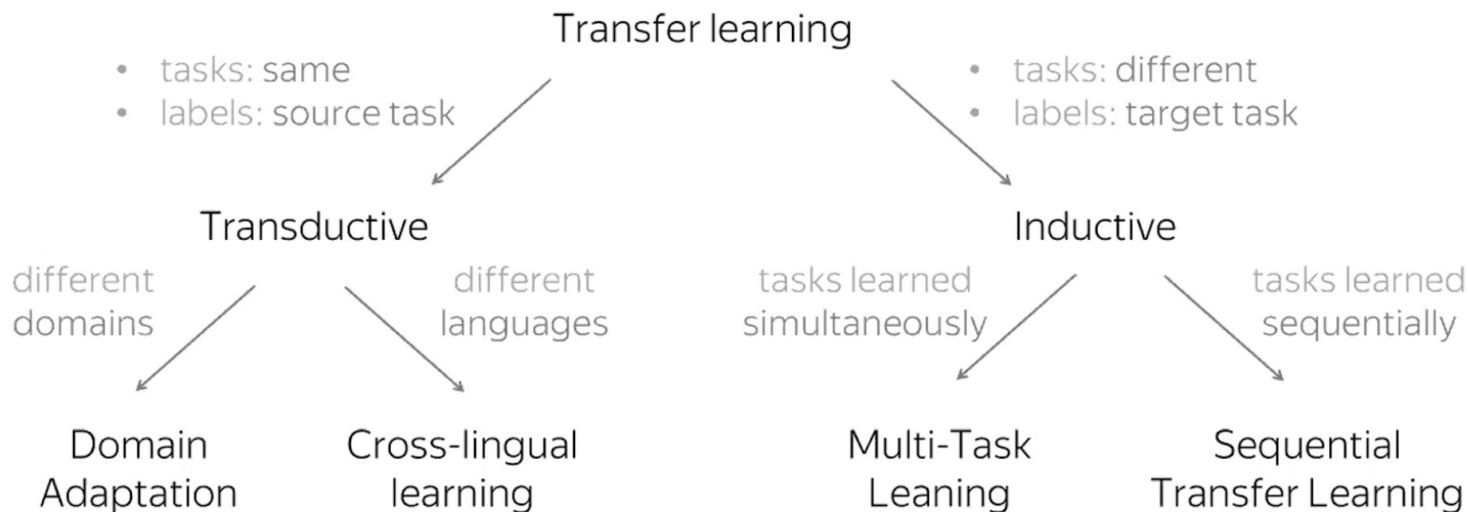


Domain Adaptation

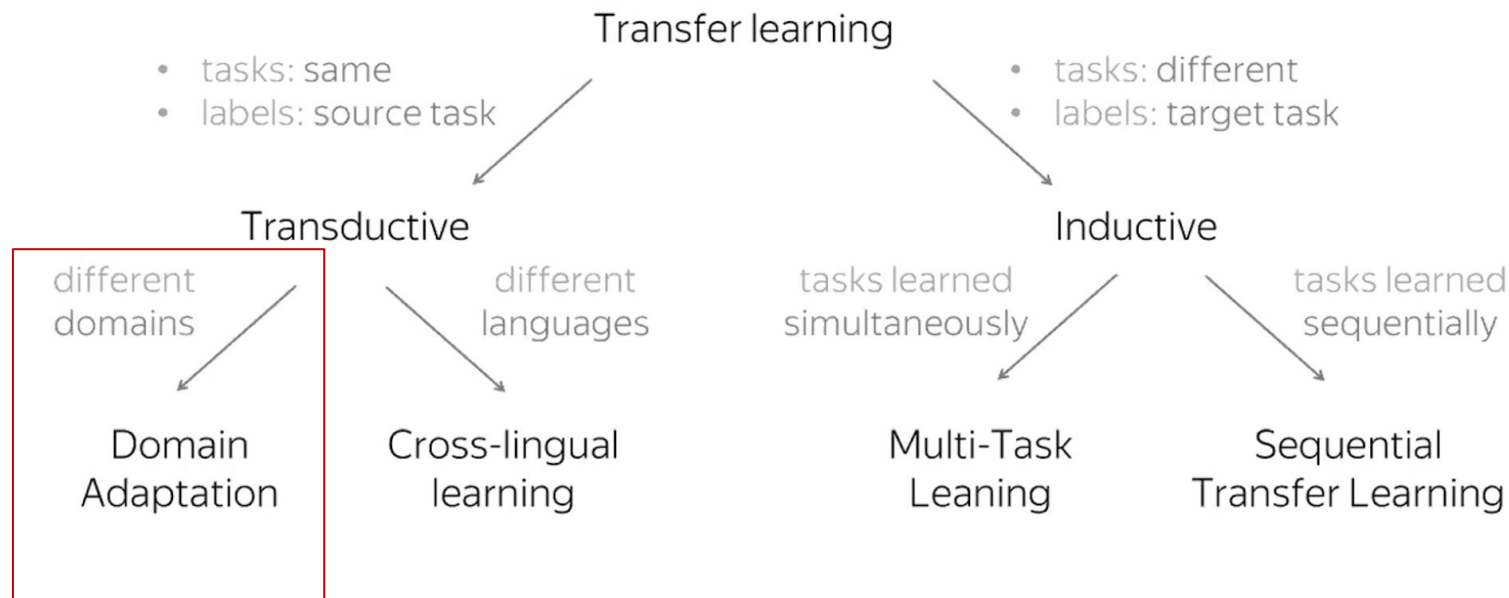
October 15, 2020

A Taxonomy of Transfer Learning in NLP



This taxonomy is from Ruder, 2019

A Taxonomy of Transfer Learning in NLP



This taxonomy is from Ruder, 2019

Standard supervised learning setting

“model” = model + data + task

Standard supervised learning setting

“model” = model + data + task

- Model (CNN vs RNN, GNMT vs Transformer and so on) $\mathcal{H} = \{h_{\theta} : X \rightarrow Y\}$

Standard supervised learning setting

“model” = model + data + task

- Model (CNN vs RNN, GNMT vs Transformer and so on) $\mathcal{H} = \{h_\theta : X \rightarrow Y\}$

- Data (source and target marginal distributions)
source ~ training data
target ~ testing data

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$

Standard supervised learning setting

“model” = model + data + task

- Model (CNN vs RNN, GNMT vs Transformer and so on) $\mathcal{H} = \{h_\theta : X \rightarrow Y\}$
- Data (source and target marginal distributions)
source ~ training data
target ~ testing data
$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$
$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$
- Task (loss function, risk)
$$h_\theta^* = \arg \min_{\theta} R^S(h_\theta) = \arg \min_{\theta} \mathbb{E}_{(x,y) \sim P_S} [L(h_\theta(x), y)]$$

Standard learning assumptions

Standard learning assumptions

- Sample points are IID

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \stackrel{iid}{\sim} P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \stackrel{iid}{\sim} P_T(x, y)$$

Standard learning assumptions

- Sample points are IID

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \stackrel{iid}{\sim} P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \stackrel{iid}{\sim} P_T(x, y)$$

- Same distributions for training and test

$$P_S(x, y) = P_T(x, y)$$

Standard learning assumptions

- Sample points are IID

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \stackrel{iid}{\sim} P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \stackrel{iid}{\sim} P_T(x, y)$$

- Same distributions for training and test

$$P_S(x, y) = P_T(x, y)$$

- Fixed distributions (they don't change with time)

Standard learning assumptions

- Sample points are IID

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \stackrel{iid}{\sim} P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \stackrel{iid}{\sim} P_T(x, y)$$

- Same distributions for training and test

$$P_S(x, y) = P_T(x, y)$$

- Fixed distributions (they don't change with time)

$$\epsilon_T = \mathbb{E}_{(x,y) \sim P_T} [L(h_\theta(x), y)]$$

$$\epsilon_T \leq \hat{\epsilon}_S + O\left(\frac{\text{complexity}(h)}{\sqrt{N}}\right)$$

$$\hat{\epsilon}_S = \frac{1}{N} \sum_{i=1}^N L(h_\theta(x_i^S), y_i)$$

Standard learning assumptions

- Sample points are IID

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \stackrel{iid}{\sim} P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \stackrel{iid}{\sim} P_T(x, y)$$

- Same distributions for training and test $P_S(x, y) = P_T(x, y)$
- Fixed distributions (they don't change with time)

EXPECTATIONS

Real-world problems

- Training sample is biased
- Noisy labels for training sample
- Sample points are not drawn IID
- Distributions drifts with time

REALITY

Domain shift:

$$P_S(x, y) \neq P_T(x, y)$$

Real-world examples

- Data from different sources
- Absence of good in-domain data
- Temporal (both short term and long term) changes of audience structure
- Synthetic training data

Real-world examples

- Data from different sources
- Absence of good in-domain data
- Temporal (both short term and long term) changes of audience structure
- Synthetic training data

Domain shift:

$$P_S(x, y) \neq P_T(x, y)$$

In-domain vs. out-of-domain generalization

On SQuAD Dani Yogatama et al., “Learning and Evaluating General Linguistic Intelligence,”

<http://arxiv.org/abs/1901.11373>

	SQuAD	Trivia	QuAC	QA-SRL	QA-ZRE
BERT	86.5 (78.5)	35.6 (13.4)	56.2 (43.9)	77.5 (65.0)	55.3 (40.0)
ELMo	81.8 (72.2)	32.9 (12.6)	45.0 (34.5)	68.7 (52.3)	60.2 (42.0)

Table 2: F_1 (exact match) scores of the best BERT and ELMo models trained on SQuAD and evaluated on other question answering datasets.

On MNLI: R. Thomas McCoy, Junghyun Min, and Tal Linzen, “BERTs of a Feather Do Not Generalize Together: Large Variability in Generalization across Models with Similar Test Set Performance,” <http://arxiv.org/abs/1911.02969>

Heuristic	Definition	Example
Lexical overlap	Assume that a premise entails all hypotheses constructed from words in the premise	The doctor was paid by the actor. → The doctor paid the actor. WRONG
Subsequence	Assume that a premise entails all of its contiguous subsequences.	The doctor near the actor danced. → The actor danced. WRONG
Constituent	Assume that a premise entails all complete subtrees in its parse tree.	If the artist slept , the actor ran. → The artist slept. WRONG

Figure 1: The heuristics targeted by the HANS dataset, along with examples of incorrect entailment predictions that these heuristics would lead to. (Figure from McCoy et al. (2019).)

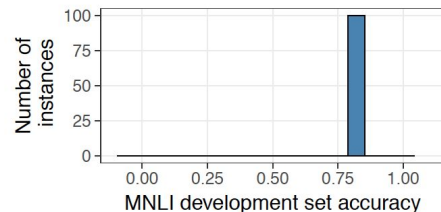


Figure 2: In-distribution generalization: Performance on the MNLI development set

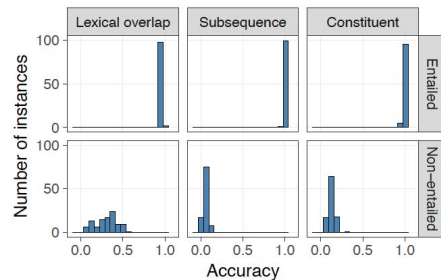


Figure 3: Out-of-distribution generalization: Performance on the HANS evaluation set, broken down into six categories of examples based on which syntactic heuristic each example targets and whether the correct label is *entailment* or *non-entailment*. The non-entailed lexical overlap cases (lower left plot) display a large degree of variability across instances.

In-domain vs. out-of-domain generalization

On MNLI: R. Thomas McCoy, Junghyun Min, and Tal Linzen, “BERTs of a Feather Do Not Generalize Together: Large Variability in Generalization across Models with Similar Test Set Performance,” <http://arxiv.org/abs/1911.02969>

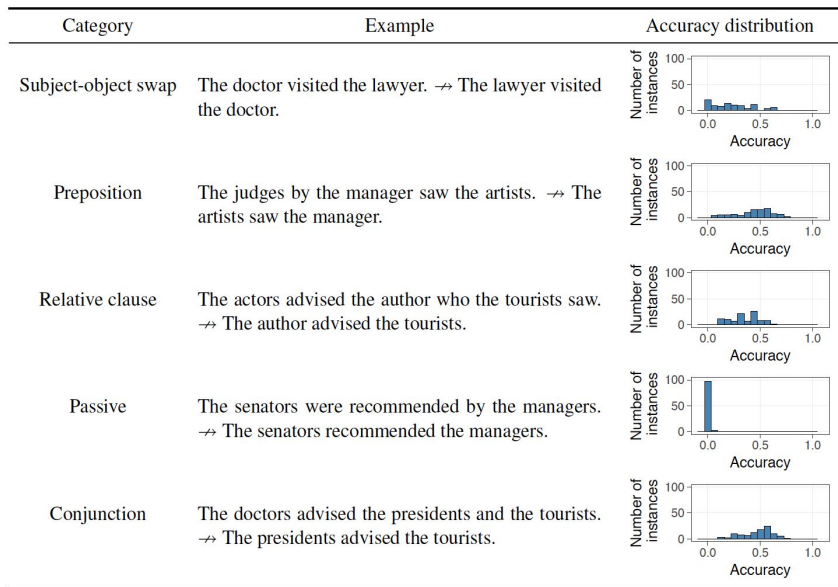


Figure 4: Results on the various subcategories within the non-entailed lexical overlap examples of the HANS dataset. We do not include the other 25 subcategories of the HANS dataset in this figure as there was little variability across instances for those subcategories.

How to check whether there is a domain shift or not?



How to check whether there is a domain shift or not?

If domain distribution are different we can train classifier to discriminate them.

$$P_S(x, y) \neq P_T(x, y)$$

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \epsilon_S + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

$$P_S(x, y) \neq P_T(x, y)$$

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\boxed{\epsilon_T} \leq \epsilon_S + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

Expected target error

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \boxed{\epsilon_S} + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

Expected source error

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \epsilon_S + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

The divergence between
source and target domain

$$d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) = 2 \sup_{(h, h')} \left| \mathbb{E}_{(x, y) \sim P_S} [h(x) \neq h'(x)] - \mathbb{E}_{(x, y) \sim P_T} [h(x) \neq h'(x)] \right|$$

H-divergence measures the worst case of the disagreement between a pair of hypothesis.
How much features discriminative for S and T

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \epsilon_S + \boxed{\frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T)} + \lambda$$

The divergence between
source and target domain

$$d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) \sim d_{\mathcal{A}} = 2(1 - 2\epsilon)$$

As H-divergence approximation proxy A-distance (PAD) can be used.
 ϵ is the generalization error of the domain classifier.

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \epsilon_S + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

Error of ideal joint hypothesis

$$\lambda = \min_h [\epsilon_S(h) + \epsilon_T(h)]$$

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \epsilon_S + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

Generalization error under domain shift

[Ben-David, 2010](#)

Ben-David theory:

$$\epsilon_T \leq \epsilon_S + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda$$

Domain adaptation tend to minimize second and (third?) terms.

Domain adaptation

- Unsupervised DA:
No labels for target domain

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{T} = \{x_i^T\}_{i=1}^M \sim P_T(x)$$

- (Semi)supervised DA:
Target domain dataset is (partially) labeled

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$

Covariate shift:

$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

Why does a covariate shift make problems?

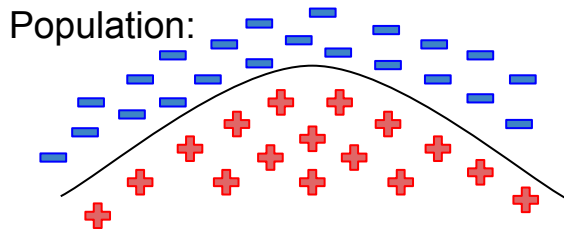
$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

Why does a covariate shift make problems?

$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

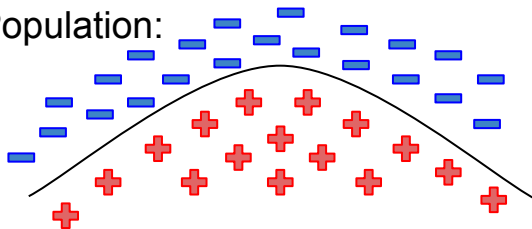


Why does a covariate shift make problems?

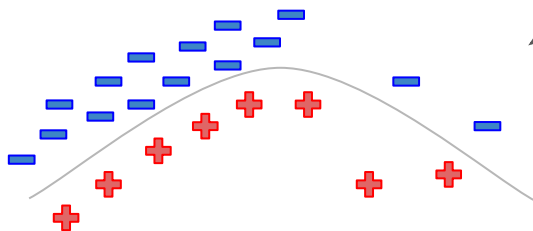
$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

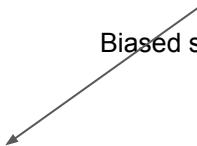
Population:



Source distribution:



Biased sampling

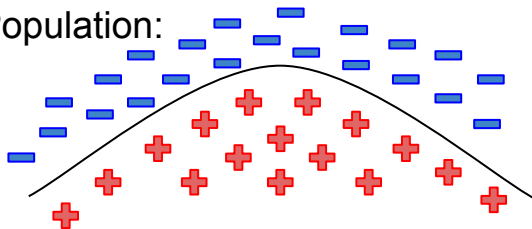


Why does a covariate shift make problems?

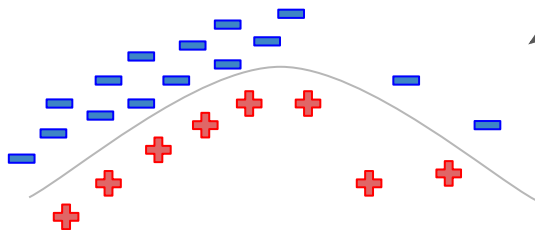
$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

Population:

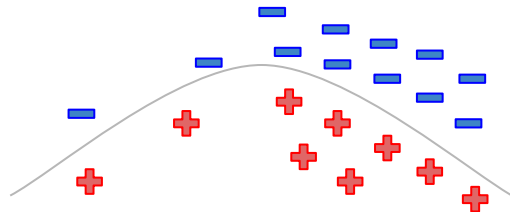


Source distribution:



Biased sampling

Target distribution:

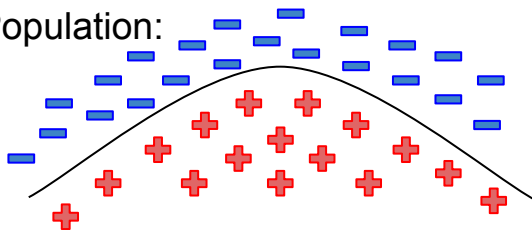


Why does a covariate shift make problems?

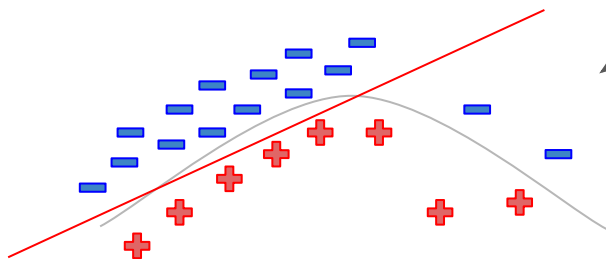
$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

Population:

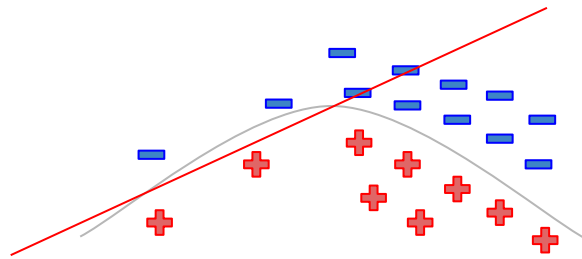


Source distribution:



Biased sampling

Target distribution:

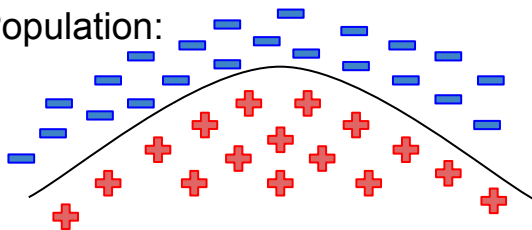


Why does a covariate shift make problems?

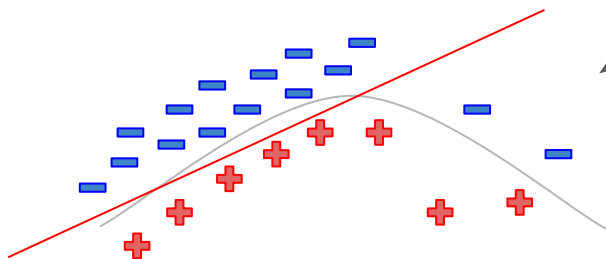
$$P_S(x) \neq P_T(x)$$

$$P_S(y|x) = P_T(y|x)$$

Population:

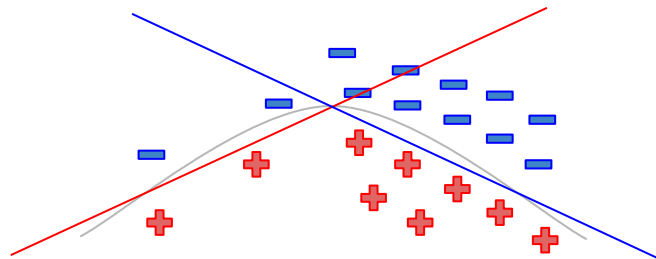


Source distribution:



Biased sampling

Target distribution:



Not optimal decision rule!

Unsupervised domain adaptation

Unsupervised domain adaptation

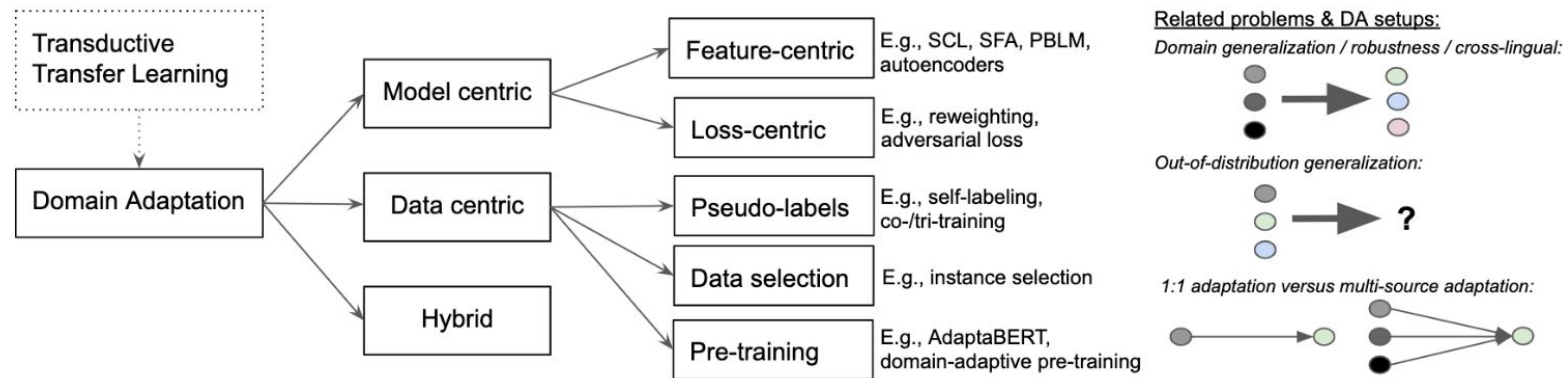


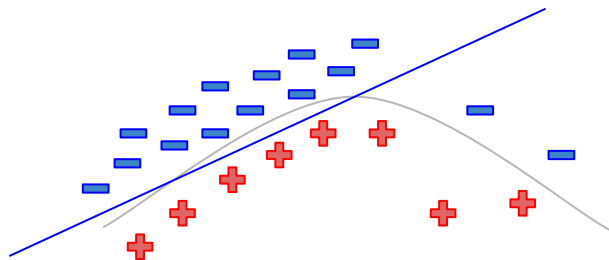
Figure 1: Taxonomy of DA as special case of transductive transfer learning, DA setups (1:1 versus multi-source adaptation) and related problems (domain and out-of-distribution generalization).

[Neural Unsupervised Domain Adaptation in NLP—A Survey](#)

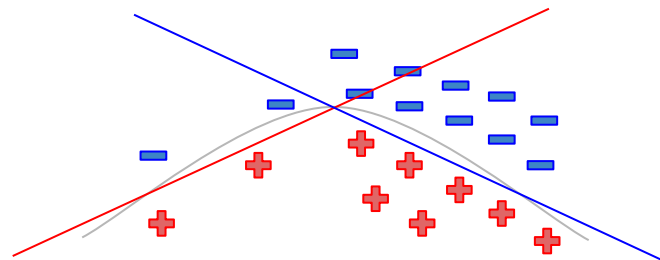
Instance weighting

$$h_{\theta}^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N L(h_{\theta}(x), y)$$

Source distribution:



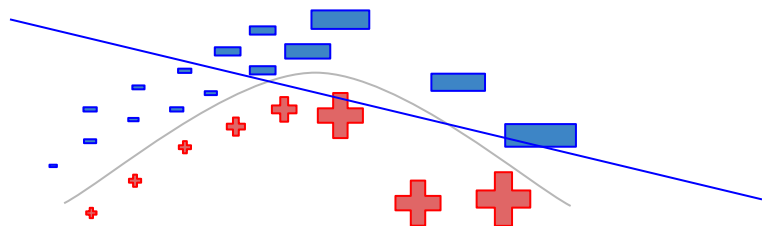
Target distribution:



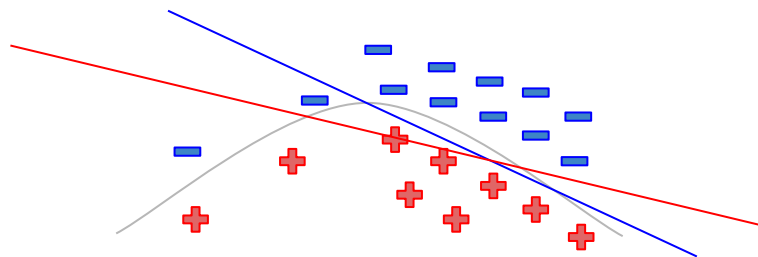
Not optimal decision rule!

Instance weighting

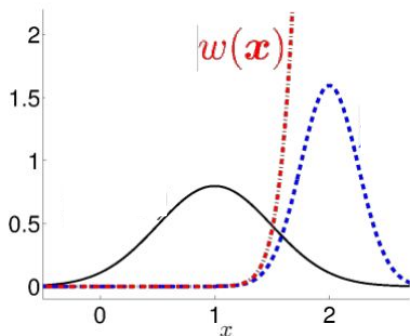
Source distribution:



Target distribution:



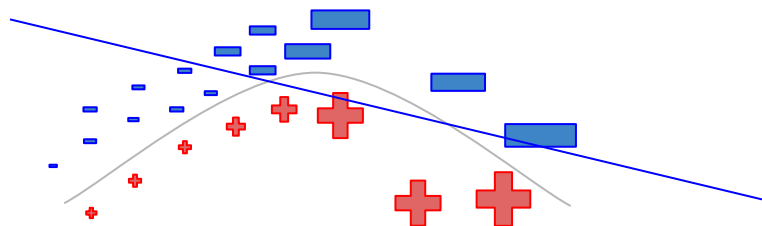
Instance weighting



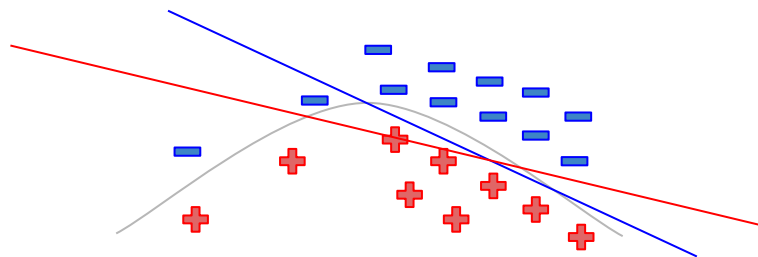
$$h_{\theta}^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N w(x_i) L(h_{\theta}(x_i), y_i)$$

$$w(x) = \frac{p_T(x)}{p_S(x)}$$

Source distribution:



Target distribution:



Instance weighting: NLP applications

- Neural Machine Translation

$x = (v_1, v_2, \dots, v_k)$: Input sequence of tokens

$P_T(x), P_S(x) = ?$

Instance weighting: NLP applications

- Neural Machine Translation

$x = (v_1, v_2, \dots, v_k)$: Input sequence of tokens

$P_T(x), P_S(x)$ can be defined using source and target language models

Instance weighting: NLP applications

- Neural Machine Translation

$x = (v_1, v_2, \dots, v_k)$: Input sequence of tokens

$P_T(x), P_S(x)$ can be defined using source and target language models

Weighting using language models ([Wang et al. 2017](#)):

$$w_i = \delta[H_S(x) - H_T(x)]$$

$$H = -\frac{1}{k} \log P(x)$$

(δ - min-max normalization)

Instance weighting: NLP applications

- Neural Machine Translation

$x = (v_1, v_2, \dots, v_k)$: Input sequence of tokens

$P_T(x), P_S(x)$ can be defined using source and target language models

Weighting using domain classifier ([Chen et al. 2017](#)):

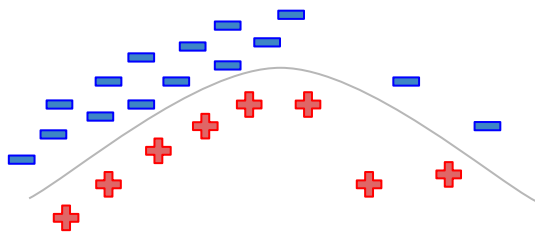
$$w_i = (1 + p_d(x_i))$$

$p_d(x)$ - probability of being from target

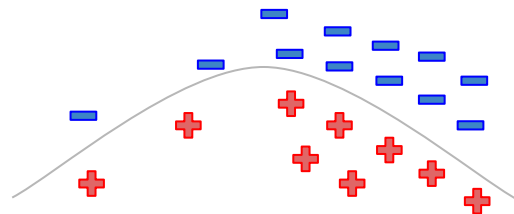
Data selection

Instead of weighting, we can train classifier and drop all observations that are too dissimilar from target domain.

Source distribution:



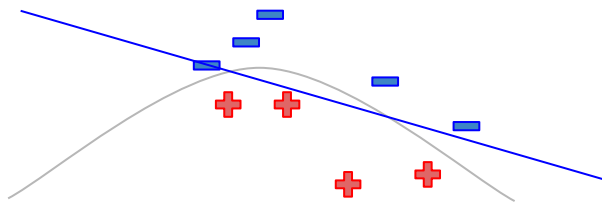
Target distribution:



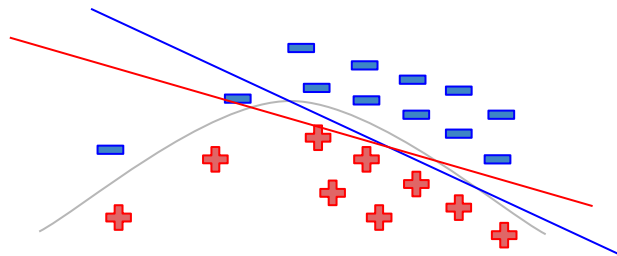
Data selection

Instead of weighting, we can train classifier and drop all observations that are too dissimilar from target domain.

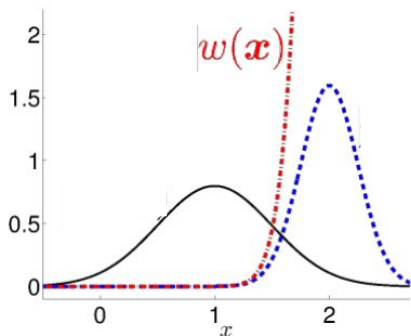
Source distribution:



Target distribution:



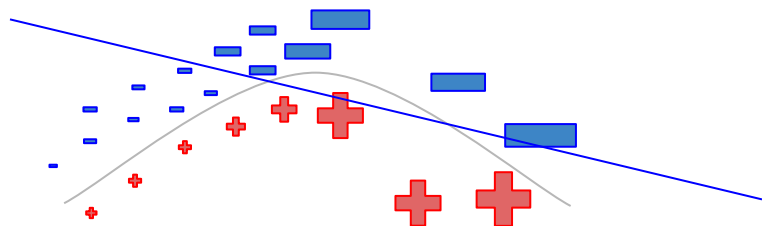
Instance weighting: problems



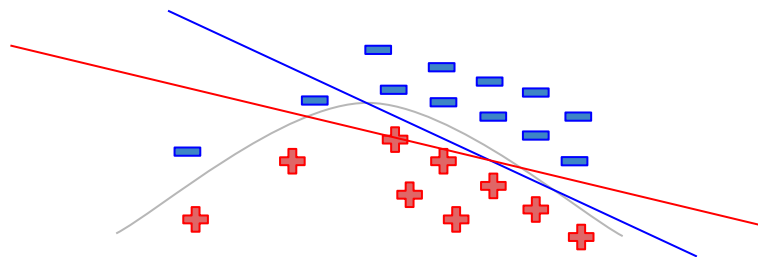
$$\hat{\epsilon}_S(w, x) = \frac{1}{N} \sum_{i=1}^N w(x_i) L(h_{\theta}(x_i^S), y_i^S)$$

$$\epsilon_T \leq \hat{\epsilon}_S(w, x) + \sqrt{\frac{O(\max_x w(x)^2)}{N}}$$

Source distribution:



Target distribution:



Proxy-labels methods

In previous case we use target domain data only for weight calculation. And do not use knowledge that target domain data contain.

It's good to utilize somehow unlabelled data for training.

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{T} = \{x_i^T\}_{i=1}^M \sim P_T(x)$$

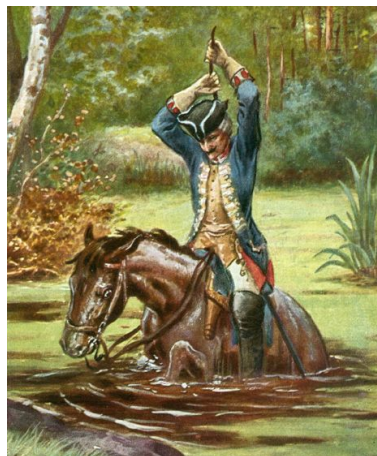
Self-training

Based on [Ruder's presentation at ACL'18](#)

1. Train model on labeled data.
2. Use confident predictions on unlabeled data as training examples. Repeat.

Algorithm 1 Self-training (Abney, 2007)

```
1: repeat  
2:    $m \leftarrow \text{train\_model}(L)$   
3:   for  $x \in U$  do  
4:     if  $\max m(x) > \tau$  then  
5:        $L \leftarrow L \cup \{(x, p(x))\}$   
6: until no more predictions are confident
```

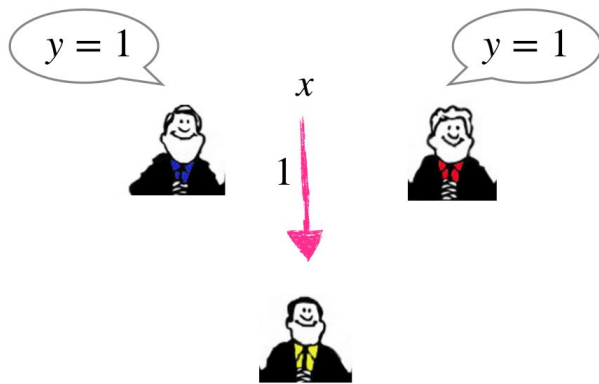


Error is amplified cause a model can not correct its own mistakes.

Tri-training

Based on [Ruder's presentation at ACL'18](#)

1. Train three models on bootstrapped samples.
2. Use predictions on unlabeled data for third if two agree.
3. Final prediction: majority voting



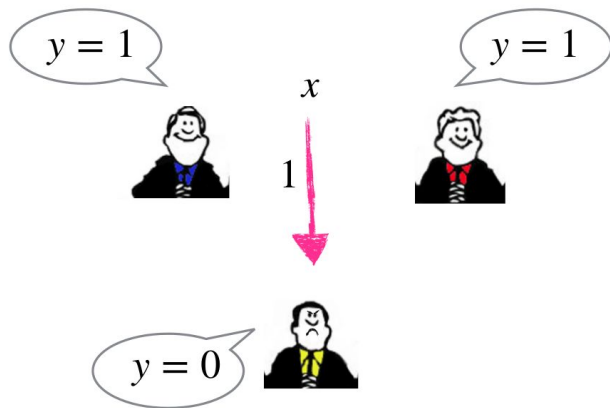
Algorithm 2 Tri-training (Zhou and Li, 2005)

```
1: for  $i \in \{1..3\}$  do
2:    $S_i \leftarrow \text{bootstrap\_sample}(L)$ 
3:    $m_i \leftarrow \text{train\_model}(S_i)$ 
4: repeat
5:   for  $i \in \{1..3\}$  do
6:      $L_i \leftarrow \emptyset$ 
7:     for  $x \in U$  do
8:       if  $p_j(x) = p_k(x) (j, k \neq i)$  then
9:          $L_i \leftarrow L_i \cup \{(x, p_j(x))\}$ 
10:         $m_i \leftarrow \text{train\_model}(L \cup L_i)$ 
11: until none of  $m_i$  changes
12: apply majority vote over  $m_i$ 
```

Tri-training with disagreement

Based on [Ruder's presentation at ACL'18](#)

1. Train three models on bootstrapped samples.
2. Use predictions on unlabeled data for third if two agree and prediction differs.
3. Final prediction: majority voting



The problem with tri-training is that training three separate models can be too expensive.

Let's share parameters somehow.

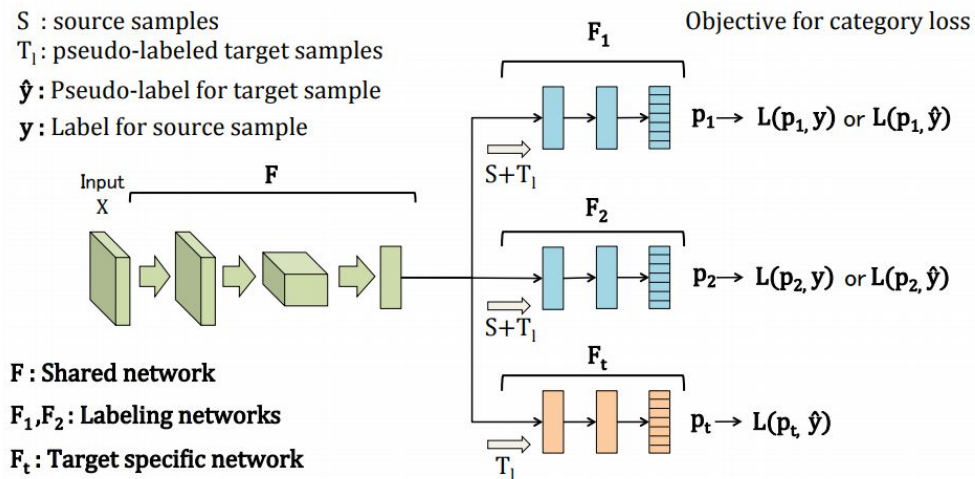
Asymmetric Tri-training

S : source samples

T_1 : pseudo-labeled target samples

\hat{y} : Pseudo-label for target sample

y : Label for source sample



Input: data

$$\mathcal{S} = \{(x_i, t_i)\}_{i=1}^m, \mathcal{T} = \{(x_j)\}_{j=1}^n$$

$$\mathcal{T}_l = \emptyset$$

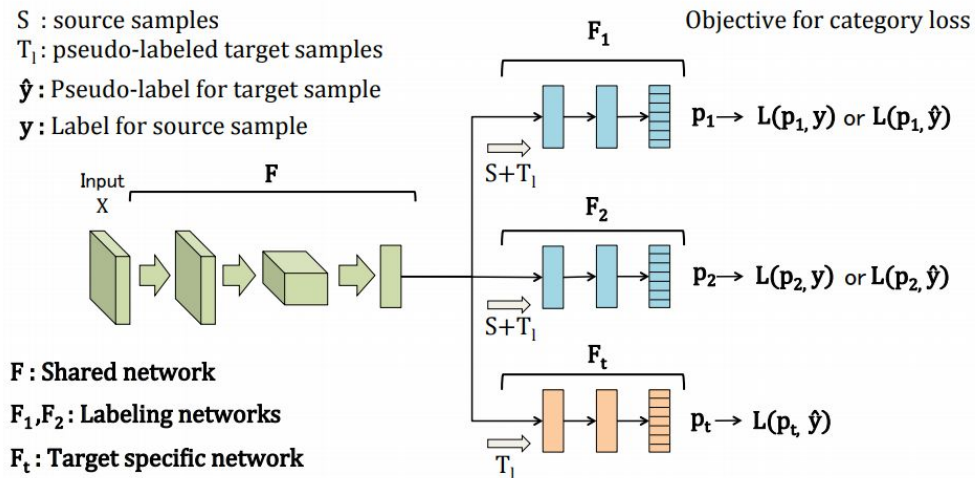
Asymmetric Tri-training

S : source samples

T_1 : pseudo-labeled target samples

\hat{y} : Pseudo-label for target sample

y : Label for source sample



$$E(\theta_F, \theta_{F_1}, \theta_{F_2}) = \frac{1}{n} \sum_{i=1}^n [L_y(F_1 \circ F(x_i)), y_i] + L_y(F_2 \circ (F(x_i)), y_i) + \lambda |W_1^T W_2|$$

Input: data

$$\mathcal{S} = \{(x_i, t_i)\}_{i=1}^m, \mathcal{T} = \{(x_j)\}_{j=1}^n$$

$$\mathcal{T}_l = \emptyset$$

for $j = 1$ **to** $iter$ **do**

Train F, F_1, F_2, F_t with mini-batch from training set \mathcal{S}

end for

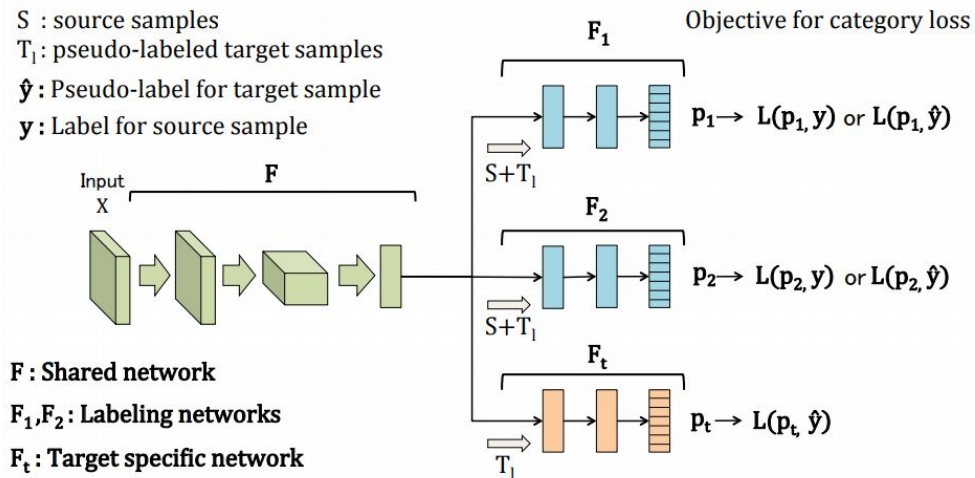
Asymmetric Tri-training

S : source samples

T_1 : pseudo-labeled target samples

\hat{y} : Pseudo-label for target sample

y : Label for source sample



$$E(\theta_F, \theta_{F_1}, \theta_{F_2}) = \frac{1}{n} \sum_{i=1}^n [L_y(F_1 \circ F(x_i)), y_i] + L_y(F_2 \circ (F(x_i)), y_i) + \lambda |W_1^T W_2|$$

Input: data

$$\mathcal{S} = \{(x_i, t_i)\}_{i=1}^m, \mathcal{T} = \{(x_j)\}_{j=1}^n$$

$$\mathcal{T}_l = \emptyset$$

for $j = 1$ **to** $iter$ **do**

Train F, F_1, F_2, F_t with mini-batch from training set \mathcal{S}

end for

To force F1 and F2 to learn from different features.

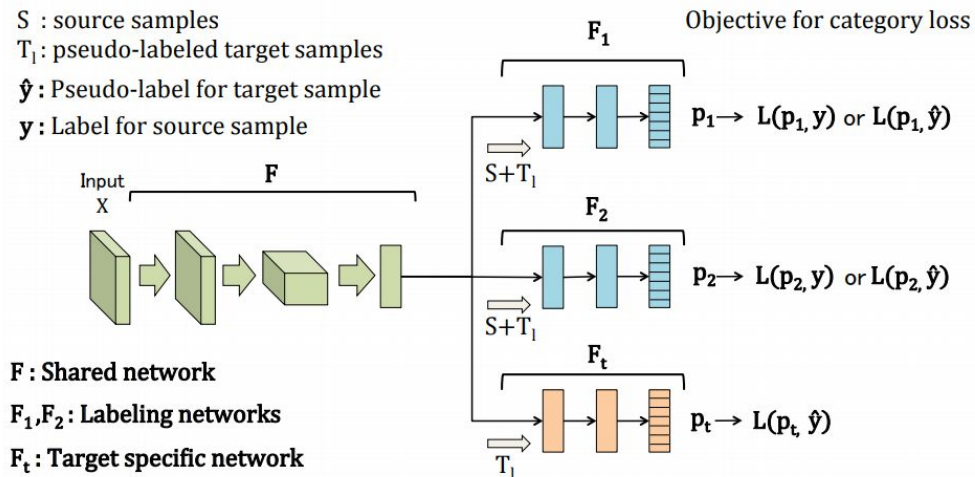
Asymmetric Tri-training

S : source samples

T_l : pseudo-labeled target samples

\hat{y} : Pseudo-label for target sample

y : Label for source sample



F : Shared network

F_1, F_2 : Labeling networks

F_t : Target specific network

$$E(\theta_F, \theta_{F_1}, \theta_{F_2}) = \frac{1}{n} \sum_{i=1}^n [L_y(F_1 \circ F(x_i)), y_i] + L_y(F_2 \circ (F(x_i)), y_i) + \lambda |W_1^T W_2|$$

Input: data

$$\mathcal{S} = \{(x_i, t_i)\}_{i=1}^m, \mathcal{T} = \{(x_j)\}_{j=1}^n$$

$$\mathcal{T}_l = \emptyset$$

for $j = 1$ **to** $iter$ **do**

Train F, F_1, F_2, F_t with mini-batch from training set \mathcal{S}

end for

$$N_t = N_{init}$$

$$\mathcal{T}_l = \text{Labeling}(F, F_1, F_2, \mathcal{T}, N_t)$$

$$\mathcal{L} = \mathcal{S} \cup \mathcal{T}_l$$

for k steps **do**

for $j = 1$ **to** $iter$ **do**

Train F, F_1, F_2 with mini-batch from training set \mathcal{L}

Train F, F_t with mini-batch from training set \mathcal{T}_l

end for

$$\mathcal{T}_l = \emptyset, N_t = k/20 * n$$

$$\mathcal{T}_l = \text{Labeling}(F, F_1, F_2, \mathcal{T}, N_t)$$

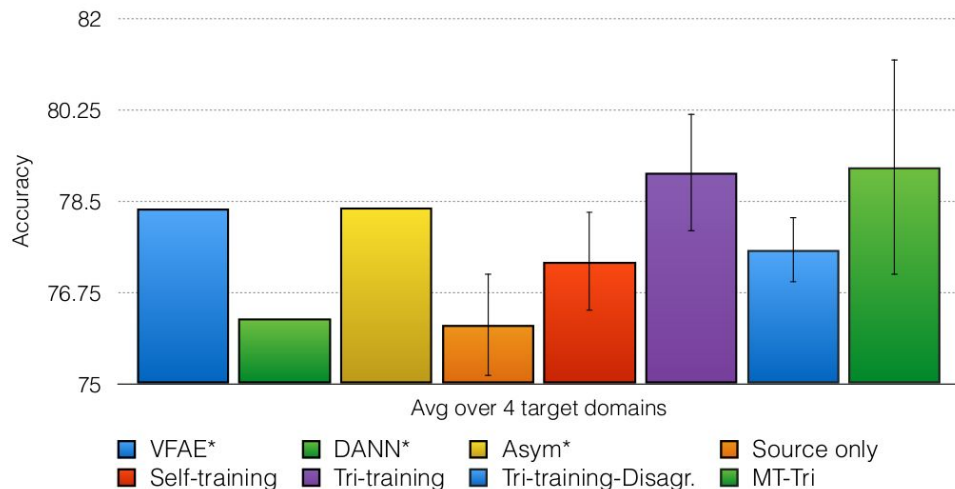
$$\mathcal{L} = \mathcal{S} \cup \mathcal{T}_l$$

end for

On effectiveness of proxy-label methods

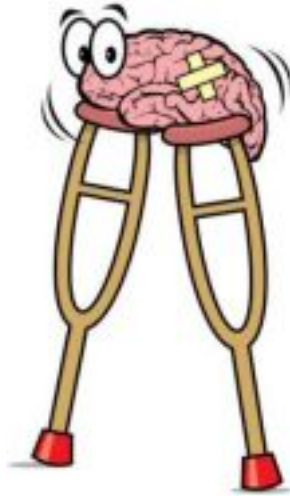
Sentiment Analysis Results

Based on [Ruder's presentation at ACL'18](#)



Sentiment analysis on Amazon reviews dataset (Blitzer et al, 2006)

Hack of the day...



Hack of the day: back-translation

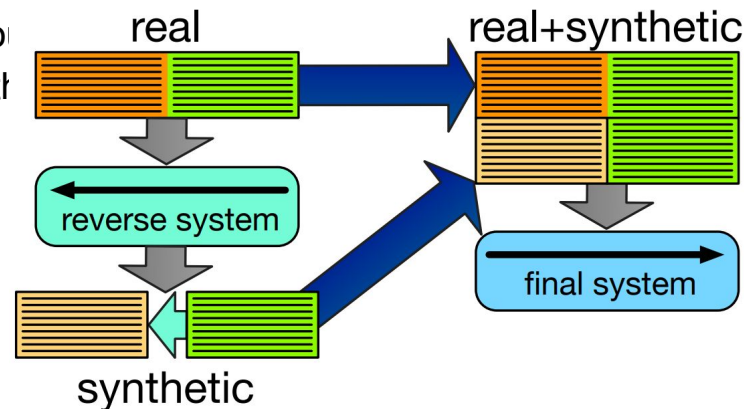
[Sennrich, 2016b](#)

In proxy-label methods we generate output labels based on unlabeled input data.

But what if we want to generate input data based on output? Does it make sense?

NMT is trained in the reverse translation direction (target-to-source) then used to translate target-side monolingual data back into the source language.

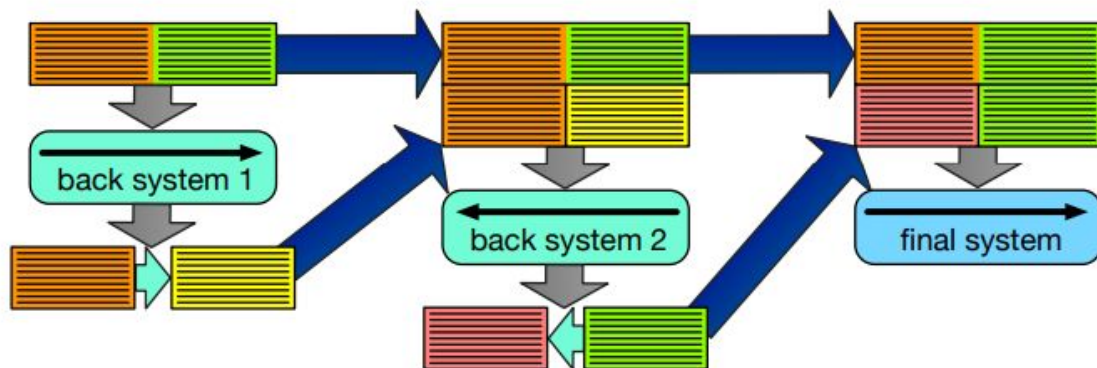
Picture from [Hoang et al., 2018](#)



Step further: iterative back-translation

[Hoang et al., 2018](#)

Iterative back-translation: back-translated data is used to build better translation systems in forward and backward directions, which in turn is used to reback-translate monolingual data



Wait a second...

Why do we need proxy-labels?

What about **unsupervised pre-training**?

Unsupervised in-domain pre-training is effective

Don't Stop Pretraining: Adapt Language Models to Domains and Tasks

[Gururangan et al., 2020](#)

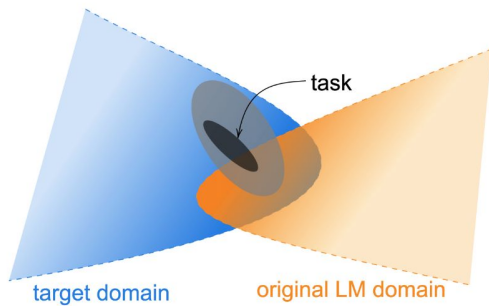


Figure 1: An illustration of data distributions. Task data is comprised of an observable task distribution, usually non-randomly sampled from a wider distribution (light grey ellipsis) within an even larger target domain, which is not necessarily one of the domains included in the original LM pretraining domain – though overlap is possible. We explore the benefits of continued pretraining on data from the task distribution and the domain distribution.

Domain	Task	RoBERTa	Additional Pretraining Phases		
			DAPT	TAPT	DAPT + TAPT
BIO MED	CHEMPROT	81.9 _{1.0}	84.2 _{0.2}	82.6 _{0.4}	84.4 _{0.4}
	†RCT	87.2 _{0.1}	87.6 _{0.1}	87.7 _{0.1}	87.8 _{0.1}
CS	ACL-ARC	63.0 _{5.8}	75.4 _{2.5}	67.4 _{1.8}	75.6 _{3.8}
	SCIERC	77.3 _{1.9}	80.8 _{1.5}	79.3 _{1.5}	81.3 _{1.8}
NEWS	HYPERPARTISAN	86.6 _{0.9}	88.2 _{5.9}	90.4 _{5.2}	90.0 _{6.6}
	†AGNEWS	93.9 _{0.2}	93.9 _{0.2}	94.5 _{0.1}	94.6 _{0.1}
REVIEWS	†HELPLEFULNESS	65.1 _{3.4}	66.5 _{1.4}	68.5 _{1.9}	68.7 _{1.8}
	†IMDB	95.0 _{0.2}	95.4 _{0.1}	95.5 _{0.1}	95.6 _{0.1}

Table 5: Results on different phases of adaptive pretraining compared to the baseline RoBERTa (col. 1). Our approaches are DAPT (col. 2, §3), TAPT (col. 3, §4), and a combination of both (col. 4). Reported results follow the same format as Table 3. State-of-the-art results we can compare to: CHEMPROT (84.6), RCT (92.9), ACL-ARC (71.0), SCIERC (81.8), HYPERPARTISAN (94.8), AGNEWS (95.5), IMDB (96.2); references in §A.2.

Unsupervised in-domain pre-training is effective

Don't Stop Pretraining: Adapt Language Models to Domains and Tasks

[Gururangan et al., 2020](#)

Pretraining	BIO MED		CS
	CHEMPROT	RCT-500	ACL-ARC
ROBERTA	81.9 _{1.0}	79.3 _{0.6}	63.0 _{5.8}
TAPT	82.6 _{0.4}	79.8 _{1.4}	67.4 _{1.8}
RAND-TAPT	81.9 _{0.6}	80.6 _{0.4}	69.7 _{3.4}
50NN-TAPT	83.3 _{0.7}	80.8 _{0.6}	70.7 _{2.8}
150NN-TAPT	83.2 _{0.6}	81.2 _{0.8}	73.3 _{2.7}
500NN-TAPT	83.3 _{0.7}	81.7 _{0.4}	75.5 _{1.9}
DAPT	84.2 _{0.2}	82.5 _{0.5}	75.4 _{2.5}

Table 8: Mean test set micro- F_1 (for CHEMPROT and RCT) and macro- F_1 (for ACL-ARC), across five random seeds, with standard deviations as subscripts, comparing RAND-TAPT (with 50 candidates) and k NN-TAPT selection. Neighbors of the task data are selected from the domain data.

Unsupervised in-domain pre-training is effective

Don't Stop Pretraining: Adapt Language Models to Domains and Tasks

[Han et al., 2020](#)

AdaptaBERT:

1. Domain tuning: (50/50 source+target data unsupervised pretraining.
2. Task tuning using source data

System	WNUT				CoNLL
	domain adaptation data	Precision	Recall	F1	F1
<i>Unsupervised domain adaptation</i>					
1. Task-tuned BERT	n/a	50.9	66.6	57.7	97.8
2. AdaptaBERT	WNUT training	52.8	66.7	58.9	97.6
3. AdaptaBERT	+ 1M tweets	53.6	68.3	60.0	97.8
4. AdaptaBERT	WNUT train+test	57.7	68.9	62.8	97.8
<i>Supervised in-domain training</i>					
5. Fine-tuned BERT	n/a	66.3	62.3	64.3	80.9
6. Limsopatham and Collier (2016)	n/a	73.5	59.7	65.9	

Table 4: Named entity segmentation performance on the WNUT test set and CoNLL test set A. [Limsopatham and Collier \(2016\)](#) had the winning system at the 2016 WNUT shared task. Their results are reprinted from their paper, which did not report performance on the CoNLL dataset.

(Semi)-supervised domain adaptation

Both in source and target domain labels are known.

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$

Fine-tuning

- Train model on the target domain data
- Train with lower learning rate on the source domain

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$

Fine-tuning

- Train model on the target domain data
- Train with lower learning rate on the source domain

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$

Problems?

Fine-tuning

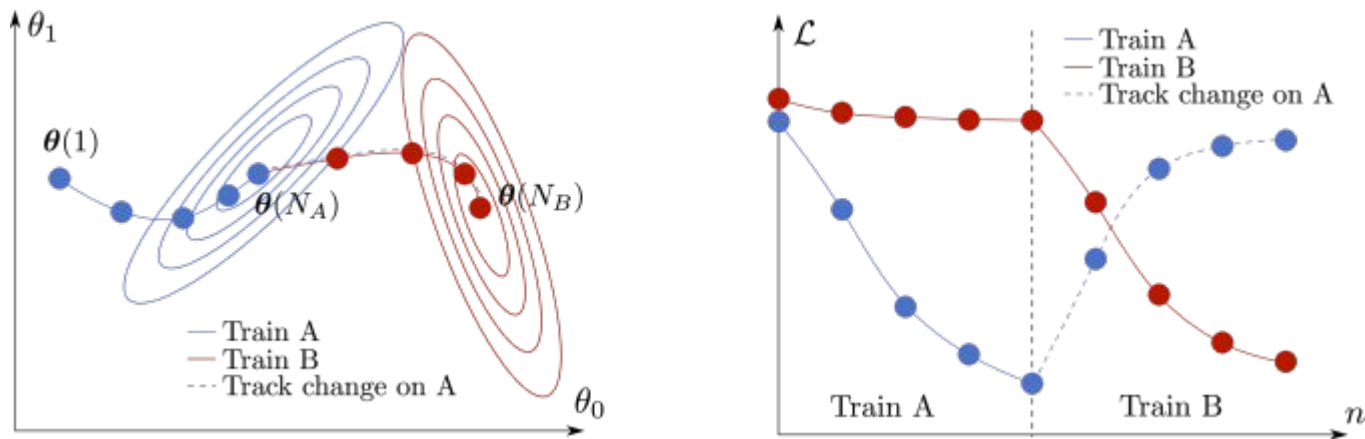
- Train model on the labeled target domain data
- Train with lower learning rate on the labeled source domain data

$$\mathcal{D}_S = \{x_i^S, y_i^S\}_{i=1}^N \sim P_S(x, y)$$

$$\mathcal{D}_T = \{x_i^T, y_i^T\}_{i=1}^M \sim P_T(x, y)$$

Ordinarily, in-domain data is limited, therefore we have generalization problems (over-fitting).

Fine-tuning: catastrophic forgetting



Picture from [Felix Wiewel](#), [Bin Yang](#), 2019

Fine-tuning: catastrophic forgetting

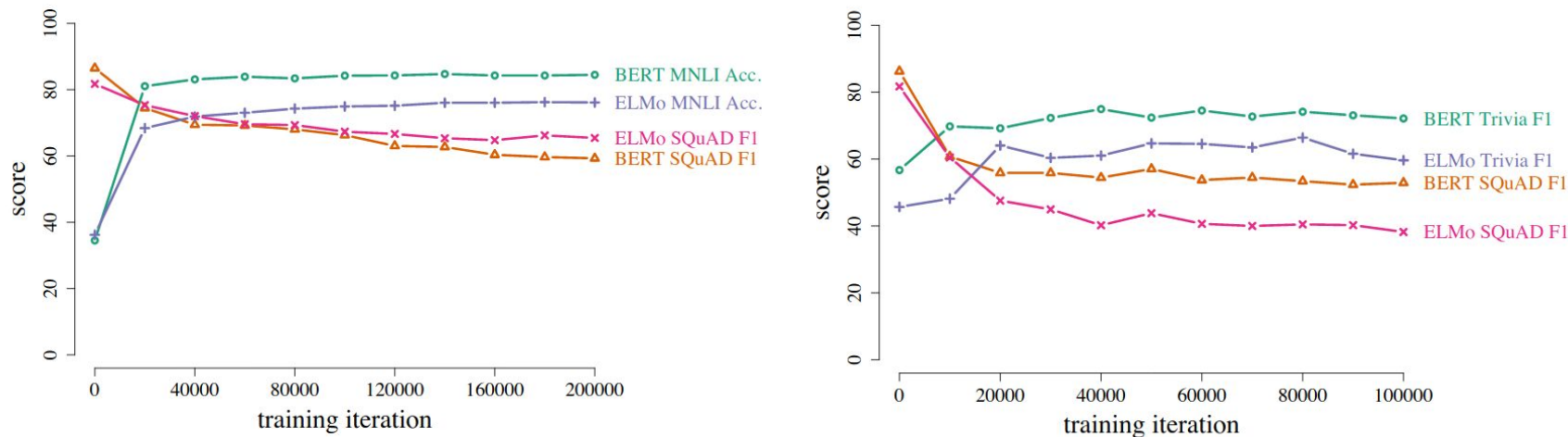


Figure 4: Catastrophic forgetting in a continual learning setup on unsupervised \rightarrow SQuAD \rightarrow MNLI (top) and unsupervised \rightarrow SQuAD \rightarrow TriviaQA (bottom).

Fine-tuning: regularization techniques

Ordinarily, in-domain data is limited, therefore we have generalization problems (over-fitting). To prevent this:

Fine-tuning: regularization techniques

Ordinarily, in-domain data is limited, therefore we have generalization problems (over-fitting). To prevent this:

1. **Weight decay** $L_W = ||W||_{L_p}$

W is the in-domain parameter matrix to be learned

\hat{W} is the corresponding fixed out-of-domain parameter matrix.

Fine-tuning: regularization techniques

Ordinarily, in-domain data is limited, therefore we have generalization problems (over-fitting). To prevent this:

1. Weight decay $L_W = ||W||_{L_p}$

2. Dropout

W is the in-domain parameter matrix to be learned

\hat{W} is the corresponding fixed out-of-domain parameter matrix.

Fine-tuning: regularization techniques

Ordinarily, in-domain data is limited, therefore we have generalization problems (over-fitting). To prevent this:

1. Weight decay $L_W = ||W||_{L_p}$
2. Dropout
3. L2-distance from out-of-domain penalization (MAP-L2) $L_W = \lambda \cdot ||W - \hat{W}||_2^2$

W is the in-domain parameter matrix to be learned

\hat{W} is the corresponding fixed out-of-domain parameter matrix.

Fine-tuning: regularization techniques

[Barone et al., 2017](#)

Out-of-domain: WMT

In-domain: TED talks

Table 1: English-to-German translation BLEU scores

System	valid	test			avg
	tst2010	tst2011	tst2012	tst2013	
Out-of-domain only	27.19	29.65	25.78	27.85	27.76
In-domain only	25.95	27.84	23.68	25.83	25.78
Fine-tuning	30.53	32.62	28.86	32.11	31.20
Fine-tuning + dropout	30.63	33.06	28.90	32.02	31.33
Fine-tuning + MAP-L2	30.81	32.87	28.99	31.88	31.25
Fine-tuning + tuneout	30.49	32.07	28.66	31.60	30.78†
Fine-tuning + dropout + MAP-L2	30.80	33.19	29.13	32.13	31.48†

†: different from the fine-tuning baseline at 5% significance.

Limited degradation on out-of-domain

L2-distance from out-of-domain penalization (MAP-L2) limits quality degradation on the out-of-domain.

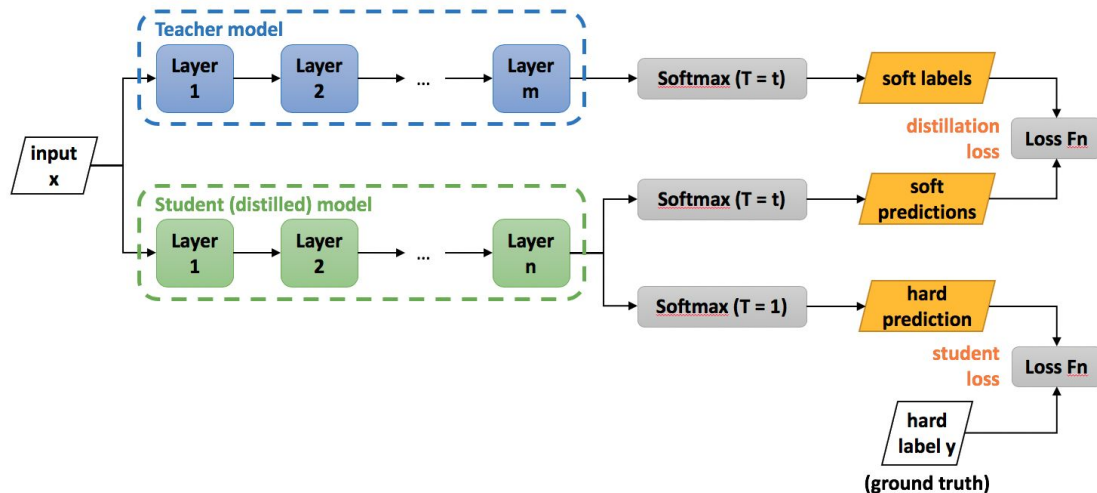
Limited degradation on out-of-domain

L2-distance from out-of-domain penalization (MAP-L2) limits quality degradation on the out-of-domain.

But we can directly force predictions of in-domain model and out-of-domain model to be closer to each other!

Distillation-like domain adaptation

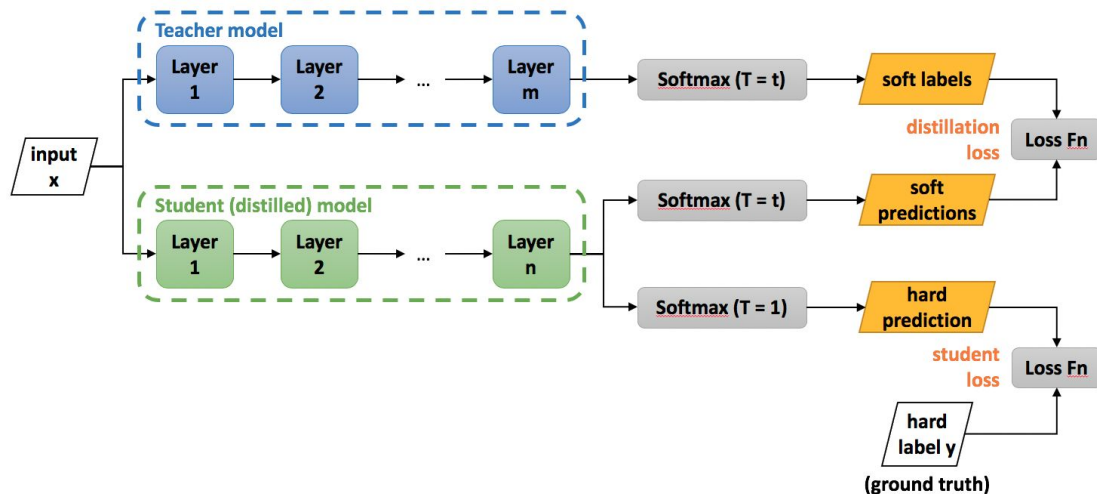
“Knowledge distillation” framework of (Hinton et al., 2014): a smaller “student” network learns to mimic a large “teacher” network by minimizing the loss between the output distributions of the two networks.



Picture from nervanasystems.github.io

Distillation-like domain adaptation

“Knowledge distillation” framework of (Hinton et al., 2014): a smaller “student” network learns to mimic a large “teacher” network by minimizing the loss between the output distributions of the two networks.



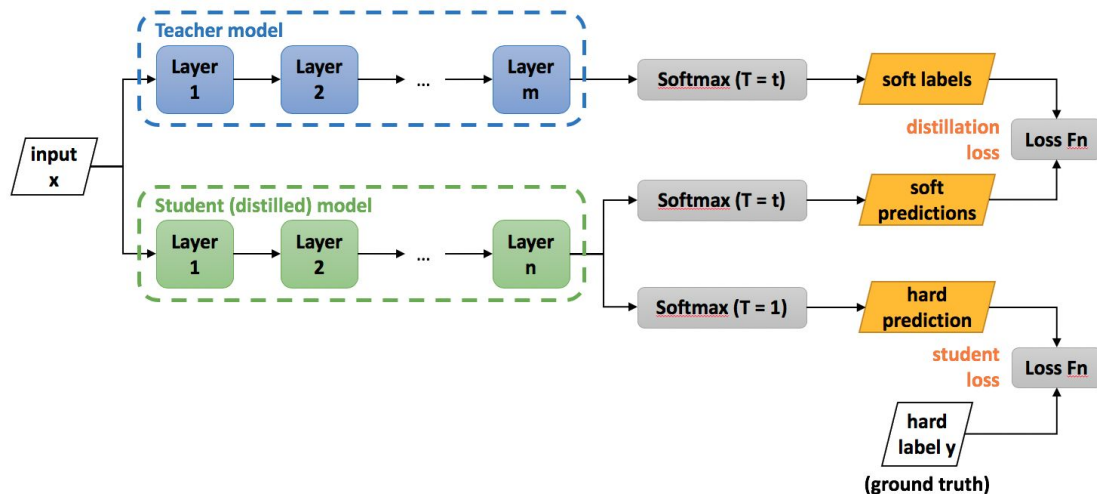
Originally, to compress multiple models (ensemble) to a smaller model.

Picture from nervanasystems.github.io

Distillation-like domain adaptation

[Hinton et al., 2015](#)

“Knowledge distillation” framework of (Hinton et al., 2014): a smaller “student” network learns to mimic a large “teacher” network by minimizing the loss between the output distributions of the two networks.



But this framework can be adapted for domain adaptation.

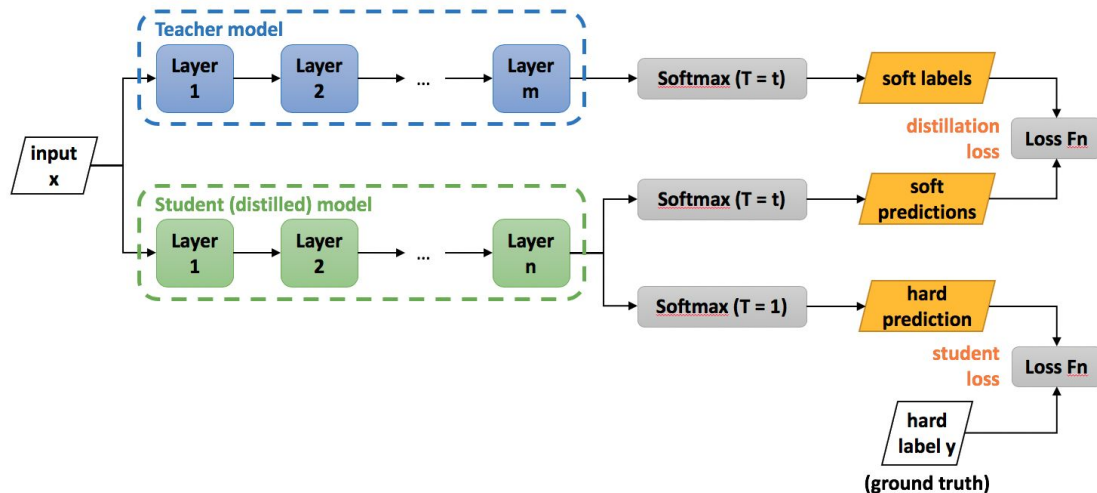
Picture from [nervanasystems.github.io](#)

Distillation-like domain adaptation

[Dakwale et al., 2017](#)

- Train teacher network
- Initialize student network by weights of teacher
- Train student with composite loss

Teacher produces distribution p ,
Student produces distribution q .

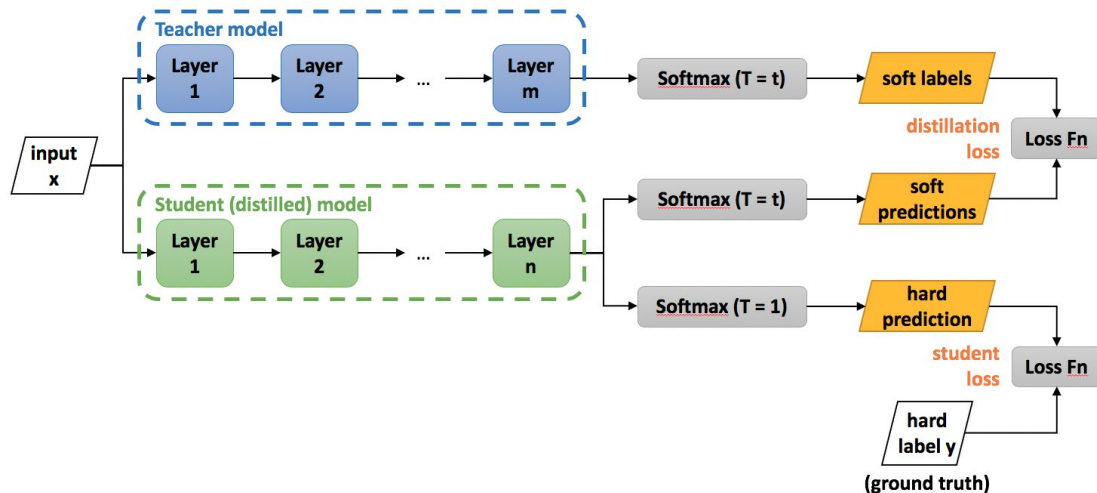


$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

Picture from [nervanasystems.github.io](#)

Distillation-like domain adaptation

- Train teacher network
- Initialize student network by weights of teacher
- Train student with composite loss



Composite loss consists of two terms:

1. Cross-entropy loss on q
2. KL-divergence between p and q

$$L_{total} = (1 - \lambda)L(q_{\theta}) + \lambda D_{KL}(p||q_{\theta})$$

Picture from nervanasystems.github.io

Distillation-like domain adaptation

Composite loss consists of two terms:

1. Cross-entropy loss between hard-label distribution and q
2. KL-divergence between p and q

$$L_{total} = (1 - \lambda)L(q_\theta) + \lambda D_{KL}(p||q_\theta)$$

Distillation-like domain adaptation

Composite loss consists of two terms:

1. Cross-entropy loss between hard-label distribution and q
2. KL-divergence between p and q

$$L_{total} = (1 - \lambda) \boxed{L(q_\theta)} + \lambda D_{KL}(p || q_\theta)$$

$$L(q_\theta) = - \sum_k [l = y_k] \log q_\theta(y_k) = H(l, q_\theta)$$

Distillation-like domain adaptation

Composite loss consists of two terms:

1. Cross-entropy loss between hard-label distribution and q
2. KL-divergence between p and q

$$L_{total} = (1 - \lambda)L(q_\theta) + \lambda D_{KL}(p||q_\theta)$$

$$D_{KL}(p||q_\theta) = \sum_k p(y_k)(\log p(y_k) - \log q_\theta(y_k)) = H(p, q_\theta) - H(p)$$

Distillation-like domain adaptation

Composite loss consists of two terms:

1. Cross-entropy loss between hard-label distribution and q
2. KL-divergence between p and q

$$L_{total} = (1 - \lambda)L(q_\theta) + \lambda D_{KL}(p||q_\theta)$$

$$L_{total} \sim (1 - \lambda)H(l, q_\theta) + \lambda(H(p, q_\theta) - H(p)) \sim (1 - \lambda)H(l, q_\theta) + \lambda H(p, q_\theta)$$

Distillation-like domain adaptation

Using composite loss we force internal representation to be tolerant to domain.
Thus we implicitly minimize discrepancy between domain distributions in internal representation space.

$$\epsilon_T \leq \epsilon_S + \boxed{\frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T)} + \lambda$$

The divergence between
source and target domain

Should we minimize distribution divergence directly?

$$\epsilon_T \leq \epsilon_S + \boxed{\frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T)} + \lambda$$

The divergence between
source and target domain

Should we minimize distribution divergence directly?

Deep distribution alignment! -> adversarial methods, variational bayes methods and other cool stuff.

$$\epsilon_T \leq \epsilon_S + \boxed{\frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T)} + \lambda$$

The divergence between
source and target domain

Batch normalization revisited

[Li et al., 2017](#)

Algorithm 1 Adaptive Batch Normalization (AdaBN)

for neuron j in DNN **do**

Concatenate neuron responses on all images of target domain t : $\mathbf{x}_j = [\dots, x_j(m), \dots]$

Compute the mean and variance of the target domain: $\mu_j^t = \mathbb{E}(\mathbf{x}_j^t)$, $\sigma_j^t = \sqrt{\text{Var}(\mathbf{x}_j^t)}$.

end for

for neuron j in DNN, testing image m in target domain **do**

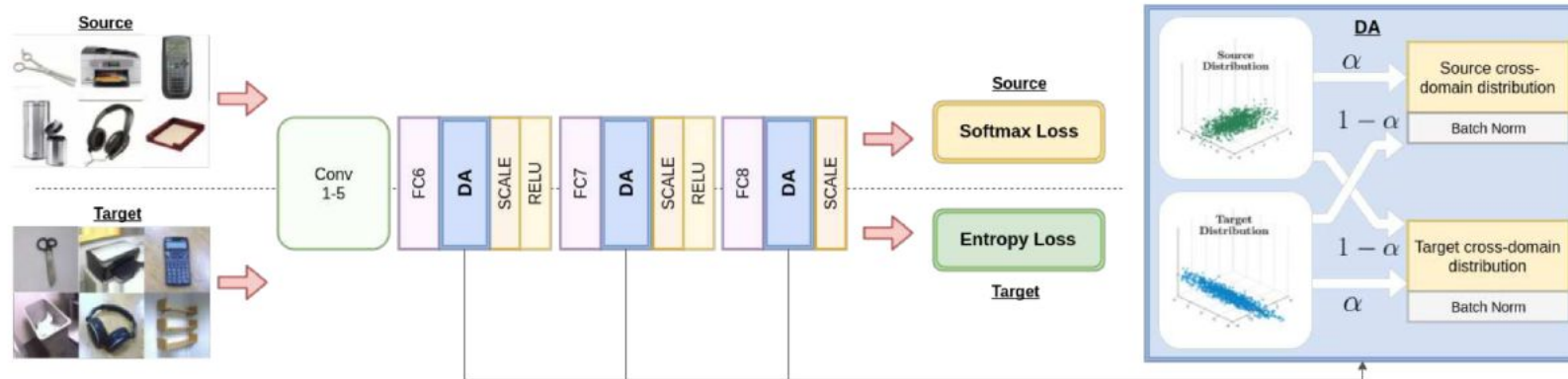
Compute BN output $y_j(m) := \gamma_j \frac{(x_j(m) - \mu_j^t)}{\sigma_j^t} + \beta_j$

end for

$$\text{DA}(x_s; \alpha) = \frac{x_s - \mu_{st, \alpha}}{\sqrt{\epsilon + \sigma_{st, \alpha}^2}}, \quad \text{DA}(x_t; \alpha) = \frac{x_t - \mu_{ts, \alpha}}{\sqrt{\epsilon + \sigma_{ts, \alpha}^2}},$$

Batch normalization revisited

[Carlucci et al., 2017](#)



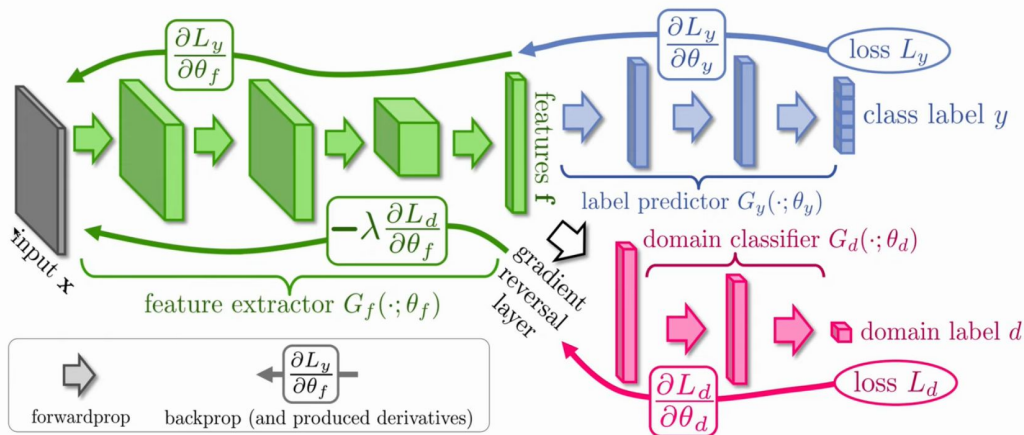
$$DA(x_s; \alpha) = \frac{x_s - \mu_{st, \alpha}}{\sqrt{\epsilon + \sigma_{st, \alpha}^2}}, \quad DA(x_t; \alpha) = \frac{x_t - \mu_{ts, \alpha}}{\sqrt{\epsilon + \sigma_{ts, \alpha}^2}},$$

Assume q^s and q^t to be the distribution of x_s and x_t , respectively, and let $q_{\alpha}^{st} = \alpha q^s + (1 - \alpha) q^t$ and, symmetrically, $q_{\alpha}^{ts} = \alpha q^t + (1 - \alpha) q^s$ be cross-domain distributions mixed by a factor $\alpha \in [0.5, 1]$.

Auto-DIAL

Deep distribution alignment

- Train a classifier for the **domain** of an image based on deep convolutional features.
- Try to maximize the loss of this classifier when training the CNN (**confusion loss**).
- Simultaneously, minimize the classification loss on the source domain using the same convolutional features.
- Train the digit classifier with source domain data, and the domain classifier with both domains' data.



Deep distribution alignment

[Long et al., 2015](#)

Domain adaptation network utilizes Maximum Mean Discrepancy (MMD)

$$d_k^2(p, q) \triangleq \|\mathbf{E}_p[\phi(\mathbf{x}^s)] - \mathbf{E}_q[\phi(\mathbf{x}^t)]\|_{\mathcal{H}_k}^2$$

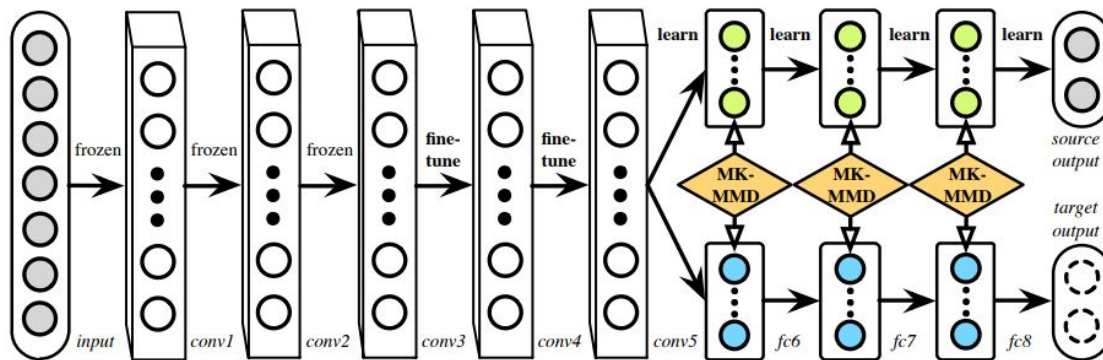


Figure 1. The DAN architecture for learning transferable features.