

# Apprentissage de vote de majorité pour la classification supervisée et l'adaptation de domaine

approches PAC-Bayésiennes et combinaison de similarités

Emilie Morvant

Laboratoire d'Informatique Fondamentale, QARMA Group, Aix\*Marseille Université, France



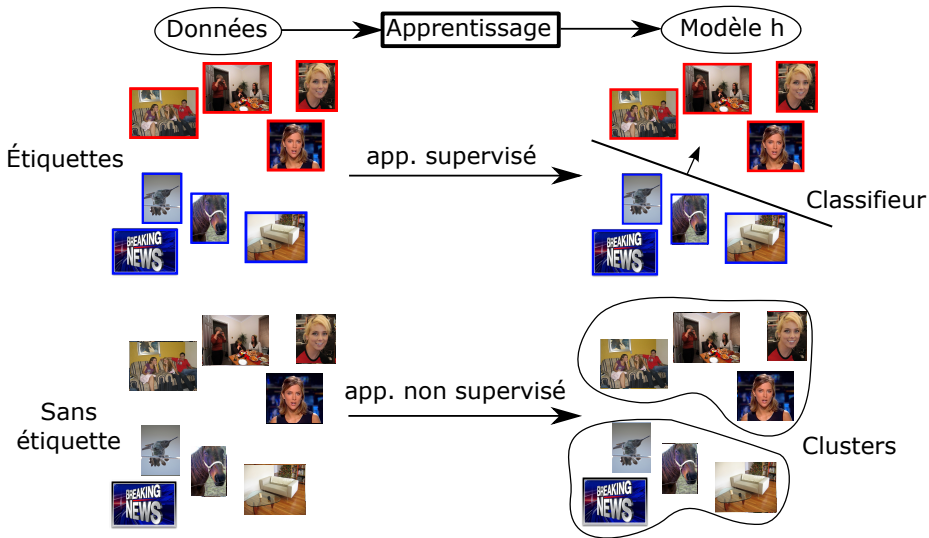
Rapporteurs : Michèle Sebag, Mario Marchand  
Examineurs : Antoine Cornuéjols, Rémi Gilleron, Liva Ralaivola  
Directeur : Amaury Habrard  
Co-directeur : Stéphane Ayache

Soutenance de thèse  
18 septembre 2013

# Contexte de la thèse

Apprentissage automatique

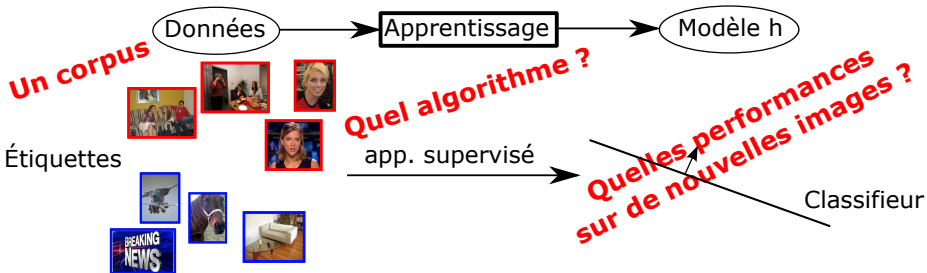
**Tâche :** Y-a-t-il une Personne dans l'image/la vidéo ?



# Contexte de la thèse

Apprentissage automatique

**Tâche :** Y-a-t-il une Personne dans l'image/la vidéo ?

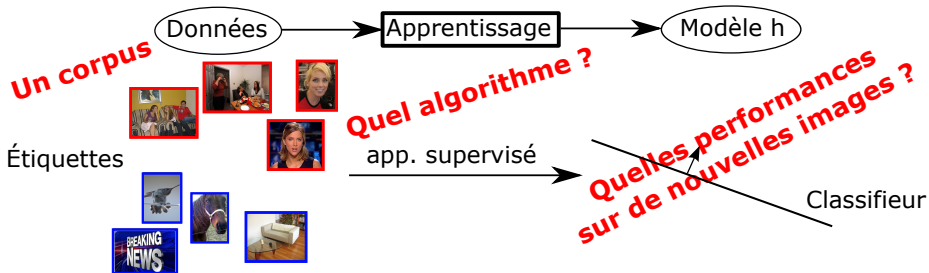


**Comment apprendre  $h$  pour qu'il se trompe le moins possible sur de nouvelles images ?**

# Contexte de la thèse

Apprentissage automatique

**Tâche** : Y-a-t-il une Personne dans l'image/la vidéo ?



**Comment apprendre  $h$  pour qu'il se trompe le moins possible sur de nouvelles images ?**

**Solution** : minimiser l'erreur empirique mesurée sur les données d'apprentissage

⇒ Requiert des garanties ⇒ Borne en généralisation

$$\text{erreur réelle} \leq \text{erreur empirique} + f(\text{complexité, nb de données})$$

### Problématiques

- ❶ Comment tirer bénéfice de différentes descriptions ? *ex : son, texte, image, couleur, ...*
- ❷ Plusieurs objets différents à reconnaître. *ex : personne, canapé, cheval, ...*
- ❸ Les nouvelles images proviennent d'un corpus différent. *ex : photos → vidéos*

### Problématiques

- 1 Comment tirer bénéfice de différentes descriptions? *ex : son, texte, image, couleur, ...*  
⇒ **multimodalité, multivue, combinaisons/fusion de modèles** ⇒ **vote de majorité**
- 2 Plusieurs objets différents à reconnaître. *ex : personne, canapé, cheval, ...*
- 3 Les nouvelles images proviennent d'un corpus différent. *ex : photos → vidéos*

### Problématiques

- ❶ Comment tirer bénéfice de différentes descriptions? *ex : son, texte, image, couleur, ...*  
⇒ **multimodalité, multivue, combinaisons/fusion de modèles** ⇒ **vote de majorité**

Qu'est ce qu'un vote de majorité?

$\mathcal{H}$  : un ensemble de modèles/votants retournant  $-1$  ou  $+1$

**Objectif** : Construire un vote de majorité pondéré sur  $\mathcal{H}$  :  $\text{sign} \left[ \sum_{h \in \mathcal{H}} \rho(h) h(\mathbf{x}) \right]$

**Question** : Comment apprendre les poids  $\rho(h)$ ?

- ❷ Plusieurs objets différents à reconnaître. *ex : personne, canapé, cheval, ...*
- ❸ Les nouvelles images proviennent d'un corpus différent. *ex : photos  $\rightarrow$  vidéos*

### Problématiques

- 1 Comment tirer bénéfice de différentes descriptions? *ex* : *son, texte, image, couleur, ...*  
⇒ **multimodalité, multivue, combinaisons/fusion de modèles** ⇒ **vote de majorité**

Qu'est ce qu'un vote de majorité?

$\mathcal{H}$  : un ensemble de modèles/votants retournant  $-1$  ou  $+1$

**Objectif** : Construire un vote de majorité pondéré sur  $\mathcal{H}$  :  $\text{sign} \left[ \sum_{h \in \mathcal{H}} \rho(h) h(\mathbf{x}) \right]$

**Question** : Comment apprendre les poids  $\rho(h)$ ?

- 2 Plusieurs objets différents à reconnaître. *ex* : *personne, canapé, cheval, ...*  
⇒ **classification multiclasse ou multilabel**
- 3 Les nouvelles images proviennent d'un corpus différent. *ex* : *photos* → *vidéos*



### Problématiques

- ❶ Comment tirer bénéfice de différentes descriptions? *ex* : *son, texte, image, couleur, ...*  
⇒ **multimodalité, multivue, combinaisons/fusion de modèles** ⇒ **vote de majorité**

Qu'est ce qu'un vote de majorité?

$\mathcal{H}$  : un ensemble de modèles/votants retournant  $-1$  ou  $+1$

**Objectif** : Construire un vote de majorité pondéré sur  $\mathcal{H}$  :  $\text{sign} \left[ \sum_{h \in \mathcal{H}} \rho(h) h(\mathbf{x}) \right]$

**Question** : Comment apprendre les poids  $\rho(h)$ ?

- ❷ Plusieurs objets différents à reconnaître. *ex* : *personne, canapé, cheval, ...*  
⇒ **classification multiclasse ou multilabel**
- ❸ Les nouvelles images proviennent d'un corpus différent. *ex* : *photos* → *vidéos*  
⇒ **on doit adapter le modèle** ⇒ **adaptation de domaine**

## Classification supervisée et théorie PAC-Bayésienne

- Vote de majorité contraint et classification binaire
  - ▶ Application à des classifieurs de type  $k$  plus proches voisins (CAp'13)
  - ▶ Spécialisation à la fusion de classifieurs en multimédia
- Classification multiclasse
  - ▶ Borne PAC-Bayésienne sur la confusion du classifieur de Gibbs (ICML'12, CAp'12)
  - ▶ Borne sur le risque du vote de majorité pondéré

## Adaptation de domaine

- Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes (ICDM'11, CAp'11, KAIS'12)
- Analyse PAC-Bayésienne de l'adaptation de domaine (ICML'13, CAp'13)

## Classification supervisée et théorie PAC-Bayésienne

- Vote de majorité contraint et classification binaire
  - ▶ Application à des classifieurs de type  $k$  plus proches voisins (CAp'13)
  - ▶ Spécialisation à la fusion de classifieurs en multimédia
- Classification multiclasse
  - ▶ Borne PAC-Bayésienne sur la confusion du classifieur de Gibbs (ICML'12, CAp'12)
  - ▶ Borne sur le risque du vote de majorité pondéré

## Adaptation de domaine

- Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes (ICDM'11, CAp'11, KAIS'12)
- Analyse PAC-Bayésienne de l'adaptation de domaine (ICML'13, CAp'13)

- 1 La théorie de l'adaptation de domaine
- 2 Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes
- 3 Analyse PAC-Bayésienne de l'adaptation de domaine
- 4 Conclusion et perspectives générales

- 1 La théorie de l'adaptation de domaine
- 2 Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes
- 3 Analyse PAC-Bayésienne de l'adaptation de domaine
- 4 Conclusion et perspectives générales

# La théorie de l'adaptation de domaine

Motivation

Quand a-t-on besoin d'adaptation de domaine (DA) ?

Lorsque la distribution **d'apprentissage** diffère de la distribution **de test**

Exemple



Personne pas de Personne

Y-a-t'il une Personne ?

Corpus de **Photos** étiquetées

Corpus de **Videos** non étiquetées

Domaine **source**

Domaine **cible**

⇒ Comment apprendre, à partir d'une **distribution**,  
un classifieur performant sur une **distribution différente** ?

# La théorie de l'adaptation de domaine

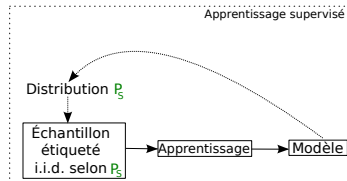
## Description du problème et notations

$X \in \mathbb{R}^d$	espace d'entrée
$Y = \{-1, +1\}$	espace de sortie
$\mathcal{H}$	ensemble de classifieurs
$P_S$ domaine source $P_T$ domaine cible	} distributions sur $X \times Y$
$D_S, D_T$	

# La théorie de l'adaptation de domaine

## Description du problème et notations

$X \in \mathbb{R}^d$	espace d'entrée
$Y = \{-1, +1\}$	espace de sortie
$\mathcal{H}$	ensemble de classifieurs
$P_S$ domaine source	} distributions sur $X \times Y$
$P_T$ domaine cible	
$D_S, D_T$	distributions marginales sur $X$



## Classification supervisée

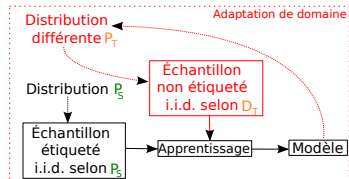
**Objectif :** Trouver  $h \in \mathcal{H}$  minimisant l'**erreur source** :  $R_{P_S}(h) = \mathbf{E}_{(x^s, y^s) \sim P_S} \mathbf{I}[h(x^s) \neq y^s]$



# La théorie de l'adaptation de domaine

## Description du problème et notations

$X \in \mathbb{R}^d$	espace d'entrée
$Y = \{-1, +1\}$	espace de sortie
$\mathcal{H}$	ensemble de classifieurs
$P_S$ domaine source	} distributions sur $X \times Y$
$P_T$ domaine cible	
$D_S, D_T$	distributions marginales sur $X$



## Classification supervisée

**Objectif :** Trouver  $h \in \mathcal{H}$  minimisant **l'erreur source** :  $R_{P_S}(h) = \mathbf{E}_{(x^s, y^s) \sim P_S} \mathbf{I}[h(x^s) \neq y^s]$

## Adaptation de domaine

### Contexte

$S = \{(x_i^s, y_i^s)\}_{i=1}^{m_s}$  Échantillon **source** tiré *i.i.d.* selon  $P_S$

$T = \{x_i^t\}_{i=1}^{m_t}$  Échantillon **cible** tiré *i.i.d.* selon  $D_T$

**Objectif :** Trouver  $h \in \mathcal{H}$  minimisant **l'erreur cible** :  $R_{P_T}(h) = \mathbf{E}_{(x^t, y^t) \sim P_T} \mathbf{I}[h(x^t) \neq y^t]$

# Contexte de la thèse

## Adaptation de domaine

### Trois grands types d'approches

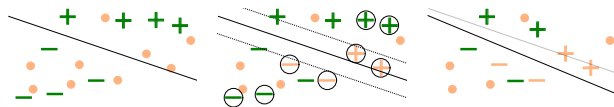
#### Repondérations des données

[Huang et al., 2007, Jiang and Zhai, 2007, Mansour et al., 2009b]



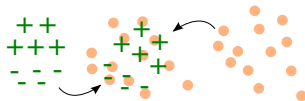
#### Auto-étiquetage

[Bruzzone and Marconcini, 2010]



#### Recherche d'un espace de représentation commun

[Blitzer et al., 2006, Blitzer et al., 2011, Chen et al., 2011, Daumé III, 2007, Daumé III et al., 2010]

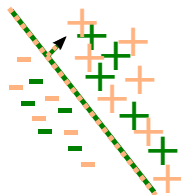
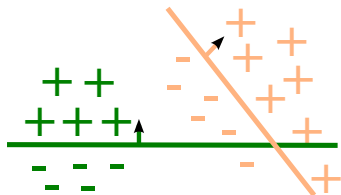


# La théorie de l'adaptation de domaine

Nécessité d'une mesure de divergence entre les domaines

## Problématique principale

Si  $h$  est appris depuis le domaine **source**,  
quelle sera sa performance sur le domaine **cible** ?



⇒ Si les domaines sont “**proches**”  
alors un classifieur d'erreur **source** faible peut être un classifieur d'erreur **cible** faible

# La théorie classique de l'adaptation de domaine

Les résultats de S. Ben-David *et al.* et Mansour *et al.*

Théorème classique [Ben-David *et al.*, 2010, Mansour *et al.*, 2009a]

Soit  $\mathcal{H}$  un espace d'hypothèses. Si  $D_S$  et  $D_T$  sont deux distributions sur  $X$ , alors :

$$\forall h \in \mathcal{H}, \overbrace{R_{P_T}(h)}^{\text{erreur cible}} \leq R_{P_S}(h) + \frac{1}{2}d_{\mathcal{H}}(D_S, D_T) + \nu$$

# La théorie classique de l'adaptation de domaine

Les résultats de S. Ben-David *et al.* et Mansour *et al.*

Théorème classique [Ben-David *et al.*, 2010, Mansour *et al.*, 2009a]

Soit  $\mathcal{H}$  un espace d'hypothèses. Si  $D_S$  et  $D_T$  sont deux distributions sur  $X$ , alors :

$$\forall h \in \mathcal{H}, \quad \overbrace{R_{P_T}(h)}^{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \frac{1}{2}d_{\mathcal{H}}(D_S, D_T) + \nu$$

$R_{P_S}(h)$  : erreur classique sur le domaine source

Minimisable via une méthode de classification supervisée sans adaptation

# La théorie classique de l'adaptation de domaine

Les résultats de S. Ben-David *et al.* et Mansour *et al.*

Théorème classique [Ben-David *et al.*, 2010, Mansour *et al.*, 2009a]

Soit  $\mathcal{H}$  un espace d'hypothèses. Si  $D_S$  et  $D_T$  sont deux distributions sur  $X$ , alors :

$$\forall h \in \mathcal{H}, \quad \overbrace{R_{P_T}(h)}^{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \underbrace{\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)}_{\text{divergences}} + \nu$$

$\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)$  : la  $\mathcal{H}$ -divergence entre  $D_S$  et  $D_T$

$$\begin{aligned} \frac{1}{2}d_{\mathcal{H}}(D_S, D_T) &= \sup_{(h, h') \in \mathcal{H}^2} \left| R_{D_T}(h, h') - R_{D_S}(h, h') \right| \\ &= \sup_{(h, h') \in \mathcal{H}^2} \left| \mathbf{E}_{\mathbf{x}^t \sim D_T} \mathbf{I}[h(\mathbf{x}^t) \neq h'(\mathbf{x}^t)] - \mathbf{E}_{\mathbf{x}^s \sim D_S} \mathbf{I}[h(\mathbf{x}^s) \neq h'(\mathbf{x}^s)] \right| \end{aligned}$$

# La théorie classique de l'adaptation de domaine

Les résultats de S. Ben-David *et al.* et Mansour *et al.*

Théorème classique [Ben-David *et al.*, 2010, Mansour *et al.*, 2009a]

Soit  $\mathcal{H}$  un espace d'hypothèses. Si  $D_S$  et  $D_T$  sont deux distributions sur  $X$ , alors :

$$\forall h \in \mathcal{H}, \quad \overbrace{R_{P_T}(h)}^{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \underbrace{\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)}_{\text{divergences}} + \nu$$

$\nu$  : divergence entre les étiquetages

$$\nu = \inf_{h' \in \mathcal{H}} (R_{P_S}(h') + R_{P_T}(h')),$$

erreur jointe optimale [Ben-David *et al.*, 2010]

ou  $\nu = R_{P_T}(h_T^*) + R_{P_T}(h_T^*, h_S^*)$ ,  
 $h_{\mathcal{X}}^*$  est la meilleure hypothèse sur le domaine  $\mathcal{X}$  [Mansour *et al.*, 2009a]

# La théorie classique de l'adaptation de domaine

Les résultats de S. Ben-David *et al.* et Mansour *et al.*

Théorème classique [Ben-David *et al.*, 2010, Mansour *et al.*, 2009a]

Soit  $\mathcal{H}$  un espace d'hypothèses. Si  $D_S$  et  $D_T$  sont deux distributions sur  $X$ , alors :

$$\forall h \in \mathcal{H}, \quad \overbrace{R_{P_T}(h)}^{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \underbrace{\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)}_{\text{divergence}} + \psi$$

## Idée

Minimiser cette borne pour construire un nouvel espace de projection

- rapprochant les deux distributions marginales
- tout en gardant de bonnes performances sur le domaine source
- en supposant que les deux domaines sont liés



# La théorie classique de l'adaptation de domaine

Les résultats de S. Ben-David *et al.* et Mansour *et al.*

Théorème classique [Ben-David *et al.*, 2010, Mansour *et al.*, 2009a]

Soit  $\mathcal{H}$  un espace d'hypothèses. Si  $D_S$  et  $D_T$  sont deux distributions sur  $X$ , alors :

$$\forall h \in \mathcal{H}, \quad \overbrace{R_{P_T}(h)}^{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \underbrace{\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)}_{\text{divergence}} + \psi$$

## Première contribution

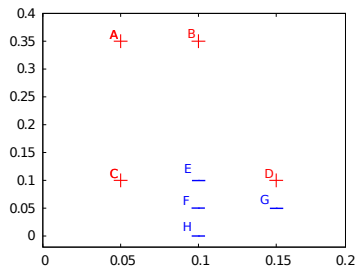
Travailler sur l'espace de projection explicite défini à partir de fonctions de similarités  $(\epsilon, \gamma, \tau)$ -bonnes

- 1 La théorie de l'adaptation de domaine
- 2 Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes**
- 3 Analyse PAC-Bayésienne de l'adaptation de domaine
- 4 Conclusion et perspectives générales

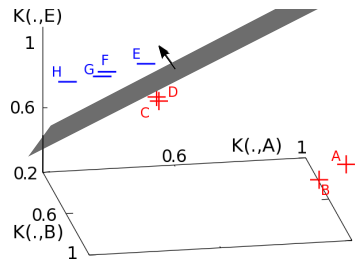
# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Apprentissage supervisé avec des fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes [Balcan et al., 2008] (1/2)

$$K : X \times X \rightarrow [-1, +1] \quad \Rightarrow \quad \phi^R : \begin{cases} X & \rightarrow \mathbb{R}^r \\ \mathbf{x} & \mapsto \langle K(\mathbf{x}, \mathbf{x}'_1), \dots, K(\mathbf{x}, \mathbf{x}'_r) \rangle \end{cases}$$
$$R = \{\mathbf{x}'_j\}_{j=1}^r$$



$\phi^R \rightarrow$



un classifieur linéaire dans le  $\phi^R$ -espace  $\Leftrightarrow$  un vote de majorité pondéré sur  $\{K(\cdot, \mathbf{x}'_j)\}_{\mathbf{x}'_j \in R}$

$$h(\mathbf{x}) = \text{sign} \left[ \sum_{j=1}^r \alpha_j K(\mathbf{x}, \mathbf{x}'_j) \right]$$

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Apprentissage supervisé avec des fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes [Balcan et al., 2008] (2/2)

Définition [Balcan et al., 2008]

$K : X \times X \rightarrow [-1, +1]$  est dite  $(\epsilon, \gamma, \tau)$ -bonne sur un domaine  $P$  si

- (i) une proportion de  $1 - \epsilon$  des exemples  $(\mathbf{x}, y)$  satisfont :  $\mathbf{E}_{(\mathbf{x}', y') \sim P} [yy' K(\mathbf{x}, \mathbf{x}') | \mathbf{x}' \in R] \geq \gamma$
- (ii)  $\Pr_{\mathbf{x}'} [\mathbf{x}' \in R] \geq \tau$

Intuitivement

Pour  $(\mathbf{x}_1, y_1) \sim P$ , on veut qu'en moyenne pour  $(\mathbf{x}'_2, y'_2) \in R$

$$\begin{array}{l} \text{si } y_1 = y'_2 \\ \mathbf{x}_1 \text{ est similaire à } \mathbf{x}'_2 \\ K(\mathbf{x}_1, \mathbf{x}'_2) \geq \gamma \end{array}$$

$$\begin{array}{l} \text{si } y_1 \neq y'_2 \\ \mathbf{x}_1 \text{ est similaire à } \mathbf{x}'_2 \\ K(\mathbf{x}_1, \mathbf{x}'_2) \leq -\gamma \end{array}$$

Remarque : Moins de contraintes qu'une fonction noyau

$K$  ni symétrique, ni SDP  $\implies$  Facilite la modification l'espace pour adapter

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Minimiser l'erreur sur le domaine source

## Rappel

Minimiser  $R_{P_T}(h)$  à l'aide de la borne

$$\underbrace{R_{P_T}(h)}_{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \underbrace{\frac{1}{2} d_{\mathcal{H}}(D_S, D_T)}_{\text{divergence}} + \psi$$

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Minimiser l'erreur sur le domaine source

## Rappel

Minimiser  $R_{P_T}(h)$  à l'aide de la borne

$$\underbrace{R_{P_T}(h)}_{\text{erreur cible}} \leq \underbrace{R_{P_S}(h)}_{\text{erreur source}} + \underbrace{\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)}_{\text{divergence}} + \psi$$

On minimise  $R_{P_S}(h)$  via le problème d'optimisation pour les SF

avec  $S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{m_s}$  et  $R = \{\mathbf{x}'_j\}_{j=1}^r$

$$\left\{ \begin{array}{l} \min_{\alpha} \frac{1}{m_s} \sum_{i=1}^{m_s} L(h, (\mathbf{x}_i^s, y_i^s)) + \lambda \|\alpha\|_1, \\ \text{avec } L(h, (\mathbf{x}_i^s, y_i^s)) = \left[ 1 - y_i^s \sum_{j=1}^r \alpha_j K(\mathbf{x}_i^s, \mathbf{x}'_j) \right]_+ \end{array} \right.$$

$$\|\alpha\|_1 = \sum_{j=1}^r |\alpha_j|$$

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Minimiser la divergence entre les domaines

Minimiser  $d_{\mathcal{H}}(D_S, D_T) \iff$  Rapprocher les marginales

**Problème** : minimiser simultanément la divergence et l'erreur source est difficile...

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Minimiser la divergence entre les domaines

Minimiser  $d_{\mathcal{H}}(D_S, D_T) \iff$  Rapprocher les marginales

**Problème** : minimiser simultanément la divergence et l'erreur source est difficile...

## Notre solution

$\mathcal{C}_{ST}$  un ensemble de couples  $(\mathbf{x}^s, \mathbf{x}^t) \in S \times T$

Construction d'un nouvel espace de projection  $\phi_{new}^R$

t.q. la différence entre les pertes de  $\mathbf{x}^s$  et  $\mathbf{x}^t$  est faible

$$\left| \left[ 1 - y \sum_{j=1}^r \alpha_j K(\mathbf{x}^s, \mathbf{x}'_j) \right]_+ - \left[ 1 - y \sum_{j=1}^r \alpha_j K(\mathbf{x}^t, \mathbf{x}'_j) \right]_+ \right| \leq \underbrace{\left\| ({}^t\phi^R(\mathbf{x}^s) - {}^t\phi^R(\mathbf{x}^t)) \text{diag}(\boldsymbol{\alpha}) \right\|_1}_{\left\| {}^t\phi_{new}^R(\mathbf{x}^s) - {}^t\phi_{new}^R(\mathbf{x}^t) \right\|_1}$$



# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Minimiser la divergence entre les domaines

Minimiser  $d_{\mathcal{H}}(D_S, D_T) \iff$  Rapprocher les marginales

**Problème** : minimiser simultanément la divergence et l'erreur source est difficile...

## Notre solution

$\mathcal{C}_{ST}$  un ensemble de couples  $(\mathbf{x}^s, \mathbf{x}^t) \in S \times T$

Construction d'un nouvel espace de projection  $\phi_{new}^R$

t.q. la différence entre les pertes de  $\mathbf{x}^s$  et  $\mathbf{x}^t$  est faible

$$\left| \left[ 1 - y \sum_{j=1}^r \alpha_j K(\mathbf{x}^s, \mathbf{x}'_j) \right]_+ - \left[ 1 - y \sum_{j=1}^r \alpha_j K(\mathbf{x}^t, \mathbf{x}'_j) \right]_+ \right| \leq \underbrace{\left\| ({}^t\phi^R(\mathbf{x}^s) - {}^t\phi^R(\mathbf{x}^t)) \text{diag}(\boldsymbol{\alpha}) \right\|_1}_{\left\| {}^t\phi_{new}^R(\mathbf{x}^s) - {}^t\phi_{new}^R(\mathbf{x}^t) \right\|_1}$$

$$\Rightarrow \phi_{new}^R(\cdot) = \left\langle \underbrace{\alpha_1 K(\cdot, \mathbf{x}'_1)}_{K_{new}(\cdot, \mathbf{x}'_1)}, \dots, \underbrace{\alpha_r K(\cdot, \mathbf{x}'_r)}_{K_{new}(\cdot, \mathbf{x}'_r)} \right\rangle$$

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Problème d'optimisation global

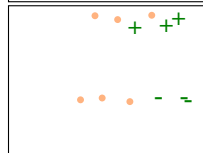
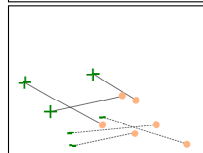
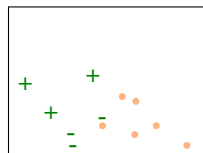
## Algorithme DASF

Avec  $S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{m_s} \sim (P_S)^{m_s}$  et  $R = \{\mathbf{x}'_j\}_{j=1}^r$

Construction du  $\phi^R$ -espace à l'aide  $\alpha$  appris par

$$\min_{\alpha} \frac{1}{m_s} \sum_{i=1}^{m_s} L(h, (\mathbf{x}_i^s, y_i^s)) + \lambda \|\alpha\|_1$$
$$+ \beta \sum_{(\mathbf{x}^s, \mathbf{x}^t) \in \mathcal{C}_{ST}} \left\| ({}^t\phi^R(\mathbf{x}^s) - {}^t\phi^R(\mathbf{x}^t)) \text{diag}(\alpha) \right\|_1$$

$$\text{avec } L(h, (\mathbf{x}_i^s, y_i^s)) = \left[ 1 - y_i \sum_{j=1}^r \alpha_j K(\mathbf{x}_i^s, \mathbf{x}'_j) \right]_+$$



# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Problème d'optimisation global

## Algorithme DASF

Avec  $S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{m_s} \sim (P_S)^{m_s}$  et  $R = \{\mathbf{x}'_j\}_{j=1}^r$

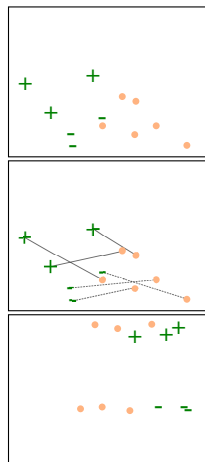
Construction du  $\phi^R$ -espace à l'aide  $\alpha$  appris par

$$\min_{\alpha} \frac{1}{m_s} \sum_{i=1}^{m_s} L(h, (\mathbf{x}_i^s, y_i^s)) + \lambda \|\alpha\|_1$$
$$+ \beta \sum_{(\mathbf{x}^s, \mathbf{x}^t) \in \mathcal{C}_{ST}} \left\| ({}^t\phi^R(\mathbf{x}^s) - {}^t\phi^R(\mathbf{x}^t)) \text{diag}(\alpha) \right\|_1$$

$$\text{avec } L(h, (\mathbf{x}_i^s, y_i^s)) = \left[ 1 - y_i \sum_{j=1}^r \alpha_j K(\mathbf{x}_i^s, \mathbf{x}'_j) \right]_+$$

Tester tous les couples  $(\mathbf{x}^s, \mathbf{x}^t) \in \mathcal{C}_{ST}$  est intraitable

⇒ Procédure itérative + Validation “circulaire”



## L'intuition portée par la théorie

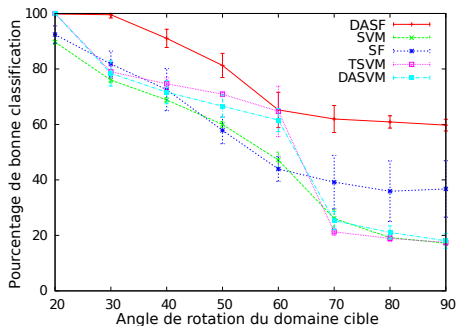
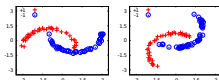
- Borne en généralisation
  - ↪ DASF est robuste sur le domaine **source**
- Analyse de la parcimonie
  - ↪ la parcimonie dépend de l'écart entre les coordonnées
  - ⇒ Les domaines sont éloignés
  - ⇒ La différence entre les coordonnées est élevée
  - ⇒ augmentation de la parcimonie

# AD par pondération de fonctions de similarité ( $\epsilon, \gamma, \tau$ )-bonnes

Experimentations - lunes jumelles

Problème jouet : "lunes jumelles"

- 1 domaine source
- 8 domaines cibles selon 8 angles de rotations
- 10 tirages aléatoire pour chaque angle
- Performances sur 1500 exemples cibles
- Noyau gaussien ou renormalisation (non SDP, non symétrique) du noyau



Taille 10 10 9 8 4 4 4 3

# AD par pondération de fonctions de similarité ( $\epsilon, \gamma, \tau$ )-bonnes

Experimentations - Annotation d'images

## Annotation d'images

- Domaine source : PascalVOC 2007 avec un ratio +/- de 1/3
- Deux domaines cibles :
  - ▶ Ratio +/- différent : PascalVOC 2007 Test
  - ▶ De même ratio +/- : TrecVid 2007
- F-mesure sur le domaine cible
- Noyau gaussien ou renormalisation (non SDP, non symétrique) du noyau

		SVM	SF	TSVM	DASVM	DASF
<b>VOC vs VOC</b>						
Moy. sur	F-mes.	0.22	0.19	0.17	0.20	<b>0.25</b>
20 conc.	Taille	642	210	705	622	<b>200</b>
<b>VOC vs Trec</b>						
BOAT	F-mes.	0.56	0.49	0.56	0.52	<b>0.57</b>
	Taille	351	214	498	202	<b>120</b>
CAR	F-mes.	0.43	0.50	0.52	<b>0.55</b>	<b>0.55</b>
	Taille	1096	<b>176</b>	631	627	254
MONITOR	F-mes.	0.19	0.34	0.37	0.30	<b>0.42</b>
	Taille	698	246	741	523	<b>151</b>
PERSON	F-mes.	0.52	0.45	0.46	0.54	<b>0.57</b>
	Taille	951	226	1024	274	<b>19</b>
PLANE	F-mes.	0.32	0.54	0.61	0.52	<b>0.66</b>
	Taille	428	178	259	450	<b>7</b>
Moy. sur	F-mes.	0.40	0.47	0.50	0.49	<b>0.55</b>
les 5 conc.	Taille	705	208	631	415	<b>110</b>

**VOC vs VOC** : Points raisonnables pour le concept Personne



- Idée : Contrôle de la borne d'adaptation de domaine théorique
- Algorithme d'adaptation itératif
  - ▶ construit un nouvel espace de projection
  - ▶ rapproche les domaines tout en gardant de bonnes garanties sur le domaine source
  - ▶ en repondérant une fonction de similarité  $(\epsilon, \gamma, \tau)$ -bonne
- Pas d'étiquette cible (mais a été généralisé à l'utilisation d'étiquettes cibles)
- Avec des garanties en généralisation
  - ▶ l'algorithme est robuste sur le domaine source
- La parcimonie des modèles dépend de la difficulté du problème

# AD par pondération de fonctions de similarité $(\epsilon, \gamma, \tau)$ -bonnes

Et alors ?

## Les inconvénients de cette approche

- le couplage  $\implies$  les points proches sont de même étiquette
- procédure itérative  $\implies$  compliquée et lourde à mettre en œuvre
- pas de justification théorique de la minimisation de la divergence
- la  $\mathcal{H}$ -divergence est difficile à optimiser en même temps que l'erreur source
- l'analyse classique peut être vue comme une analyse dans le pire cas



Les inconvénients de cette approche

- le couplage  $\implies$  les points proches sont de même étiquette
- procédure itérative  $\implies$  compliquée et lourde à mettre en œuvre
- pas de justification théorique de la minimisation de la divergence
- la  $\mathcal{H}$ -divergence est difficile à optimiser en même temps que l'erreur source
- l'analyse classique peut être vue comme une analyse dans le pire cas

## Seconde contribution

Une première analyse PAC-Bayésienne de  
l'adaptation de domaine

Analyse en “moyenne” sur  $\mathcal{H} \implies$  Garanties pour les votes de majorité

- 1 La théorie de l'adaptation de domaine
- 2 Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes
- 3 Analyse PAC-Bayésienne de l'adaptation de domaine**
- 4 Conclusion et perspectives générales

# Analyse PAC-Bayésienne de l'adaptation de domaine

La théorie PAC-Bayésienne - Description du problème

$X \in \mathbb{R}^d$ ;  $Y = \{-1, +1\}$  espace de sortie;  $\mathcal{H}$  ensemble de classifieurs;  $P_S$  domaine sur  $X \times Y$

## Classification supervisée

**Objectif** : Trouver  $h \in \mathcal{H}$  minimisant l'**erreur** :  $R_{P_S}(h) = \mathbf{E}_{(\mathbf{x}^s, y^s) \sim P_S} \mathbf{I}[h(\mathbf{x}^s) \neq y^s]$

## Approche PAC-Bayésienne [McAllester, 1999]

**Objectif** : Trouver le **vote de majorité  $\rho$ -pondéré**  $B_\rho$  sur  $\mathcal{H}$  minimisant l'erreur  $R_{P_S}(B_\rho)$

$$B_\rho(\mathbf{x}) = \text{sign} \left[ \sum_{h \in \mathcal{H}} \rho(h) h(\mathbf{x}) \right]$$

où  $\rho$  est la distribution **posterior** sur  $\mathcal{H}$  apprise à partir d'une distribution **prior**  $\pi$  sur  $\mathcal{H}$

# Analyse PAC-Bayésienne de l'adaptation de domaine

## La théorie PAC-Bayésienne - Le classifieur stochastique de Gibbs

- La théorie PAC-Bayésienne **ne** borne **pas** directement  $R_{P_S}(B_\rho)$
- mais l'erreur du **classifieur de Gibbs**  $G_\rho$   
qui prédit l'étiquette de  $\mathbf{x} \in X$  en
  - ▶ choisissant aléatoirement selon  $\rho$  un  $h$
  - ▶ puis en retournant  $h(\mathbf{x})$

# Analyse PAC-Bayésienne de l'adaptation de domaine

La théorie PAC-Bayésienne - Le classifieur stochastique de Gibbs

- La théorie PAC-Bayésienne **ne** borne **pas** directement  $R_{P_S}(B_\rho)$
- mais l'erreur du **classifieur de Gibbs**  $G_\rho$   
qui prédit l'étiquette de  $\mathbf{x} \in X$  en
  - ▶ choisissant aléatoirement selon  $\rho$  un  $h$
  - ▶ puis en retournant  $h(\mathbf{x})$

**IMPORTANT** — l'erreur de  $G_\rho$  correspond à l'espérance selon  $\rho$  des erreurs de  $\mathcal{H}$  :

$$R_{P_S}(G_\rho) = \mathbf{E}_{h \sim \rho} R_{P_S}(h)$$

On a :  $R_{P_S}(B_\rho) \leq 2R_{P_S}(G_\rho)$

$\Rightarrow$  une majoration de  $R_{P_S}(G_\rho) \Rightarrow$  une majoration de  $R_{P_S}(B_\rho)$

# Analyse PAC-Bayésienne de l'adaptation de domaine

La théorie PAC-Bayésienne - La borne de Catoni

Théorème ([Catoni, 2007], comme énoncé dans [Germain et al., 2009])

Pour tout  $P_S$  sur  $X \times Y$ , pour tout  $\mathcal{H}$ , pour tout  $\pi$  sur  $\mathcal{H}$ , pour tout  $\delta \in (0, 1]$ , et pour tout  $c > 0$ , avec une proba. d'au moins  $1 - \delta$  sur  $S \sim (P_S)^{m_S}$ , pour tout  $\rho$  sur  $\mathcal{H}$ , on a

$$\overbrace{R_{P_S}(G_\rho)}^{\text{erreur réelle}} \leq \frac{c}{1 - e^{-c}} \left[ \underbrace{R_S(G_\rho)}_{\text{erreur empirique}} + \underbrace{\frac{\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta}}{m_S c}}_{f(\text{"complexité", nb données})} \right]$$

avec  $\text{KL}(\rho \parallel \pi) = \mathbf{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)}$  la divergence de Kullback-Leibler

- $c$  contrôle le compromis entre  $R_S(G_\rho)$  et  $\frac{\text{KL}(\rho \parallel \pi)}{m_S}$

# Analyse PAC-Bayésienne de l'adaptation de domaine

La théorie PAC-Bayésienne - La borne de Catoni

Théorème ([Catoni, 2007], comme énoncé dans [Germain et al., 2009])

Pour tout  $P_S$  sur  $X \times Y$ , pour tout  $\mathcal{H}$ , pour tout  $\pi$  sur  $\mathcal{H}$ , pour tout  $\delta \in (0, 1]$ , et pour tout  $\mathbf{c} > 0$ , avec une proba. d'au moins  $1 - \delta$  sur  $S \sim (P_S)^{m_S}$ , pour tout  $\rho$  sur  $\mathcal{H}$ , on a

$$\underbrace{R_{P_S}(G_\rho)}_{\text{erreur réelle}} \leq \frac{\mathbf{c}}{1 - e^{-\mathbf{c}}} \left[ \underbrace{R_S(G_\rho)}_{\text{erreur empirique}} + \underbrace{\frac{\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta}}{m_S \mathbf{c}}}_{f(\text{"complexité", nb données})} \right]$$

avec  $\text{KL}(\rho \parallel \pi) = \mathbf{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)}$  la divergence de Kullback-Leibler

- $\mathbf{c}$  contrôle le compromis entre  $R_S(G_\rho)$  et  $\frac{\text{KL}(\rho \parallel \pi)}{m_S}$
- **L'algorithme PBGD** : ensemble des classifieurs linéaires  $\rho_{\mathbf{w}}$  gaussienne isotropique centrée en  $\mathbf{w}$   
Étant donné  $S \sim P_S^{m_S}$

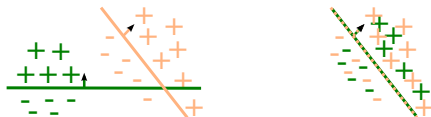
$$\operatorname{argmin}_{\rho_{\mathbf{w}}} \mathbf{C} m_S R_S(G_{\rho_{\mathbf{w}}}) + \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{KL}(\rho_{\mathbf{w}} \parallel \pi_0)}$$

# Analyse PAC-Bayésienne de l'adaptation de domaine

Revenons à l'adaptation de domaine

## Adaptation de domaine

Nécessité d'une divergence entre les domaines



## Théorie PAC-Bayésienne

Vote de majorité  $\rho$ -pondéré

Analyse en “moyenne” / en espérance sur  $\mathcal{H}$



# Analyse PAC-Bayésienne de l'adaptation de domaine

Une divergence entre domaines pour la théorie PAC-Bayésienne

Définition :  $\rho$ -désaccord entre  $D_S$  et  $D_T$

$$\text{dis}_\rho(D_S, D_T) = \left| \mathbf{E}_{(h, h') \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')] \right|$$

où :  $\mathbf{E}_{(h, h') \sim \rho^2} R_D(h, h') = \mathbf{E}_{(h, h') \sim \rho^2} \mathbf{E}_{\mathbf{x} \sim D} \mathbf{1}[h(\mathbf{x}) \neq h'(\mathbf{x})]$

# Analyse PAC-Bayésienne de l'adaptation de domaine

Une divergence entre domaines pour la théorie PAC-Bayésienne

Définition :  $\rho$ -désaccord entre  $D_S$  et  $D_T$

$$\text{dis}_\rho(D_S, D_T) = \left| \mathbf{E}_{(h, h') \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')] \right|$$

$$\text{où : } \mathbf{E}_{(h, h') \sim \rho^2} R_D(h, h') = \mathbf{E}_{(h, h') \sim \rho^2} \mathbf{E}_{\mathbf{x} \sim D} \mathbb{I}[h(\mathbf{x}) \neq h'(\mathbf{x})]$$

En suivant l'idée portée par la **C-borne** [Lacasse et al., 2007] :

$$R_{P_S}(B_\rho) \leq 1 - \frac{(1 - 2R_{P_S}(G_\rho))^2}{1 - 2 \mathbf{E}_{(h, h') \sim \rho^2} R_{D_S}(h, h')}$$

L'intuition est

Si  $R_{P_S}(G_\rho) \simeq R_{P_T}(G_\rho)$

$\Rightarrow R_{P_S}(B_\rho) \simeq R_{P_T}(B_\rho)$  quand  $\mathbf{E}_{(h, h') \sim \rho^2} R_{D_S}(h, h') \simeq \mathbf{E}_{(h, h') \sim \rho^2} R_{D_T}(h, h')$

# Analyse PAC-Bayésienne de l'adaptation de domaine

La borne d'adaptation de domaine pour le classifieur de Gibbs

## Théorème

Soit  $\mathcal{H}$  un espace d'hypothèses. Pour toute distribution  $\rho$  sur  $\mathcal{H}$ , on a

$$R_{P_T}(G_\rho) \leq R_{P_S}(G_\rho) + \text{dis}_\rho(D_S, D_T) + \nu_\rho$$

avec  $\rho_T^* = \operatorname{argmin}_\rho R_{P_T}(G_\rho)$  est le meilleur posterior sur le domaine cible et  $\nu_\rho = R_{P_T}(G_{\rho_T^*}) + \mathbf{E}_{h \sim \rho} \mathbf{E}_{h' \sim \rho_T^*} [R_{D_T}(h, h') + R_{D_S}(h, h')]$

# Analyse PAC-Bayésienne de l'adaptation de domaine

Différences avec la borne classique

$$R_{P_T}(G_\rho) \leq R_{P_S}(G_\rho) + \text{dis}_\rho(D_S, D_T) + \nu_\rho \quad \mathbf{Vs} \quad R_{P_T}(h) \leq R_{P_S}(h) + \frac{1}{2}d_{\mathcal{H}}(D_S, D_T) + \nu$$

- notre borne porte sur l'erreur du classifieur de Gibbs :  $R_{P_T}(G_\rho) = \mathbf{E}_{h \sim \rho} R_{P_T}(h)$
- la notion de divergence diffère

**Rappel :**

$$\begin{aligned} \text{dis}_\rho(D_S, D_T) &= \left| \mathbf{E}_{(h, h') \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')] \right| \\ \frac{1}{2}d_{\mathcal{H}}(D_S, D_T) &= \sup_{(h, h') \in \mathcal{H}^2} |R_{D_T}(h, h') - R_{D_S}(h, h')| \end{aligned}$$

- ▶  $\frac{1}{2}d_{\mathcal{H}}(D_S, D_T)$  est une divergence “dans le pire cas”
- ▶  $\text{dis}_\rho(D_S, D_T)$  est spécifique au classifieur  $G_\rho$  considéré
- ▶ on a :  $\frac{1}{2}d_{\mathcal{H}}(D_S, D_T) \geq \text{dis}_\rho(D_S, D_T)$

# Analyse PAC-Bayésienne de l'adaptation de domaine

La borne d'adaptation de domaine pour le classifieur de Gibbs

## Théorème

Soit  $\mathcal{H}$  un espace d'hypothèses. Pour toute distribution  $\rho$  sur  $\mathcal{H}$ , on a

$$R_{P_T}(G_\rho) \leq R_{P_S}(G_\rho) + \text{dis}_\rho(D_S, D_T) + \cancel{\nu_\rho}$$

avec  $\rho_T^* = \operatorname{argmin}_\rho R_{P_T}(G_\rho)$  est le meilleur posterior sur le domaine cible et  $\nu_\rho = R_{P_T}(G_{\rho_T^*}) + \mathbf{E}_{h \sim \rho} \mathbf{E}_{h' \sim \rho_T^*} [R_{D_T}(h, h') + R_{D_S}(h, h')]$

## Notre solution pour l'adaptation PAC-Bayésienne

Minimisation **simultanée** de  $R_{P_S}(G_\rho)$  et  $\text{dis}_\rho(D_S, D_T)$  avec des justifications théoriques

# Analyse PAC-Bayésienne de l'adaptation de domaine

La borne d'adaptation de domaine pour le classifieur de Gibbs — borne de consistance

## Théorème : Borne en généralisation PAC-Bayésienne (à la Catoni)

Pour tout  $P_S$  et  $P_T$  sur  $X \times Y$ , pour tout  $\mathcal{H}$ , pour tout  $\pi$  sur  $\mathcal{H}$ , pour tout  $\delta \in (0, 1]$ , pour tout  $\mathbf{a} > 0$ ,  $\mathbf{c} > 0$ , avec une proba. d'au moins  $1 - \delta$  sur  $S \times T \sim (P_S \times D_T)^m$ , on a

$$\forall \rho \sim \mathcal{H}, \quad \overbrace{R_{P_T}(G_\rho)}^{\text{erreur réelle cible}} \leq \underbrace{\mathbf{c}' R_S(G_\rho)}_{\text{erreur empirique source}} + \underbrace{\mathbf{a}' \text{dis}_\rho(S, T)}_{\rho\text{-desaccord empirique}} + \left(\frac{\mathbf{c}'}{\mathbf{c}} + \frac{2\mathbf{a}'}{\mathbf{a}}\right) \underbrace{\frac{\text{KL}(\rho \parallel \pi) + \ln \frac{3}{\delta}}{m}}_{f(\text{"complexité", nb données})} + \nu_\rho + \alpha' - 1$$

$$\text{où } \mathbf{c}' = \frac{\mathbf{c}}{1 - e^{-\mathbf{c}}} \text{ et } \mathbf{a}' = \frac{2\mathbf{a}}{1 - e^{-2\mathbf{a}}}$$

⇒ Comme pour PBGD : spécialisation aux classifieurs linéaires

### L'algorithme PBDA

Étant donné  $S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^m \sim (P_S)^m$ ,  $T = \{\mathbf{x}_i^t\}_{i=1}^m \sim (P_T)^m$

$$\operatorname{argmin}_{\rho_{\mathbf{w}}} C m R_S(G_{\rho_{\mathbf{w}}}) + A m \operatorname{dis}_{\rho_{\mathbf{w}}}(S, T) + \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\operatorname{KL}(\rho_{\mathbf{w}} \parallel \pi_0)}$$

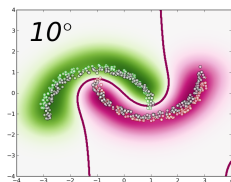
$$\text{où } \operatorname{dis}_{\rho_{\mathbf{w}}}(S, T) = \left| \mathbf{E}_{(h, h') \sim \rho_{\mathbf{w}}^2} R_S(h, h') - \mathbf{E}_{(h, h') \sim \rho_{\mathbf{w}}^2} R_T(h, h') \right|$$

# Analyse PAC-Bayésienne de l'adaptation de domaine

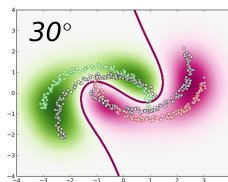
Expérimentations - lunes jumelles

Problème joué : "les lunes jumelles"

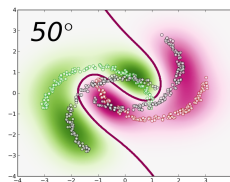
- 1 domaine source
- 8 domaines cibles selon 8 angles de rotations
- 10 tirages aléatoire pour chaque angle
- Performances sur un ensemble de test de 1 500 exemples cibles
- Noyau gaussien



(a)



(b)



(c)

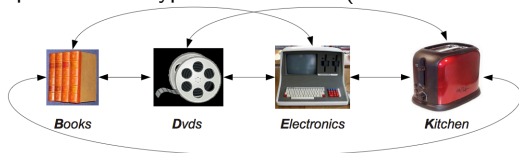


# Analyse PAC-Bayésienne de l'adaptation de domaine

Expérimentation - Analyse d'avis

Avis sur des produits *Amazon.com*

- 12 tâches d'adaptation d'un type vers un autre (ex. books → kitchen)



- ▶ Domaine source : 2 000 exemples étiquetés
- ▶ Domaine cible : 2 000 exemples non étiquetés
- ▶ Performances sur l'échantillon cible de test : entre 3 000 et 6 000 exemples
- Dimension des données :  $\sim 40,000$
- Noyau linéaire

	Erreur moyenne
PBGD3 [Germain et al., 2009]	0.226
SVM	0.231
DASF	0.275
DASVM [Bruzzone and Marconcini, 2010]	0.204
CODA [Chen et al., 2011]	0.210
PBDA (5× plus rapide)	0.208

### Conclusion

- La première analyse PAC-Bayésienne de l'adaptation de domaine
  - ▶ exprimée comme une espérance selon  $\rho$  sur une classe d'hypothèses
  - ▶ une divergence qui dépend de  $\rho$
  - ▶ directement optimisable (avec des justifications théoriques)
- Un premier algorithme spécialisé aux classifieurs linéaires
  - ▶ résultats prometteurs

### Points forts

- Minimisation directe divergence/erreur source
- Nouvelles pistes pour l'adaptation de domaine

# Analyse PAC-Bayésienne de l'adaptation de domaine

## Conclusion et Perspectives

### Conclusion

- La première analyse PAC-Bayésienne de l'adaptation de domaine
  - ▶ exprimée comme une espérance selon  $\rho$  sur une classe d'hypothèses
  - ▶ une divergence qui dépend de  $\rho$
  - ▶ directement optimisable (avec des justifications théoriques)
- Un premier algorithme spécialisé aux classifieurs linéaires
  - ▶ résultats prometteurs

### Points forts

- Minimisation directe divergence/erreur source
- Nouvelles pistes pour l'adaptation de domaine

### Perspectives

- Utilité de distribution(s) prior
- Traitement du terme  $\nu_\rho$  (*une première amélioration de la borne*)

- 1 La théorie de l'adaptation de domaine
- 2 Adaptation de domaine par pondération de fonctions de similarité  $(\epsilon, \gamma, \tau)$ -bonnes
- 3 Analyse PAC-Bayésienne de l'adaptation de domaine
- 4 Conclusion et perspectives générales

## Conclusion

- DASF → un algorithme “classique” d'adaptation
- PBDA → un premier algorithme pour l'adaptation PAC-Bayésienne

Mais aussi :

- P-MinCq → extension d'un algorithme PAC-Bayésien pour considérer un prior  
→ minimise la  $C$ -borne
- Étude PAC-Bayésienne du multiclasse → met en jeu la matrice de confusion  
→ généralisation de la  $C$ -borne

Comment dériver un algorithme pour **plusieurs classes, des sorties structurées** ?

- ↪ la matrice de confusion
- ↪ vote de majorité multiclasse
- ↪ code correcteur

## Perspectives générales

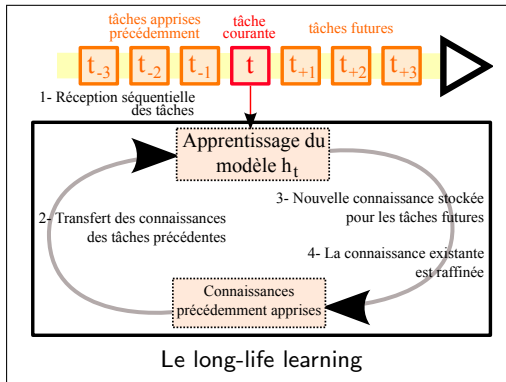
Comment dériver un algorithme pour **plusieurs classes, des sorties structurées** ?

- ↪ la matrice de confusion
- ↪ vote de majorité multiclasse
- ↪ code correcteur

Comment utiliser des modèles de **tâches** précédentes/simultanées ?

Comment tirer parti de **connaissances** ?

Comment **valider** les hyperparamètres ?



## Perspectives générales

Comment dériver un algorithme pour **plusieurs classes, des sorties structurées** ?

- ↪ la matrice de confusion
- ↪ vote de majorité multiclasse
- ↪ code correcteur

Comment utiliser des modèles de **tâches** précédentes/simultanées ?

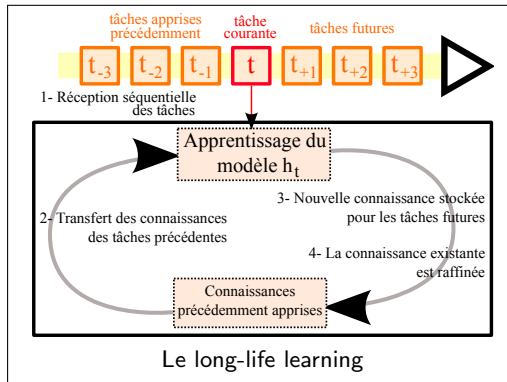
- ↪ vote de majorité
- ↪ apprentissage en ligne

Comment tirer parti de **connaissances** ?

- ↪ notion de prior et PAC-Bayes
- ↪ divergence entre tâches

Comment **valider** les hyperparamètres ?

- ↪ une validation "PAC-Bayésienne"





Merci pour votre attention.

## References



Balcan, M. F., Blum, A. and Srebro, N. (2008).  
Improved Guarantees for Learning via Similarity Functions.  
In *Proceedings of Annual Conference on Computational Learning Theory* pp. 287–298,.



Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F. and Vaughan, J. (2010).  
A theory of learning from different domains.  
*Machine Learning Journal* 79, 151–175.



Blitzer, J., Foster, D. and Kakade, S. (2011).  
Domain adaptation with coupled subspaces.  
*Journal of Machine Learning Research-Proceedings Track* 15, 173–181.



Blitzer, J., McDonald, R. and Pereira, F. (2006).  
Domain Adaptation with Structural Correspondence Learning.  
In *Proceedings of Conference on Empirical Methods on Natural Language Processing* pp. 120–128,.



Bruzzone, L. and Marconcini, M. (2010).  
Domain Adaptation Problems : A DASVM Classification Technique and a Circular Validation Strategy.  
*IEEE Transactions on Pattern Analysis and Machine Intelligence* 32, 770–787.



Catoni, O. (2007).  
PAC-Bayesian supervised classification : the thermodynamics of statistical learning, vol. 56.,  
Institute of Mathematical Statistic.



Chen, M., Weinberger, K. Q. and Blitzer, J. (2011).  
Co-Training for Domain Adaptation.  
In *Proceedings of Annual Conference on Neural Information Processing Systems* pp. 2456–2464,.



Daumé III, H. (2007).  
Frustratingly Easy Domain Adaptation.  
In *ACL*.



Daumé III, H., Kumar, A. and Saha, A. (2010).  
Co-regularization Based Semi-supervised Domain Adaptation.  
In *NIPS* pp. 478–486,.



Germain, P., Lacasse, A., Laviolette, F. and Marchand, M. (2009).  
PAC-Bayesian Learning of Linear Classifiers.  
In *Proceedings of International Conference on Machine Learning*.



Huang, J., Smola, A. J., Gretton, A., Borgwardt, K. M. and Scholkopf, B. (2007).  
Correcting sample selection bias by unlabeled data.  
*Advances in neural information processing systems* 19, 601.



Jiang, J. and Zhai, C. (2007).  
Instance Weighting for Domain Adaptation in NLP.  
In *Proceedings of Association for Computational Linguistics*.



Lacasse, A., Laviolette, F., Marchand, M., Germain, P. and Usunier, N. (2007).  
PAC-Bayes Bounds for the Risk of the Majority Vote and the Variance of the Gibbs Classifier.  
In *Proceedings of annual conference on Neural Information Processing Systems*.



Mansour, Y., Mohri, M. and Rostamizadeh, A. (2009a).  
Domain Adaptation : Learning Bounds and Algorithms.  
In *Proceedings of Annual Conference on Learning Theory* pp. 19–30,.



Mansour, Y., Mohri, M. and Rostamizadeh, A. (2009b).  
Multiple Source Adaptation and the Rényi Divergence.  
In *Proceedings of annual Conference on Uncertainty in Artificial Intelligence* pp. 367–374,.



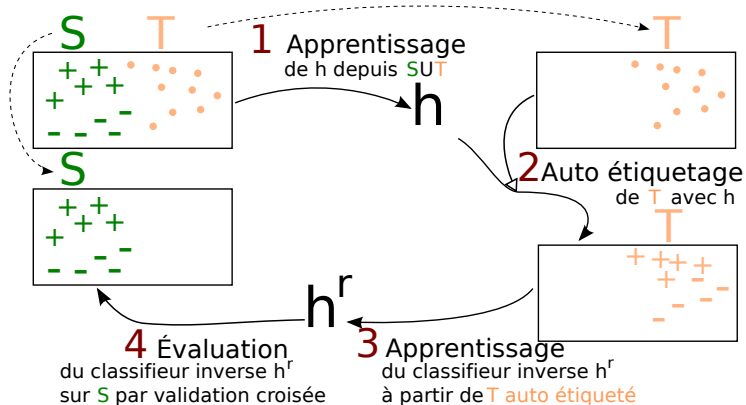
McAllester, D. A. (1999).  
Some PAC-Bayesian Theorems.  
*Machine Learning Journal* 37, 355–363.



Xu, H. and Mannor, S. (2010).  
Robustness and Generalization.  
In *Proceedings of Annual Conference on Computational Theory* pp. 503–515,.

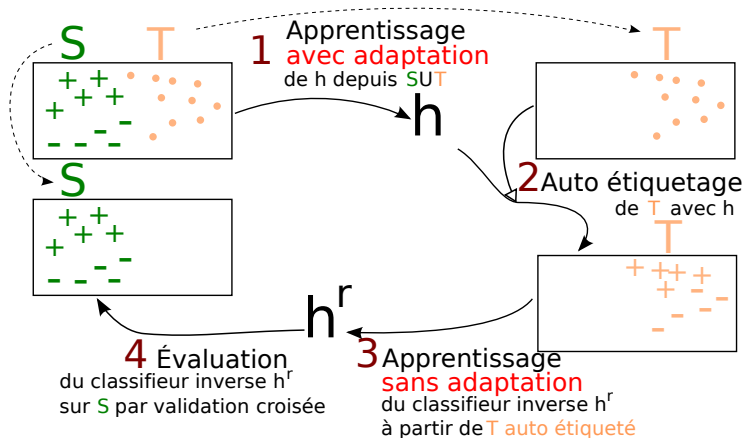
# Annexes

## Validation inverse/circulaire



# Annexes

## Validation inverse/circulaire pour DASF



- Algorithme robuste [Xu and Mannor, 2010]
  - ▶ Idée : “if a testing sample is similar to a training sample then the testing error is close to the training error” (sans adaptation)

⇒ DASF est robuste sur le domaine source

$$R_{P_T}^L(h) \leq R_S^L(h) + \frac{N_\eta}{\beta B_R + \lambda} + \sqrt{\frac{4M_\eta \ln 2 + 2 \ln \frac{1}{\delta}}{m_s}} + \frac{1}{2} d_{\mathcal{H}}(D_S, D_T) + \nu,$$

## Lemme

Soit  $B_R = \min_{\mathbf{x}'_j \in R} \left\{ \max_{(\mathbf{x}^s, \mathbf{x}^t) \in C_{ST}} |K(\mathbf{x}^s, \mathbf{x}'_j) - K(\mathbf{x}^t, \mathbf{x}'_j)| \right\} > 0$ .

Si  $\alpha^*$  est la solution optimale du problème, alors

$$\|\alpha^*\|_1 \leq \frac{1}{\beta B_R + \lambda}.$$

La parcimonie dépend des **hyperparameters** et de  $B_R$

- ⇒ Les domaines sont éloignés
- ⇒ La différence entre les coordonnées est élevée
- ⇒  $B_R$  tend à croître
- ⇒ augmentation de la parcimonie