

# Transfer learning / Multi-task learning ideas in a nutshell ①

Traditional setting: given dataset  $D$ , learn the underlying distribution, then test it on test data coming from the SAME distribution  
(ML)

Transfer learning: I have source  $D_S$  and target  $D_T$ , where:  
(TL)  
① some relatedness between  $D_S$  and  $D_T$  exists  
② lot of info in  $D_S$ , little info available in  $D_T$

Multi-task learning: I have several sources  $D_1, \dots, D_K$ , and/or different tasks (regression, classification, clustering...)  
(MTL)

Note: TL/MTL do not ask for different dataset and different tasks at the same time.

## TL/MTL ideas

① use  $D_S, D_T$  to train model together:

1) naive approach: reweight loss on  $D_S$ :

$$\theta^* = \operatorname{argmin} L(\theta; D_S) + \alpha L(\theta; D_T) + R(\theta),$$

$\theta$  is model parameter,  $\alpha$  controls the weighting

$R(\theta)$  regularize the model

problem: how to choose  $\alpha$ ?

2) covariate shift: remember often we define  $\mathcal{L}(\theta; \mathcal{D}_S) = \mathbb{E}_{x \sim \mathcal{D}_S} [\mathcal{L}(\theta; x)]$  (same for  $\mathcal{L}(\theta; \mathcal{D}_T$ ), then we first "shift" the datapoints in  $\mathcal{D}_S$  to ~~be~~ be approx. distributed like  $\mathcal{D}_T$ , then train the model (see Shimodaira 2000 paper):

$$\theta^* = \operatorname{argmin} \mathcal{L}(\theta; \mathcal{D}_T) + \mathbb{E}_{x \sim \mathcal{D}_S} \left[ \frac{\mathcal{P}_T(x)}{\mathcal{P}_S(x)} \mathcal{L}(\theta; x) \right]$$

where  $\mathcal{P}_T, \mathcal{P}_S$  are (empirical) distributions of  $\mathcal{D}_T, \mathcal{D}_S$

(now we solve the problem of choosing  $\alpha$  in the naive way)

now new problem: how to estimate  $\frac{\mathcal{P}_T(x)}{\mathcal{P}_S(x)}$ ?

A) estimate  $\mathcal{P}_T, \mathcal{P}_S$  separately using whatever methods you like.

problem: estimator of  $\mathcal{P}_T$  is often not reliable.

B) directly obtain IS weights is not optimal (because of problems of IS), so maybe we can accept some biases?

some ideas: (kernel)

A) feature mean matching (see Gretton's book chapter):

define  $\beta(x) = \frac{\mathcal{P}_T(x)}{\mathcal{P}_S(x)}$  and feature embedding  $\Phi(x)$ ,

then solve the following problem:

$$\underset{\beta}{\text{minimize}} \quad \mathbb{E}_{x \sim \mathcal{P}_T} [\Phi(x)] - \mathbb{E}_{x \sim \mathcal{P}_S} [\beta(x) \Phi(x)]$$

$$\text{subject to } \beta(x) \geq 0 \text{ and } \mathbb{E}_{x \sim \mathcal{P}_S} [\beta(x)] = 1$$

now  $\|\cdot\|$  is of choose, or we can do it implicitly using kernel tricks!

B) train a classifier  $f(x)$  to distinguish data from  $\mathcal{D}_S$  or  $\mathcal{D}_T$ :  
 def  $p(z=1|x) = \frac{\mathcal{P}_T(x)}{\mathcal{P}_T(x) + \mathcal{P}_S(x)}$  (prob. of  $x$  from  $\mathcal{D}_T$ )

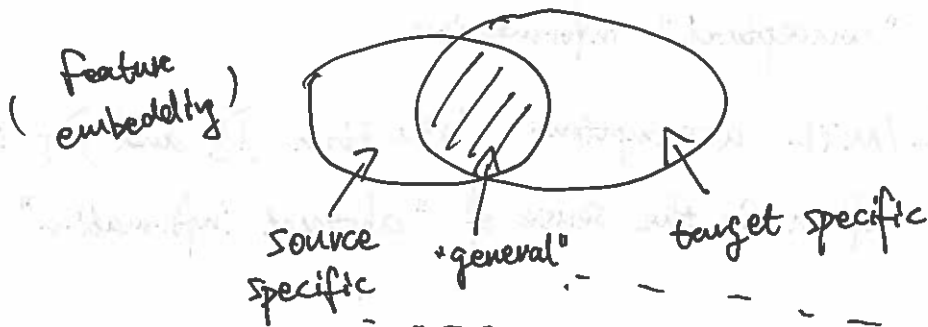
then  $\frac{P_T(x)}{P_S(x)} = \frac{P(z=1|x)}{P(z=-1|x)}$ , and if we define  $(3)$

$$P(z=1|x) = \frac{1}{1 + e^{-f(x)}} \Rightarrow \frac{P_T(x)}{P_S(x)} = e^{f(x)}$$

So the problem of covariate shift is solved in two steps:

- ① compute  $f(x)$  using maximum likelihood / MAP ...
- ② use  $e^{f(x)}$  as the weight to solve the original model training.

3) "bridging" via feature space (see Daumé II's EasyAdapt)



define two embeddings:

$$\text{for } x \sim \mathcal{D}_S: \Phi_S(x) = (\phi_S(x), \phi_g(x), 0)$$

$$\text{for } x \sim \mathcal{D}_T: \Phi_T(x) = (0, \phi_g(x), \phi_T(x))$$

then define model using input  $\Phi_S(x)$  and  $\Phi_T(x)$ :

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x \sim \mathcal{D}_S} [L(\theta; \Phi_S(x))] + \mathbb{E}_{x \sim \mathcal{D}_T} [L(\theta; \Phi_T(x))] + R(\theta)$$

example: logistic regression: model

$$F(x; \theta) = \begin{cases} \operatorname{sigmoid}(W_S^T \phi_S(x) + W_g^T \phi_g(x)), & x \sim \mathcal{D}_S \\ \operatorname{sigmoid}(W_g^T \phi_g(x) + W_T^T \phi_T(x)), & x \sim \mathcal{D}_T \end{cases}$$

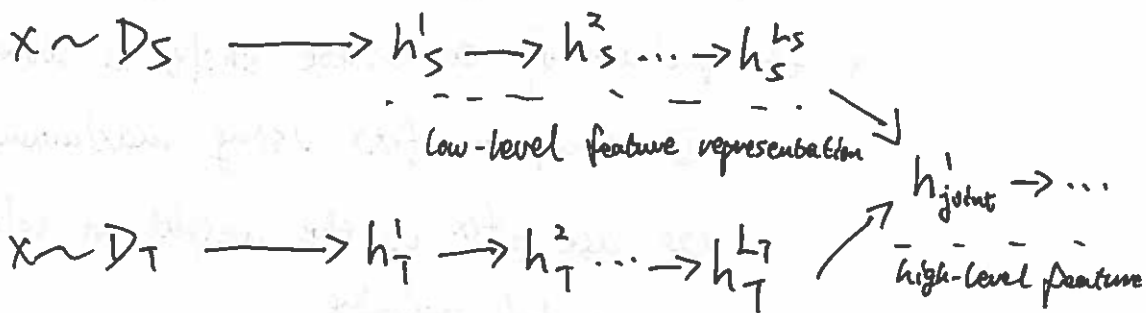
adding unlabelled data: we want the model to agree on unlabelled data no matter where it comes from:

$$\text{(logistic regression)} \quad \operatorname{sigmoid}(W_S^T \phi_S(x) + W_g^T \phi_g(x)) = \operatorname{sigmoid}(W_g^T \phi_g(x) + W_T^T \phi_T(x))$$

$\Rightarrow$  constraint  $W_S^T \phi_S(x) = W_T^T \phi_T(x)$  on unlabelled data

### 4) Joint representation (deep learning people like this!)

Learn a joint representation of data from  $D_S$  and  $D_T$



deep learning assumption: from low level to high level embedding, feature representation captures more "abstract", "conceptual" information.

TL/MTL assumption: data from  $D_S$  and  $D_T$  is similar often in the sense of "abstract information".

### TL/MTL ideas

② adapt model trained on  $D_S$  to  $D_T$ :

- 1) fine-tuning: take a pre-trained model, e.g. VGG network, then fine-tuning (a subset of) parameters and/or add ~~to~~ model to the previous one.

see Yosinski et al. NIPS 2014 for an empirical study of transferability of different layers.

2) lazy TL (Liu & Fukumizu on arxiv, 2015)

assume we have trained a classifier  $\hat{P}_S(y|x)$  on  $D_S$ , then the new classifier on  $D_T$  can be decomposed to

$\hat{P}_T(y|x) = \frac{\hat{P}_T(y|x)}{\hat{P}_S(y|x)} \hat{P}_S(y|x)$

estimator of  $P_S(y|x)$  data distribution  $\nearrow$

now assume we already have  $P_S(y|x)$  and might not be able to modify it (e.g. for privacy reasons),

then we can parameterize  $q(y, x; \theta) = \frac{\hat{P}_T(y|x)}{\hat{P}_S(y|x)}$

and minimize  $KL[P_T(y|x) \| q(y, x; \theta) \hat{P}_S(y|x)]$

now if the true ratio  $\frac{P_T(y, x)}{P_S(y, x)}$  is bounded, then

we can separate the problem to

$$KL[P_T(y|x) \| q(y, x; \theta) \hat{P}_S(y|x)]$$

$$\leq KL[P_T(y|x) \| q(y, x; \theta) P_S(y|x)]$$

$$+ (1+K) \underbrace{KL[P_S(y|x) \| \hat{P}_S(y|x)]}_{\text{only related to the pre-training step}}$$

and  $KL[P_T(y|x) \| q(y, x; \theta) P_S(y|x)] \propto -\mathbb{E}_{P_T}[\log q(y, x; \theta)]$

use this objective to train model  $q$

after training, use  $q \cdot \hat{P}_S$  as the classifier on  $D_T$

example:  $q(y, x; \theta) \propto \exp[\theta^T \Phi(y, x)]$

$$\text{and } \sum_y q(y, x; \theta) P_S(y|x) = 1$$

$$\text{and let } \hat{P}_S(y|x) \propto \exp[\beta^T \Phi(y, x)]$$

then by selecting  $\theta^T$ , we do feature selection on the target domain! (e.g. when  $\theta_i + \beta_i = 0$ , the  $i$ th feature  $\Phi_i(y, x)$  is not used)

Comment: I generally like this idea although their explanation & experiments are not very convincing.