

Data Science

Apprentissage par transfert

Hachem Kadri

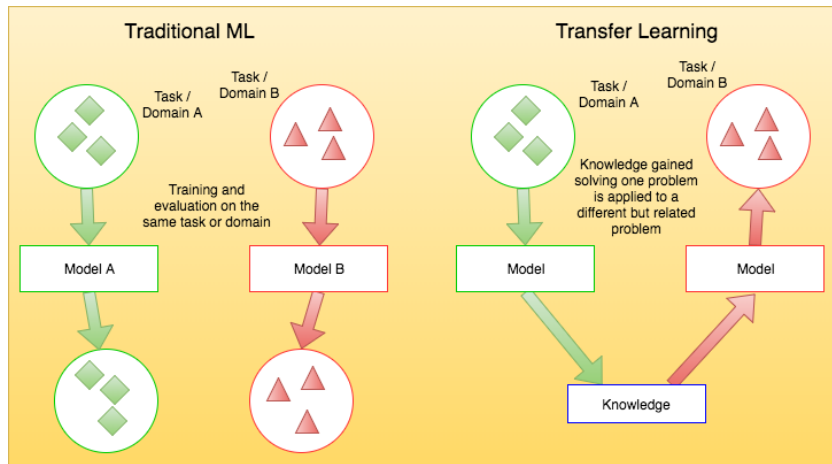
2018-2019

Apprentissage par transfert

- ▶ Dans des contextes réels, le meilleur jeu de données pour la tâche d'apprentissage n'est pas toujours disponible



Apprentissage par transfert



from <http://experiencesutra.com/experiments/deep-learning-in-fashion>

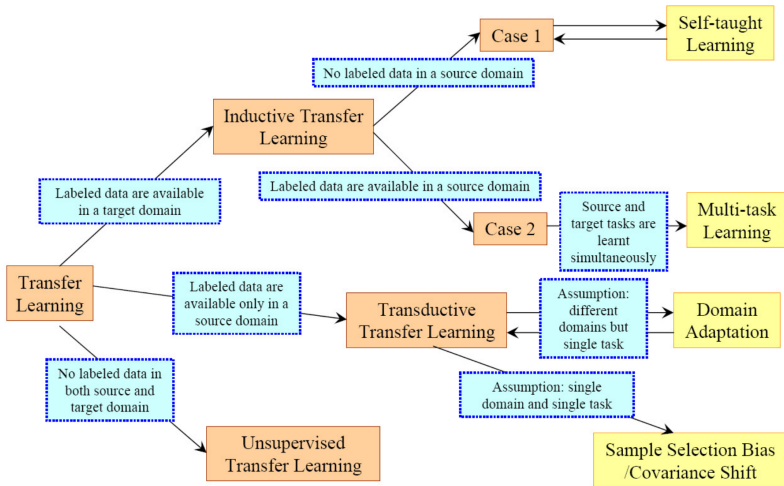
Apprentissage par transfert

- ▶ Un domaine source \mathcal{D}_S et une tâche source \mathcal{T}_S
- ▶ Un domaine cible \mathcal{D}_T et une tâche source \mathcal{T}_T

Apprentissage par transfert a pour objectif de d'améliorer l'apprentissage de la fonction cible utilisant les connaissances de \mathcal{D}_S et \mathcal{T}_S dans les cas où $\mathcal{D}_S \neq \mathcal{D}_T$ et/ou $\mathcal{T}_S \neq \mathcal{T}_T$.

- ▶ Généralement, d'autres contraintes sur le domaine cible sont imposées, par exemple peu ou pas d'étiquettes
- ▶ Dans le cas où $\mathcal{D}_S = \mathcal{D}_T$ et/ou $\mathcal{T}_S = \mathcal{T}_T$, on retrouve la cas classique d'apprentissage supervisé

Apprentissage par transfert



"survey on Transfer Learning" [Pan and Yang, TKDE 2010]

Apprentissage par transfert

- ▶ Même tâche, domaines différents \Rightarrow Adaptation de domaine
- ▶ Même domaine, plusieurs tâches \Rightarrow Apprentissage multi-tâche

Adaptation de domaine

Adaptation de domaine

Training Domain



- **Appearance Change** - - - - -



Application Domains

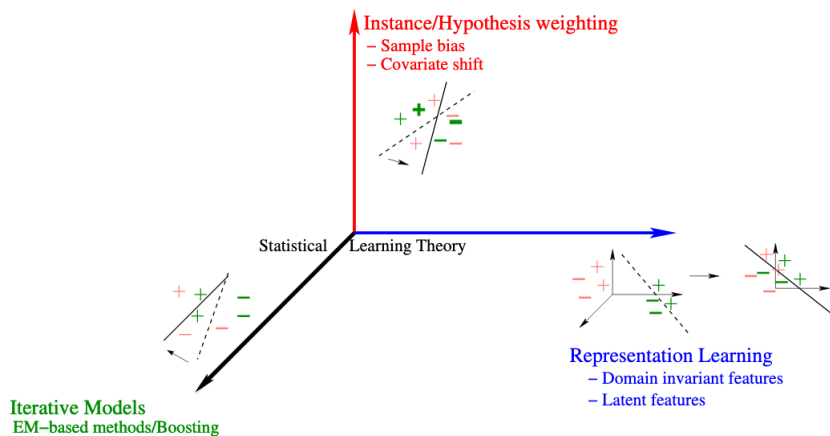
Adaptation de domaine

- ▶ un espace de features \mathcal{X} , un ensemble de labels $\mathcal{Y} = \{-1, 1\}$
- ▶ une distribution source P_s , une distribution cible P_t
- ▶ une fonction inconnue $f : \mathcal{X} \rightarrow \mathcal{Y}$ selon $P_t(y|x)$
- ▶ un ensemble de données d'apprentissage $S_s = \{(x_i, y_i)\}$ généré i.i.d selon P_s , un ensemble de données test $S_t = \{x_i\}$ généré selon D_t , la loi marginale de P_t sur \mathcal{X}
- ▶ apprendre un classifieur h le plus proche possible de la fonction inconnue f
- ▶ risque réel (source) : $R_s(h) = E_{(x,y) \sim P_s} \mathbb{I}[h(x) \neq y]$
risque empirique sur S_s : $\hat{R}_s(h) = \sum_{(x,y) \in S_s} \mathbb{I}[h(x) \neq y]$
- ▶ risque réel (cible) : $R_t(h) = E_{(x,y) \sim P_t} \mathbb{I}[h(x) \neq y]$

Borne de généralisation : $R_s(h) \leq \hat{R}_s(h) + \sqrt{\frac{\text{complexite}(h \in \mathcal{H})}{|S_s|}}$

\Rightarrow garanties théoriques pour le domaine cible ?

Adaptation de domaine - Stratégies



from https://limos.fr/media/uploads/seminaire/hebrard_slides_Nov18.pdf

Adaptation de domaine - par repondération

$$\begin{aligned}\epsilon_T(h) &= \mathbf{E}_{(\mathbf{x}^t, y^t) \sim P_T} \mathbf{I}[h(\mathbf{x}^t) \neq y^t] \\ &= \mathbf{E}_{(\mathbf{x}^t, y^t) \sim P_T} \frac{P_S(\mathbf{x}^t, y^t)}{P_S(\mathbf{x}^t, y^t)} \mathbf{I}[h(\mathbf{x}^t) \neq y^t] \\ &= \sum_{(\mathbf{x}^t, y^t)} P_T(\mathbf{x}^t, y^t) \frac{P_S(\mathbf{x}^t, y^t)}{P_S(\mathbf{x}^t, y^t)} \mathbf{I}[h(\mathbf{x}^t) \neq y^t] \\ &= \mathbf{E}_{(\mathbf{x}^t, y^t) \sim P_S} \frac{P_T(\mathbf{x}^t, y^t)}{P_S(\mathbf{x}^t, y^t)} \mathbf{I}[h(\mathbf{x}^t) \neq y^t]\end{aligned}$$

Assume similar tasks - covariate shift, $P_S(y|\mathbf{x}) = P_T(y|\mathbf{x})$, then:

$$\begin{aligned}&= \mathbf{E}_{(\mathbf{x}^t, y^t) \sim P_S} \frac{D_T(\mathbf{x}^t) P_T(y^t|\mathbf{x}^t)}{D_S(\mathbf{x}^t) P_S(y^t|\mathbf{x}^t)} \mathbf{I}[h(\mathbf{x}^t) \neq y^t] \\ &= \mathbf{E}_{(\mathbf{x}^t, y^t) \sim P_S} \frac{D_T(\mathbf{x}^t)}{D_S(\mathbf{x}^t)} \mathbf{I}[h(\mathbf{x}^t) \neq y^t]\end{aligned}$$

from https://limos.fr/media/uploads/seminaire/hebrard_slides_Nov18.pdf

Adaptation de domaine - par pondération

- ▶ covariate shift

$$R_t(h) = E_{(x^t, y^t) \sim P_s} \frac{D_t(x^t)}{D_s(x^t)} \mathbb{I}[h(x^t) \neq y^t]$$

- ▶ erreur de pondération dans le domaine source : $\omega(x^t) = \frac{D_t(x^t)}{D_s(x^t)}$
- ▶ repondérer les données sources étiquetées :

$$\sum_{(x, y) \in S_s} \hat{\omega}(x_i^s) \mathbb{I}[h(x^s) \neq y^s],$$

avec $\hat{\omega}$ une estimation de ω

Adaptation de domaine - par repondération

Estimateurs de densité

Construire des estimateurs de densité pour les domaines source et cible et estimer le ratio

- ▶ $\hat{\omega}(x) = \sum_I \alpha_I \psi_I(x)$
- ▶ apprendre α : $\arg \min_{\alpha} KL(\hat{\omega}D_s, D_t)$

Kernel mean matching

Matcher les distributions avec des noyaux

- ▶ $MMD(P_s, P_t) = \|\frac{1}{n_s} \sum_{i=1}^{n_s} \phi(x_i^s) - \frac{1}{n_t} \sum_{i=1}^{n_t} \phi(x_i^t)\|_{\mathcal{H}}$
- ▶ $\min_{\beta} MMD(\beta(x)P_s(x), P_t) \text{ s.t. } \beta(x) > 0 \text{ et } E_{P_s}[\beta(x)] = 1$

Apprentissage multi-tâche

Apprentissage multi-tâche

- ▶ Apprendre plusieurs tâches simultanément
- ▶ ... au lieu d'apprendre les tâches indépendamment
- ▶ Améliorer l'apprentissage et la prédiction en tenant en compte les dépendances entre les tâches

Apprentissage supervisé

- ▶ Classification $f : X \rightarrow Y \subseteq \{-1, 1\}$
- ▶ Régression $f : X \rightarrow Y = \mathbb{R}$

Autres tâches ?

- ▶ Multi-classe $f : X \rightarrow Y = \{1, 2, \dots, T\}$
- ▶ Régression à valeur vectorielle $f : X \rightarrow Y \subseteq \mathbb{R}^T$
- ▶ ...

Apprentissage supervisé multi-tâche

Données d'apprentissage

$$S_1 = (x_i^1, y_i^1)_{i=1}^{n_1}, \dots, S_T = (x_i^T, y_i^T)_{i=1}^{n_T}$$

Apprendre

$$f_1 : X_1 \rightarrow Y_1, \dots, f_T : X_T \rightarrow Y_T$$

- ▶ Régression à valeur vectorielle

$$S_n = (x_i, y_i)_{i=1}^n, \quad x_i \in X, \quad y_i \in \mathbb{R}^T$$

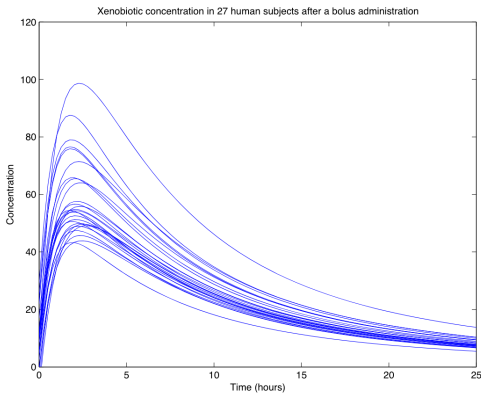
Multi-tâche avec les mêmes données d'entrée ! Chaque composante du vecteur de sortie est une "tâche"

- ▶ Multi-classe

$$S_n = (x_i, y_i)_{i=1}^n, \quad x_i \in X, \quad y_i \in \{1, \dots, T\}$$

Pourquoi apprendre plusieurs tâches simultanément ?

- Cas de la régression multiple

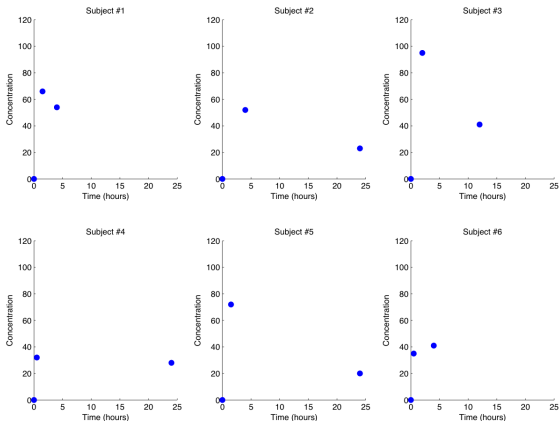


from [Dinuzzo et al., 2013]

- ▶ Les courbes de réponses ont des allures similaires
- ▶ Il y a une variabilité inter-individuelle

Pourquoi apprendre plusieurs tâches simultanément ?

- Cas de la régression multiple

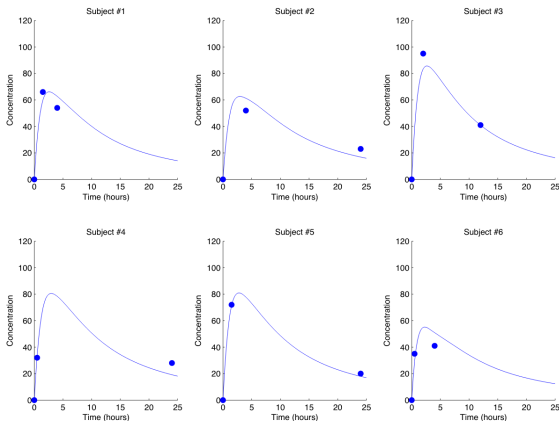


from [Dinuzzo et al., 2013]

- ▶ Peu de données par sujet et différents échantillonnages

Pourquoi apprendre plusieurs tâches simultanément ?

- Cas de la régression multiple



from [Dinuzzo et al., 2013]

- ▶ Combiner les données pour mieux estimer toutes les courbes

Apprentissage multi-tâche régularisé

$$\text{err}(w_1, \dots, w_T) + \text{pen}(w_1, \dots, w_T)$$

On considère des modèles linéaires

$$f_1(x) = w_1^\top x, \dots, f_T(x) = w_T^\top x$$

Erreur empirique

$$\text{err}(w_1, \dots, w_T) = \sum_{t=1}^T \frac{1}{n_t} \sum_{i=1}^{n_t} (y_i^t - w_t^\top x_i^t)^2$$

- ▶ Il est possible de consider d'autres fonctions coûts

Erreur empirique

$$\text{err}(w_1, \dots, w_T) = \sum_{t=1}^T \frac{1}{n_t} \sum_{i=1}^{n_t} (y_i^t - w_t^\top x_i^t)^2$$

- Cas de la régression à valeur vectorielle

$$S_n = (x_i, y_i)_{i=1}^n, \quad x_i \in \mathcal{X}, \quad y_i \in \mathbb{R}^T$$

$$\text{err}(w_1, \dots, w_T) = \frac{1}{n} \sum_{t=1}^T \sum_{i=1}^n (y_i^t - w_t^\top x_i)^2 = \frac{1}{n} \|Y - XW\|^2$$

► $Y = (y_i)_{i=1}^n \in \mathbb{R}^{n \times T}$, $X = (x_i)_{i=1}^n \in \mathbb{R}^{n \times d}$, $W \in \mathbb{R}^{d \times T}$

Régularisation multi-tâche

$$\text{pen}(w_1, \dots, w_T)$$

- ▶ Coupler les solutions obtenues pour chaque tâche par régularisation
- ▶ Exploiter la structure entre les tâches

Régularisation multi-tâche

$$\text{pen}(w_1, \dots, w_T) = \sum_{t=1}^T \|w_t\|^2$$

- ▶ Indépendance entre tâches !

$$\begin{aligned} \min_{w_1, \dots, w_T} \frac{1}{n} \sum_{t=1}^T \sum_{i=1}^n (y_i^t - w_t^\top x_i)^2 + \lambda \sum_{t=1}^T \|w_t\|^2 \\ = \sum_{t=1}^T \left(\min_{w_t} \sum_{i=1}^n (y_i^t - w_t^\top x_i)^2 + \lambda \|w_t\|^2 \right) \end{aligned}$$

Régularisation multi-tâche

$$\text{pen}(w_1, \dots, w_T) = \text{rank}(W)$$

- ▶ Régression de faible rang

$$\begin{aligned} \min_{W \in \mathbb{R}^{d \times T}} \quad & \frac{1}{n} \|Y - XW\|^2 \\ \text{s.t.} \quad & \text{rank}(W) \leq r \end{aligned}$$

Régression de faible rang

$$\hat{W}^* = \min_{W \in \mathbb{R}^{d \times T}} \frac{1}{n} \|Y - XW\|^2$$

s.t. $\text{rank}(W) \leq r$

1. Calculer la solution de la régression par moindres carrés (sans contraintes)

$$W^* = (X^T X)^{-1} X^T Y$$

2. Calculer l'estimation par moindres carrés

$$Y^* = XW^*$$

3. Calculer la SVD de Y^* et la matrice de projection P_r

$$Y^* = UDV, \quad P_r = \sum_{i=1}^r v_i v_i^T$$

4. Obtenir la solution de la régression de faible rang

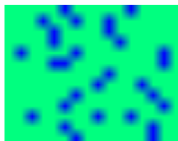
$$\hat{W}^* = W^* P_r$$

Régularisation multi-tâche

$$\text{pen}(w_1, \dots, w_T)$$

► Régularisation matricielle

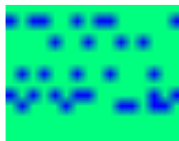
- rang
- norme nucléaire : $\|W\|_* = \text{trace}(\sqrt{W^T W})$
- norme $L_{2,1}$: $\|W\|_{2,1} = \sum_{j=1}^d \sqrt{\sum_{t=1}^T W_{jt}^2}$
- ...



#rows = 13

$\|\cdot\|_{2,1} = 19$

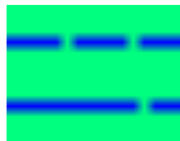
ℓ_1 -norm = 29



5

12

29



2

8

29