

Transfer Learning - мотивация из жизни

Мы часто используем в жизни знания в новых ситуациях:

- Шахматы → Шашки
- C++ → Java
- Физика/Математика → Компьютерные науки

***Transfer Learning:** Способность системы распознавать и применять знания и умения, полученные в предыдущих задачах, к новым задачам или данным.*

Различные типы Transfer Learning

тип TL	области	метки данных исходных	метки данных целевых
inductive	многозадачность	есть	есть
	самообучение	нет	есть
transductive	domain adaptation	есть	нет
unsupervised		нет	нет

Передача относительных знаний (relational-knowledge-transfer)

-
- Некоторое соотношение между данными в \mathcal{D}_T и \mathcal{D}_S аналогичны
- Знания, которые передаются, являются этими соотношениями.

Inductive Transfer Learning

Передача примеров и SVM для Inductive TL (5)

- Назначаем веса ρ_i примерам $(\mathbf{x}_i^S, y_i^S) \in \mathcal{D}_S$ в соответствии с их значимостью для $f_T(\cdot)$

$$\min_{w, \xi_i^{(T)}, \xi_i^{(S)}} J = \|w\|^2 + \lambda \sum_{i=1}^{n_T} \xi_i^{(T)} + \lambda \sum_{i=1}^{n_S} \rho_i \xi_i^{(S)}$$

при ограничениях

$$y_i^S \cdot w \cdot \mathbf{x}_i^S \geq 1 - \xi_i^{(S)}, \quad \xi_i^{(S)} \geq 0, \quad i = 1, \dots, n_S,$$

$$y_i^T \cdot w \cdot \mathbf{x}_i^T \geq 1 - \xi_i^{(T)}, \quad \xi_i^{(T)} \geq 0, \quad i = 1, \dots, n_T.$$

- Как определить веса ρ_i ?

Передача параметров и SVM для Inductive TL (2)

- **Параметры:**
 - $w_S = w_0 + v_S$ и $w_T = w_0 + v_T$, где w_S и w_T - параметры SVM для \mathcal{T}_S и \mathcal{T}_T ;
 - w_0 - общие параметры;
 - v_S и v_T - специфичные параметры SVM для \mathcal{T}_S и \mathcal{T}_T .
- **Предположение:** $f_T = w_T \cdot x$
- **Модификация SVM:**

$$\min_{w_0, v_T, \xi_{r_i}} J = \sum_{r \in \{S, T\}} \sum_{i=1}^{n_T} \xi_{r_i} + \frac{\lambda_1}{2} \sum_{r \in \{S, T\}} \|v_r\|^2 + \lambda_2 \|w_0\|^2$$

при ограничениях

$$y_i^T (w_0 + v_T) \cdot x_i^T \geq 1 - \xi_{r_i}, \quad \xi_{r_i} \geq 0, \\ i \in \{1, 2, \dots, n_T\}, \quad r \in \{S, T\}.$$

Передача представления признаков без учит. (4)

- Конструирование признаков
- Для каждой точки (\mathbf{x}_i^T, y_i^T) , вычисляем признаки $\hat{\mathbf{a}}(\cdot) \in \mathbb{R}^d$, решая задачу

$$\min_{\mathbf{a}^{(i)}} \left\| \mathbf{x}_i^T - \sum_{j=1}^s \mathbf{a}_j^{(i)} \mathbf{b}^{(j)} \right\|_2^2 + \beta \left\| \mathbf{a}^{(i)} \right\|_1$$

при ограничениях $\left\| \mathbf{b}^{(j)} \right\|_2 \leq 1$

- Разреженный вектор $\mathbf{a}_j^{(i)}$ - новое реконструированное представление вектора \mathbf{x}_i^T

Передача примеров в Transductive TL (3)

- Веса ρ_i :

$$\min_{w, \xi_i^{(T)}, \xi_i^{(S)}} J = \frac{1}{2} \|w\|^2 + \lambda \sum_{i=1}^n \rho_i \xi_i^{(S)}$$

при ограничениях

$$y_i^S \cdot w \cdot \mathbf{x}_i^S \geq 1 - \xi_i^{(S)}, \quad \xi_i^{(S)} \geq 0, \quad i = 1, \dots, n$$

- например, $\rho_i = \sigma((\mathbf{x}_i^S, y_i^S), \mathcal{D}_T)$, где

$$\sigma((\mathbf{x}_i^S, y_i^S), \mathcal{D}_T) = \frac{1}{|\mathcal{D}_T|} \sum_{j=1}^{|\mathcal{D}_T|} \exp \left\{ -\beta \|\mathbf{x}_i^S - \mathbf{x}_j^T\|^2 \right\}$$

Общий подход для Transductive TL (2)

- Так как нет целевых данных с метками классов, необходимо обучать модель uf исходных данных:

$$\begin{aligned}\theta^* &= \arg \min_{\theta \in \Theta} \sum_{(x,y) \in \mathcal{D}_S} \frac{P(\mathcal{D}_T)}{P(\mathcal{D}_S)} P(\mathcal{D}_S) l(\mathbf{x}, y, \theta) \\ &\approx \arg \min_{\theta \in \Theta} \sum_{i=1}^{n_S} \frac{P_T(\mathbf{x}_i^T, y_i^T)}{P_S(\mathbf{x}_i^S, y_i^S)} l(\mathbf{x}_i^S, y_i^S, \theta) \\ &= \arg \min_{\theta \in \Theta} \sum_{i=1}^{n_S} \frac{P(\mathbf{x}_i^S)}{P(\mathbf{x}_i^T)} l(\mathbf{x}_i^S, y_i^S, \theta)\end{aligned}$$

- Это следует из условия $P(Y_T|X_T) = P(Y_S|X_S)$. Т.о. разность между $P(\mathcal{D}_S)$ и $P(\mathcal{D}_T)$ определяется только $P(X_T)$ и $P(X_S)$ и $\frac{P(\mathbf{x}_i^S)}{P(\mathbf{x}_i^T)}$

Transductive TL и SVM (1)

Необходимо решить три задачи:

- 1 Минимизация функционала риска на области \mathcal{D}_S
- 2 Минимизация разности между двумя совместными распределениями вероятностей J_S и J_t
- 3 Максимизация согласованности маргинальных распределений P_S и P_t

Transductive TL и SVM (2)

$$f = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^{n_s} l(f(\mathbf{x}_i^s), y_i^s) + \sigma \|f\|_K^2 + \lambda D_{f,K}(J_s, J_t) + \gamma M_{f,K}(P_s, P_t)$$

K - ядро

σ, λ, γ - положительные параметры регуляризации
(ограничения на f)

Первая часть - обычный SVM для исходных (source) данных

А что такое $D_{f,K}(J_s, J_t)$ и $M_{f,K}(P_s, P_t)$?

Transductive TL и SVM (3)

Минимизация разности между двумя совместными распределениями вероятностей J_s и J_t или вычисление $D_{f,K}(J_s, J_t)$

Адаптация маргинальных распределений: (используется разность средних значений функций или MMD - maximum mean discrepancy)

$$D_{f,K}(P_s, P_t) = \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_i^s) - \frac{1}{n_t} \sum_{i=1}^{n_t} f(\mathbf{x}_i^t) \right\|_{\mathcal{H}}^2$$

\mathcal{H} определяется $\phi: \mathcal{X} \rightarrow \mathcal{H}$

Transductive TL и SVM (4)

Вычисление $D_{f,K}(J_s, J_t)$

Адаптация условных распределений:

$$D_{f,K}^{(c)}(Q_s, Q_t) = \left\| \frac{1}{n_s^{(c)}} \sum_{i \in \mathcal{D}_s^{(c)}} f(\mathbf{x}_i^s) - \frac{1}{n_t^{(c)}} \sum_{i \in \mathcal{D}_t^{(c)}} f(\mathbf{x}_i^s) \right\|_{\mathcal{H}}^2$$

$\mathcal{D}_s^{(c)}$ - множество примеров из класса s , принадлежащих \mathcal{D}_s

$\mathcal{D}_t^{(c)}$ - множество примеров из класса s , принадлежащих

\mathcal{D}_t , здесь используются псевдо метки классов (примерные)

Transductive TL и SVM (5)

Минимизация разности между двумя совместными распределениями вероятностей J_s и J_t или вычисление $D_{f,K}(J_s, J_t)$

$$D_{f,K}(J_s, J_t) = D_{f,K}(P_s, P_t) + \sum_{=1} D_{f,K}^{()}(Q_s, Q_t)$$

Transductive TL и SVM (6)

- $D_{f,K}(J_s, J_t)$ основана на использовании выборочного мат. ожидания
- Максимизация согласованности маргинальных распределений P_s и P_t основана на использовании выборочной дисперсии:

$$M_{f,K}(P_s, P_t) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_t} (f(\mathbf{x}_i^s) - f(\mathbf{x}_j^t))^2 W_{ij},$$

где

$$W_{ij} = \begin{cases} \cos(\mathbf{x}_i^s, \mathbf{x}_j^t), & \mathbf{x}_i^s \in \mathcal{N}_p(\mathbf{x}_j^t) \vee \mathbf{x}_j^t \in \mathcal{N}_p(\mathbf{x}_i^s) \\ 0, & \text{иначе} \end{cases}$$

- $\mathcal{N}_p(\mathbf{x}_i)$ - множество p ближайших соседей точки \mathbf{x}_i

Transductive TL и SVM (7)

В итоге задача

$$f = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^{n_s} l(f(\mathbf{x}_i^s), y_i^s) + \sigma \|f\|_K^2 + \lambda D_{f,K}(J_s, J_t) + \gamma M_{f,K}(P_s, P_t)$$

сводится к стандартной задаче квадратичного программирования

M.Long, J.Wang, G.Ding, S.J.Pan, P.S.Yu Adaptation Regularization: A general Framework for Transfer Learning. IEEE Trans. on Knowledge and Data Eng., vol. 26(5), pp. 1076-1089, 2014

Методы deep domain adaptation

Методы deep domain adaptation - 3 группы:

- 1 **Discrepancy-based**: подходы основаны на минимизации расстояния между векторными представлениями на исходном и целевом доменах с помощью введения этого расстояния в loss-функцию.
- 2 **Adversarial-Based**: подходы используют состязательную (adversarial) loss-функцию (из GAN), для обучения сети, инвариантной относительно домена.
- 3 **Смешанные методы**: применяют идеи из discrepancy-based семейства, а также self-ensembling, новые слои, loss-функции и т.п.

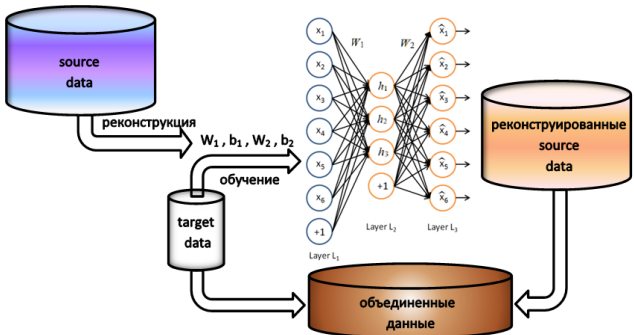
МЕТОДЫ

- **Discrepancy-based**: основан на metric learning, суть которого заключается в обучении такого векторного представления, получаемого из нейронной сети, что представители одного класса будут близки друг к другу в этом представлении по заданной метрике (чаще всего используют L^2 или косинусную метрики).

Автокодер для Transductive TL

- Glorot X, Bordes A, Bengio Y. Domain adaptation for large-scale sentiment classification: A deep learning approach. In: Proceedings of the twenty-eight international conference on machine learning, vol.27. 2011. p.97–110.
 - Chen M, Xu ZE, Weinberger KQ, Sha F (2012) Marginalized denoising autoencoders for domain adaptation. ICML. arXiv preprint arXiv:1206.4683.
- 1 Обучить стек автокодеров на основе исходных и целевых данных без меток классов. Это позволит обнаружить общие инвариантные скрытые признаки.
 - 2 Обучить классификатор, используя преобразованные скрытые признаки, добавив метки исходных данных.

Реконструирование данных



Реконструирование данных

- 1 Обучаем автокодер (веса W_1 , W_2 и параметры b_1 , b_2) на целевых данных
- 2 Для каждого класса из исходных данных \mathbf{x}_k^s на основе обученного автокодера реконструируем:

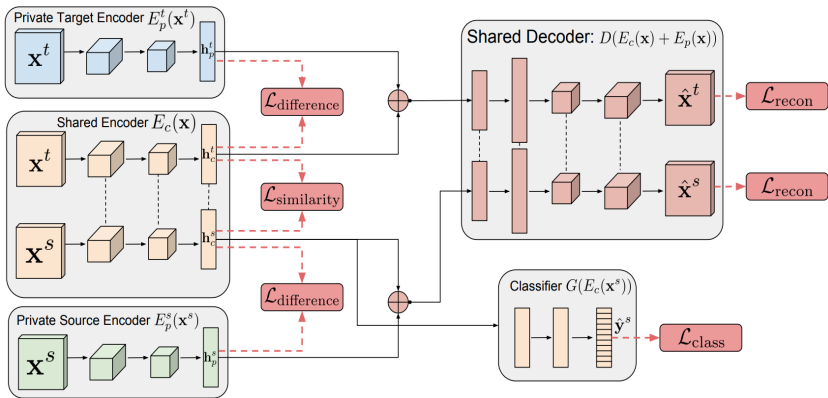
$$\mathbf{x}_k^{s \rightarrow t} = SA_{\text{Recon}}(\mathbf{x}_k^s),$$

где

$$SA_{\text{Recon}}(\mathbf{x}) = \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$

- выход автокодера.

Domain Separation Networks (1)



Domain Separation Networks (2)

- K. Bousmalis et al. Domain Separation Networks // arXiv:1608.06019
- Кодер $E_c(x)$ с shared weights учится захватывать компоненты представления для данной входной выборки, которые являются общими для доменов.
- Частный кодировщик $E_p(x)$ (по одному для каждого домена) учится захватывать специфичные для домена компоненты представления.
- Общий декодер учится восстанавливать входную выборку, используя как частные, так и исходные представления.

Domain Separation Networks (3)

- Частные и общие компоненты представления раздвигаются с помощью мягких ограничений ортогональности подпространства $L_{\text{difference}}$, тогда как общие компоненты представления остаются подобными с потерей сходства $L_{\text{similarity}}$.

$$\mathcal{L} = \mathcal{L}_{\text{task}} + \alpha \mathcal{L}_{\text{recon}} + \beta \mathcal{L}_{\text{difference}} + \gamma \mathcal{L}_{\text{similarity}}$$

- L_{task} обучает модель прогнозировать интересующие нас выходные метки:

$$\mathcal{L}_{\text{task}} = - \sum_{i=0}^{N_s} \mathbf{y}_i^s \cdot \log \hat{\mathbf{y}}_i^s$$

\mathbf{y}_i^s - one-hot кодирование меток классов для source, $\hat{\mathbf{y}}_i^s$
- softmax предсказание модели

Domain Separation Networks (4)

- L_{recon} - функция реконструкции

$$\mathcal{L}_{recon} = \sum_{i=1}^{N_s} \mathcal{L}_{si_mse}(\mathbf{x}_i^s, \hat{\mathbf{x}}_i^s) + \sum_{i=1}^{N_t} \mathcal{L}_{si_mse}(\mathbf{x}_i^t, \hat{\mathbf{x}}_i^t)$$

$$\mathcal{L}_{si_mse}(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{k} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 - \frac{1}{k^2} ([\mathbf{x} - \hat{\mathbf{x}}] \cdot \mathbf{1}_k)^2,$$

k - количество пикселей в \mathbf{x} ; $\mathbf{1}_k$ - единичный вектор длины k

Domain Separation Networks (5)

- $L_{\text{difference}}$ - функция потерь различий способствует ортогональности между общим и частным представлениями каждого домена:

$$L_{\text{difference}} = \left\| \mathbf{H}_c^s \top \mathbf{H}_p^s \right\|_F^2 + \left\| \mathbf{H}_c^t \top \mathbf{H}_p^t \right\|_F^2$$

- $L_{\text{similarity}}$:

$$\begin{aligned} \mathcal{L}_{\text{similarity}}^{\text{MMD}} &= \frac{1}{(N^s)^2} \sum_{i,j=0}^{N^s} \kappa(\mathbf{h}_{ci}^s, \mathbf{h}_{cj}^s) \\ &- \frac{2}{N^s N^t} \sum_{i,j=0}^{N^s, N^t} \kappa(\mathbf{h}_{ci}^s, \mathbf{h}_{cj}^t) + \frac{1}{(N^t)^2} \sum_{i,j=0}^{N^t} \kappa(\mathbf{h}_{ci}^t, \mathbf{h}_{cj}^t) \end{aligned}$$

Adversarial-Based методы

- Обучение нейронной сети с инвариантным по отношению к исходному и целевому доменам векторным представлением.
- Обученную сеть на размеченном source domain можно будет использовать на target domain, в идеале — практически без потери качества классификации.

Алгоритм DANN

- Y. Ganin et al. Domain-Adversarial Training of Neural Networks // arXiv:1505.07818
- Три части алгоритма:
 - Основная сеть для получения векторного представления (feature extractor);
 - "Голова отвечающая за классификацию на исходном домене;
 - "Голова которая обучается отличать данные из исходного домен от целевого.

Алгоритм ADDA (1)

- E. Tzeng et al. Adversarial Discriminative Domain Adaptation // arXiv:1702.05464
- Подразумевает разделение сети для исходного домена и сети для целевого домена
- Шаги алгоритма:
- Классифицирующая сеть обучается на исходном домене. Её векторное представление - M_s , а \mathbf{X}_s - исходный домен.
- Инициализируем нейронную сеть для целевого домена с помощью обученной сети из предыдущего шага. Обозначим её M_t , а \mathbf{X}_t - целевой домен.

Алгоритм ADDA (2)

- Перейдем к adversarial-тренировке: будем обучать дискриминатор D при фиксированных M_s и M_t с помощью следующей целевой функции:

$$\min_D L_{adv_D}(\mathbf{X}_s, \mathbf{X}_t, M_s, M_t) = -\mathbb{E}_{x_s \sim \mathbf{X}_s} [\log D(M_s(x_s))] - \mathbb{E}_{x_t \sim \mathbf{X}_t} [\log(1 - D(M_t(x_t)))]$$

- Заморозим дискриминатор и дообучим M_t на целевом домене:

$$\min_{M_s, M_t} L_{adv_M}(\mathbf{X}_s, \mathbf{X}_t, D) = -\mathbb{E}_{x_t \sim \mathbf{X}_t} [\log D(M_t(x_t))]$$

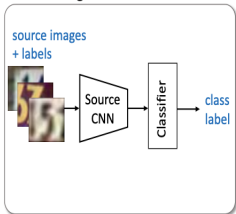
- Шаги 3 и 4 повторяются несколько раз.

Алгоритм ADDA (3)

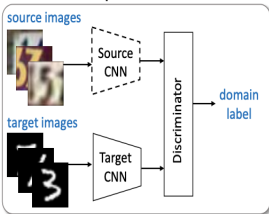
- Суть ADDA заключается в том, что мы сначала обучаем хороший классификатор на размеченном исходном домене, а затем с помощью adversarial-обучения адаптируем так, чтобы векторные представления классификатора на обоих доменах были близки.

Алгоритм ADDA (4)

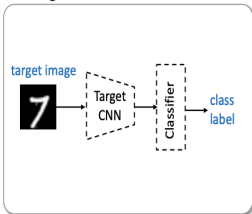
Pre-training



Adversarial Adaptation



Testing



Negative Transfer

- Происходит, когда передача знаний из \mathcal{D}_S и \mathcal{T}_S приводит к снижению качества \mathcal{T}_T
- Если задачи \mathcal{T}_S и \mathcal{T}_T слишком различны, тогда передача “в лоб” может привести к снижению качества \mathcal{T}_T
- Важно проанализировать связанность \mathcal{T}_S и \mathcal{T}_T или \mathcal{D}_S и \mathcal{D}_T , определить критерий схожести

Методы:

- Схожесть \mathcal{T}_S и \mathcal{T}_T определяется на основе схожести между распределениями вероятностей примеров
- Схожесть \mathcal{T}_S и \mathcal{T}_T определяется на основе введения характеристик задач более высокого уровня, например, признаков, которые известны заранее

Negative Transfer

- Схожесть T_S и T_T определяется на основе введения характеристик задач более высокого уровня, например, признаков, которые известны заранее
- Например, признак - пол

Мужчины



Женщины

