(t)}$$

$\mathbb{E}[x_i x_i^T] = \Sigma$

In matrix notation, we write

$$Y^{(1)} = X \cdot \beta^{(1)} + \varepsilon^{(1)}$$

$$\text{Task 2 : } Y^{(2)} = X \cdot \beta^{(2)} + \varepsilon^{(2)}$$

⋮

$$\text{Task } t : Y^{(t)} = X \cdot \beta^{(t)} + \varepsilon^{(t)}$$

Recall: for single-task learning, we can solve a linear regression problem, and get that

$$\hat{\beta}_i = (X^T X)^{-1} X^T Y^{(i)}$$

and the test loss of  $\hat{\beta}_i$  is

$$L(\hat{\beta}_i) = t^2 \cdot \text{Tr} \left[ \Sigma \cdot (X^T X)^{-1} \right], \text{ where } \Sigma \text{ is}$$

the population covariance matrix.

Problem formulation for hard parameter sharing

$$f(A, B) = \sum_{i=1}^t \left\| \underbrace{x \cdot B \cdot A_i}_{\substack{\text{shared} \\ \text{layer}}} - \underbrace{y^{(i)}}_{\substack{\text{output} \\ \text{layer}}} \right\|^2$$

Let  $(\hat{A}, \hat{B})$  denote the global minimum of  $f(A, B)$ .

Hence the hard parameter sharing predictor for task  $i$  is:

$$\hat{\beta}_i^{\text{HPS}} = \hat{B} \cdot \hat{A}_i$$

Claim. The test loss of  $\hat{\beta}_i^{\text{HPS}}$  is equal to

$$L(B^* \cdot a_i^*) + b^2 \cdot \|a_i^*\|^2 \cdot \text{Tr}[\Sigma \cdot (X^T X)^{-1}]$$

where:  $B^* = [\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(t)}]$ .

$$A^* = \underset{\text{rank}(U) = r}{\text{argmin}} \langle UU^T, B^* \Sigma B^* \rangle$$

$a_i^*$  → i-th column of  $A^* A^{*^T}$

$L(B^* a_i^*)$  = test loss of  $B^* \cdot a_i^*$ .

Remark. This is a bias-variance decomposition of HPS!!

Bias:  $L(B^* \cdot a_i^*)$  does not depend on noise.

Variance:  $b^2 \cdot \|a_i^*\|^2 \cdot \text{Tr}[\Sigma \cdot (X^T X)^{-1}]$

Ex.  $\|a_i^*\|^2 \leq 1$

$$\leq b^2. \quad \text{Tr} \left[ \Sigma \cdot (X^T X)^{-1} \right]$$

because  $A^*$  is  
a projection  
matrix

Hence, the bias of HPS is always larger than STL,

the variance of HPS is always smaller than STL.

Proof sketch:  $f(A, B) = \sum_{i=1}^t \| X B A_i - A^{(i)} \|^2$

Step 1. Derive the global minimum of  $f(A, B)$ :

~ Using local optimality condition:  $\frac{\partial f(A, B)}{\partial B} = 0$ ,

we get  $B$

— Applying  $B$  to  $f(A, B)$ , we get  $A$ .

$\Rightarrow \hat{\beta}_i^{\text{HPS}}$  predictor

Step 2.  $E \left[ L(\hat{\beta}_i^{\text{HPS}}) \right]$

$X, \varepsilon^{(1)}, \varepsilon^{(2)}, \dots, \varepsilon^{(n)}$

$$= n \cdot \left\| \Sigma^{1/2} B^* \underbrace{\left( A^* (A A^T)^{-1} A - \text{Id} \right)}_F^2 + b^2 \right\|_F^2 \text{ (nt-pr).}$$

The minimizer of this equation  
is precisely  $A^*$ .

Geometric intuition. HPS essentially finds a rank- $r$  approximation of  
all the tasks!

Recap. In a simplified linear model, we have shown that

(i) HPS' inductive bias comes from reducing variance

(ii) also introduces bias

Explaining negative transfer.

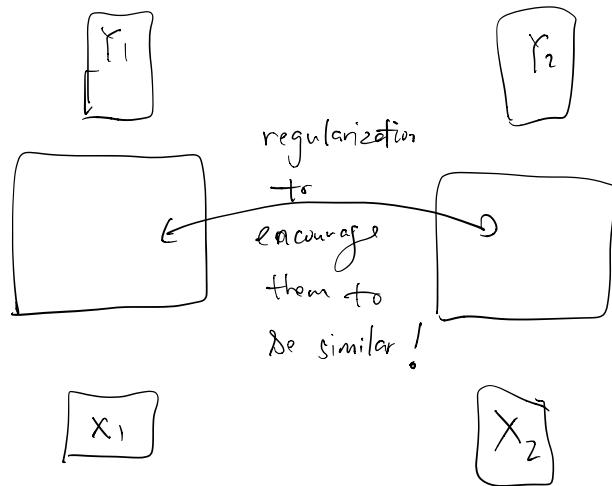
Research question: design optimization schedules for multi-task learning.

- Curriculum learning: start from easy examples, then go to more difficult examples!

Challenges: different sample sizes, covariates shifts, lots of room for improvement!

- Soft parameter sharing:

Ex.



- Transfer learning; In Th, you train a model on one task, and hope that the model can help learn the other task.

Fine-tuning :

- Train a model on the source task:  $\hat{f}^{(S)}$
- On the target task,
  - \* Initialize from  $\hat{f}^{(S)}$
  - & Run SGD
- Early stopping / Checkpointing is crucial in the use of fine-tuning (T5, google)

Many exciting work in this direction:

- Meta learning. Can a model learn to adapt to new environments quickly? What happens to a new environment with little labeled data?
- Few shot learning. Suppose we have a tiny amount of data from various tasks but lots of data from some other sources. How can we use those data to predict the "small data" tasks?