

# PAC-Bayesian Theory and Domain Adaptation Algorithms

Pascal Germain

November 25, 2015

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# Definitions

## Learning example

An example  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is a **description-label** pair.

## Data generating distribution

Each example is an **observation from distribution**  $D$  on  $\mathcal{X} \times \mathcal{Y}$ .

## Learning sample

$$S \stackrel{\text{def}}{=} \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim D^n$$

## Classifier (or hypothesis)

$$h : \mathcal{X} \rightarrow \mathcal{Y}$$

## Binary classifier

$$h : \mathcal{X} \rightarrow \{-1, +1\}$$

## Learning algorithm

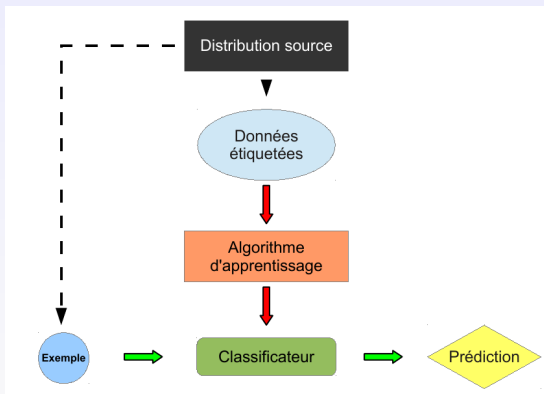
$$A(S) \rightarrow h$$

# I.I.D. Assumption

## Assumption

Examples are generated *i.i.d.* by a distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$ .

$$S = \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim D^n$$



## Risk (or generalization error)

Probability of misclassifying an example generated by distribution  $D$  :

$$\begin{aligned} R_D(h) &\stackrel{\text{def}}{=} \Pr_{(x,y) \sim D} (h(x) \neq y) \\ &= \mathbf{E}_{(x,y) \sim D} \mathbb{I}[y \cdot h(x) \leq 0] , \quad \langle \text{binary classification} \rangle \end{aligned}$$

where  $\mathbb{I}[a] = 1$  if predicate  $a$  is *true*;  $\mathbb{I}[a] = 0$  otherwise.

## Empirical risk

Error rate on the learning sample  $S \sim D^n$  :

$$\hat{R}_S(h) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i \cdot h(x_i) \leq 0] .$$

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# PAC-Bayesian Theory

Initiated by David McAllester (1999), the PAC-Bayesian theory gives generalization guarantees on **majority votes** of classifiers.

PAC guarantees (Probably Approximately Correct)

With probability at least  $\llcorner 1 - \delta \gg$ , the risk of classifier  $h$  is less than  $\llcorner \epsilon \gg$

$$\Pr_{S \sim D^n} \left( R_D(h) \leq \epsilon(R_S(h), n, \dots) \right) \geq 1 - \delta$$

Bayesian flavor

Incorporates *a priori* knowledge about the learning problem as a probability distribution over a family of classifiers.

Training bounds

- Gives generalization guarantees **not based on testing sample** ;
- Inspiration for conceiving **new learning algorithms**.



## 1 Basic Definitions

## 2 PAC-Bayesian Theory

### ■ Majority Vote Classifiers

- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# Majority Vote Classifiers

Given :

- A set of **voters**  $\mathcal{H} = \{h_1, h_2, h_3, \dots\}$  (discrete or continuous) ;
- A **weight** distribution  $Q$  on  $\mathcal{H}$ .

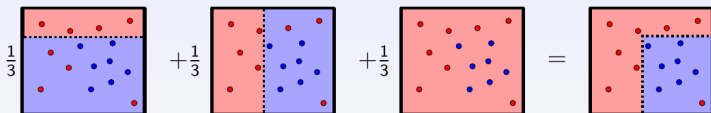
## Weighted majority vote

To predict the label of  $x \in \mathcal{X}$ , the classifier asks for the *prevailing opinion*

$$B_Q(x) \stackrel{\text{def}}{=} \text{sgn} \left( \mathbf{E}_{h \sim Q} h(x) \right)$$

Many learning algorithms output majority vote classifiers

AdaBoost, Random Forests, Bagging, ...



# A Surrogate Loss

Given

- A data distribution  $D$  on  $\mathcal{X} \times \{-1, +1\}$ ;
- A weight distribution  $Q$  on the set of voters  $\mathcal{H}$ .

Majority vote risk (or *Bayes Risk*)

$$R_D(B_Q) \stackrel{\text{def}}{=} \mathbf{E}_{(x,y) \sim D} \mathbf{I} \left[ \mathbf{E}_{h \sim Q} y \cdot h(x) \leq 0 \right]$$

Gibbs Risk

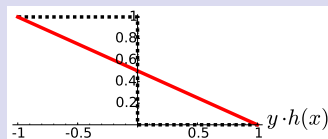
The stochastic Gibbs classifier  $G_Q(x)$  draws  $h' \in \mathcal{H}$  according to  $Q$  and output  $h'(x)$ .

$$\begin{aligned} R_D(G_Q) &\stackrel{\text{def}}{=} \mathbf{E}_{(x,y) \sim D} \mathbf{E}_{h \sim Q} \mathbf{I} \left[ y \cdot h(x) \leq 0 \right] \\ &= \mathbf{E}_{(x,y) \sim D} \left( \frac{1}{2} - \frac{1}{2} \mathbf{E}_{h \sim Q} y \cdot h(x) \right) \end{aligned}$$

Factor two

It is well-known that

$$R_D(B_Q) \leq 2 \times R_D(G_Q)$$

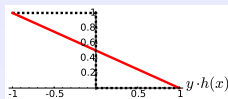


# From the *Factor 2* to the $\mathcal{C}$ -bound

From Markov's inequality (  $\Pr(X \geq a) \leq \frac{\mathbf{E}X}{a}$  ), we obtain :

## Factor 2 bound

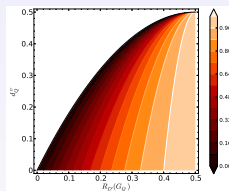
$$\begin{aligned} R_D(B_Q) &= \Pr_{(x,y) \sim D} (1 - y \cdot h(x) \geq 1) \\ &\leq \mathbf{E}_{(x,y) \sim D} (1 - y \cdot h(x)) = 2 R_D(G_Q). \end{aligned}$$



From Chebyshev's inequality (  $\Pr(X - \mathbf{E}X \geq a) \leq \frac{\text{Var } X}{a^2 + \text{Var } X}$  ), we obtain :

## The $\mathcal{C}$ -bound (Lacasse et al., 2006)

$$R_D(B_Q) \leq \mathcal{C}_Q^D \stackrel{\text{def}}{=} 1 - \frac{(1 - 2 \cdot R_D(G_Q))^2}{1 - 2 \cdot d_Q^D}$$



where  $d_Q^D$  is the **expected disagreement** :

$$d_Q^D \stackrel{\text{def}}{=} \mathbf{E}_{(x, \cdot) \sim D} \mathbf{E}_{h_1 \sim Q} \mathbf{E}_{h_2 \sim Q} \mathbb{I}[h_1(x) \neq h_2(x)] = \frac{1}{2} \left( 1 - \mathbf{E}_{(x, \cdot) \sim D'} \left[ \mathbf{E}_{h \sim Q} h(x) \right]^2 \right)$$

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# A Classical PAC-Bayesian Theorem

## Two principal components

- The **Gibbs empirical risk** :

$$\hat{R}_S(G_Q) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{2} - \frac{1}{2} \mathbf{E}_{h \sim Q} y_i \cdot h(x_i) \right)$$

- The **Kullback-Leibler divergence** between the *prior*  $P$  and the *posterior*  $Q$  :

$$\text{KL}(Q \| P) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim Q} \ln \frac{Q(h)}{P(h)}$$

## PAC-Bayesian theorem (McAllester, 2003)

For any distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$ , for any set of voters  $\mathcal{H}$ , for any distribution  $P$  on  $\mathcal{H}$ , for any  $\delta \in (0, 1]$ , we have, with probability at least  $1 - \delta$  over the choice of  $S \sim D^n$ ,

$$\forall Q \text{ on } \mathcal{H} : R_D(G_Q) \leq \hat{R}_S(G_Q) + \sqrt{\frac{1}{2n} \left[ \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right]}$$

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- **A General PAC-Bayesian Theorem**
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# A General PAC-Bayesian Theorem

$\Delta$ -function : «distance» between  $\hat{R}_S(G_Q)$  et  $R_D(G_Q)$

Convex function  $\Delta : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ .

## General theorem

*For any distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$ , for any set  $\mathcal{H}$  of voters, for any distribution  $P$  on  $\mathcal{H}$ , for any  $\delta \in (0, 1]$ , and for any  $\Delta$ -function, we have, with probability at least  $1 - \delta$  over the choice of  $S \sim D^n$ ,*

$$\forall Q \text{ on } \mathcal{H} : \Delta\left(\hat{R}_S(G_Q), R_D(G_Q)\right) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta} \right],$$

où

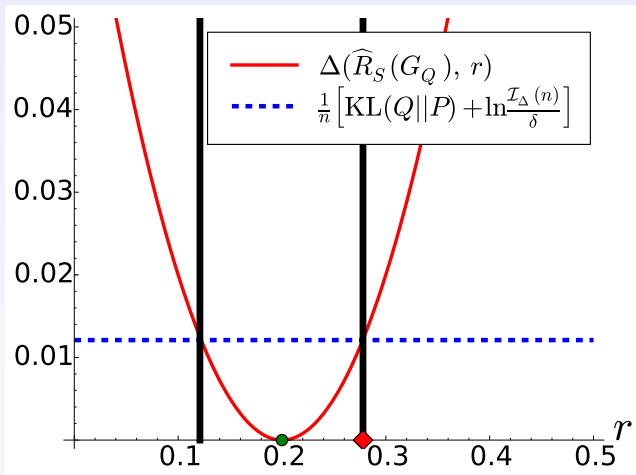
$$\mathcal{I}_\Delta(n) \stackrel{\text{def}}{=} \sup_{r \in [0, 1]} \left[ \sum_{k=0}^n \text{Bin}(k; n, r) e^{n\Delta(\frac{k}{n}, r)} \right]$$



## General theorem

$$\Pr_{S \sim D^n} \left( \forall Q \text{ on } \mathcal{H} : \Delta(\widehat{R}_S(G_Q), R_D(G_Q)) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta} \right] \right) \geq 1 - \delta.$$

**Interpretation.**



## General theorem

$$\Pr_{S \sim D^n} \left( \forall Q \text{ on } \mathcal{H} : \Delta \left( \widehat{R}_S(G_Q), R_D(G_Q) \right) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta} \right] \right) \geq 1 - \delta.$$

### Proof ideas.

#### Change of Measure Inequality

For any  $P$  and  $Q$  on  $\mathcal{H}$ , and for any measurable function  $\phi : \mathcal{H} \rightarrow \mathbb{R}$ , we have

$$\mathbf{E}_{h \sim Q} \phi(h) \leq \text{KL}(Q \| P) + \ln \left( \mathbf{E}_{h \sim P} e^{\phi(h)} \right).$$

#### Markov's inequality

$$\Pr(X \geq a) \leq \frac{\mathbf{E}X}{a} \iff \Pr(X \leq \frac{\mathbf{E}X}{\delta}) \geq 1 - \delta.$$

#### Probability of observing $k$ misclassifications among $n$ examples

Given a voter  $h$ , consider a **binomial variable** of  $n$  trials with **success**  $R_D(h)$  :

$$\begin{aligned} \text{Bin}(k; n, R_D(h)) &\stackrel{\text{def}}{=} \Pr_{S \sim D^n} \left( \widehat{R}_S(h) = \frac{k}{n} \right) \\ &= \binom{n}{k} (R_D(h))^k (1 - R_D(h))^{n-k}. \end{aligned}$$

# General theorem

$$\Pr_{S \sim D^n} \left( \forall Q \text{ on } \mathcal{H} : \Delta(\hat{R}_S(G_Q), R_D(G_Q)) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta} \right] \right) \geq 1 - \delta.$$

**Proof.**

$$n \cdot \Delta \left( \mathbf{E}_{h \sim Q} \hat{R}_S(h), \mathbf{E}_{h \sim Q} R_D(h) \right)$$

Jensen's Inequality

$$\leq \mathbf{E}_{h \sim Q} n \cdot \Delta(\hat{R}_S(h), R_D(h))$$

Change of measure

$$\leq \text{KL}(Q \| P) + \ln \mathbf{E}_{h \sim P} e^{n \Delta(\hat{R}_S(h), R_D(h))}$$

Markov's Inequality

$$\leq_{1-\delta} \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim D^n} \mathbf{E}_{h \sim P} e^{n \Delta(R_{S'}(h), R_D(h))}$$

Expectation swap

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim D^n} e^{n \Delta(R_{S'}(h), R_D(h))}$$

Binomial law

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k=0}^n \text{Bin}(k; n, R_D(h)) e^{n \Delta(\frac{k}{n}, R_D(h))}$$

Supremum over risk

$$\leq \text{KL}(Q \| P) + \ln \frac{1}{\delta} \sup_{r \in [0,1]} \left[ \sum_{k=0}^n \text{Bin}(k; n, r) e^{n \Delta(\frac{k}{n}, r)} \right]$$

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathcal{I}_\Delta(n).$$

□

## General theorem

$$\Pr_{S \sim D^n} \left( \forall Q \text{ on } \mathcal{H} : \Delta \left( \hat{R}_S(G_Q), R_D(G_Q) \right) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta} \right] \right) \geq 1 - \delta.$$

## Corollary

[...] with probability at least  $1 - \delta$  over the choice of  $S \sim D^n$ ,

$\forall Q \text{ on } \mathcal{H} :$

$$(a) \quad \text{kl} \left( \hat{R}_S(G_Q), R_D(G_Q) \right) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right], \quad (\text{Langford and Seeger, 2001})$$

$$(b) \quad R_D(G_Q) \leq \hat{R}_S(G_Q) + \sqrt{\frac{1}{2n} \left[ \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right]}, \quad (\text{McAllester, 1999})$$

$$(c) \quad R_D(G_Q) \leq \frac{1}{1 - e^{-c}} \left( c \cdot \hat{R}_S(G_Q) + \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{1}{\delta} \right] \right). \quad (\text{Catoni, 2007})$$

$$\text{kl}(q, p) \stackrel{\text{def}}{=} q \ln \frac{q}{p} + (1 - q) \ln \frac{1 - q}{1 - p} \geq 2(q - p)^2,$$

$$\Delta_c(q, p) \stackrel{\text{def}}{=} -\ln[1 - (1 - e^{-c}) \cdot p] - c \cdot q.$$

# Bounding the Expected Disagreement

## Expected disagreement

$$d_Q^D \stackrel{\text{def}}{=} \mathbf{E}_{(x, \cdot) \sim D} \mathbf{E}_{h_1 \sim Q} \mathbf{E}_{h_2 \sim Q} \mathbb{I}[h_1(x) \neq h_2(x)] = \frac{1}{2} \left( 1 - \mathbf{E}_{(x, \cdot) \sim D'} \left[ \mathbf{E}_{h \sim Q} h(x) \right]^2 \right)$$

## General theorem

[...] with probability at least  $1 - \delta$  over the choice of  $S \sim D^n$ ,

$$\forall Q \text{ on } \mathcal{H} : \quad \Delta \left( d_Q^S, d_Q^D \right) \leq \frac{1}{n} \left[ 2 \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta} \right],$$

## Corollary

- (a)  $\text{kl} \left( d_Q^S, d_Q^D \right) \leq \frac{1}{n} \left[ 2 \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right],$
- (b)  $d_Q^D \leq d_Q^S + \sqrt{\frac{1}{2n} \left[ 2 \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right]},$
- (c)  $d_Q^D \leq \frac{1}{1-e^{-c}} \left( c \cdot d_Q^S + \frac{1}{n} \left[ 2 \text{KL}(Q \| P) + \ln \frac{1}{\delta} \right] \right).$

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- **Transductive Learning**
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

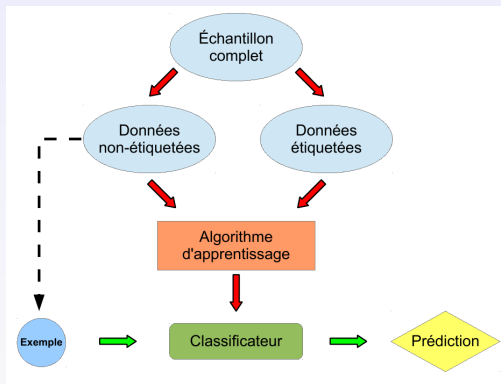
# Transductive Learning

## Assumption

Examples are drawn *without replacement* from a finite set  $Z$  of size  $N$ .

$$S = \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \subset Z$$

$$U = \{ (x_{n+1}, \cdot), (x_{n+2}, \cdot), \dots, (x_N, \cdot) \} = Z \setminus S$$



# General Theorem for Transductive Learning

## Observation

Inductive learning :  $n$  draws with replacement according to  $D \Rightarrow$  Binomial law.

Transductive learning :  $n$  draws without replacement in  $Z \Rightarrow$  Hypergeometric law.

## Theorem

*For any set  $Z$  of  $N$  examples, for any set  $\mathcal{H}$  of voters, for any distribution  $P$  on  $\mathcal{H}$ , for any  $\delta \in (0, 1]$ , and for any  $\Delta$ -function, we have, with probability at least  $1 - \delta$  over the choice of  $n$  examples among  $Z$ ,*

$$\forall Q \text{ on } \mathcal{H} : \quad \Delta(\hat{R}_S(G_Q), \hat{R}_Z(G_Q)) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{T}_\Delta(n, N)}{\delta} \right],$$

where

$$\mathcal{T}_\Delta(n, N) \stackrel{\text{def}}{=} \max_{K=0 \dots N} \left[ \sum_{k \in \mathcal{K}_{n, N, K}} \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}} e^{n \Delta(\frac{k}{n}, \frac{K}{N})} \right],$$

and  $\mathcal{K}_{n, N, K} \stackrel{\text{def}}{=} \{ \max[0, K + n - N], \dots, \min[n, K] \}.$



# Theorem

$$\Pr_{S \sim [Z]^n} \left( \forall Q \text{ on } \mathcal{H} : \Delta(\hat{R}_S(G_Q), \hat{R}_Z(G_Q)) \leq \frac{1}{n} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{T}_\Delta(n, N)}{\delta} \right] \right) \geq 1 - \delta.$$

**Poof.**

$$n \cdot \Delta \left( \mathbf{E}_{h \sim Q} \hat{R}_S(h), \mathbf{E}_{h \sim Q} \hat{R}_Z(h) \right)$$

Jensen's inequality

$$\leq \mathbf{E}_{h \sim Q} n \cdot \Delta(\hat{R}_S(h), \hat{R}_Z(h))$$

Change of measure

$$\leq \text{KL}(Q \| P) + \ln \mathbf{E}_{h \sim P} e^{n \Delta(\hat{R}_S(h), \hat{R}_Z(h))}$$

Markov's inequality

$$\leq_{1-\delta} \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim [Z]^n} \mathbf{E}_{h \sim P} e^{n \cdot \Delta(\hat{R}_{S'}(h), \hat{R}_Z(h))}$$

Expectations swap

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim [Z]^n} e^{n \cdot \Delta(\hat{R}_{S'}(h), \hat{R}_Z(h))}$$

Hypergeometric law

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k \in \mathcal{K}_{n, N, N \cdot R_Z(h)}} \frac{\binom{N \cdot \hat{R}_Z(h)}{k} \binom{N - N \cdot \hat{R}_Z(h)}{n - k}}{\binom{N}{n}} e^{n \cdot \Delta(\frac{k}{n}, \hat{R}_Z(h))}$$

Supremum over risk

$$\leq \text{KL}(Q \| P) + \ln \frac{1}{\delta} \max_{K=0 \dots N} \left[ \sum_{k \in \mathcal{K}_{n, N, K}} \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}} e^{n \Delta(\frac{k}{n}, \frac{K}{N})} \right]$$

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathcal{T}_\Delta(n, N).$$

□

# A New Transductive Bound for the Gibbs Risk

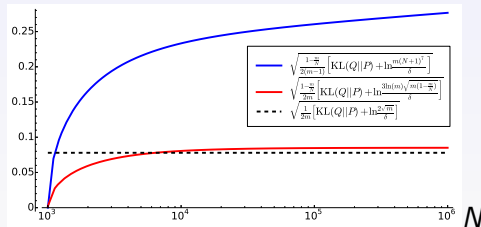
## Corollary

[...] with probability at least  $1-\delta$  over the choice of  $n$  examples among  $Z$ ,

$$\forall Q \text{ on } \mathcal{H} : \hat{R}_Z(G_Q) \leq \hat{R}_S(G_Q) + \sqrt{\frac{1-\frac{n}{N}}{2n} \left[ \text{KL}(Q\|P) + \ln \frac{3 \ln(n) \sqrt{n(1-\frac{n}{N})}}{\delta} \right]}.$$

## Theorem (Derbeko et al., 2004)

$$\hat{R}_Z(G_Q) \leq \hat{R}_S(G_Q) + \sqrt{\frac{1-\frac{n}{N}}{2(n-1)} \left[ \text{KL}(Q\|P) + \ln \frac{n(N+1)^7}{\delta} \right]}.$$



# A New Transductive Bound for the Bayes Risk

## Majority Vote Bound

[...] with probability at least  $1-\delta$  over the choice of  $n$  examples among  $Z$ ,

$\forall Q$  on  $\mathcal{H}$  :

$$(a) \quad \hat{R}_Z(B_Q) \leq 2 \times \bar{r} \quad (\text{Factor 2})$$

$$(b) \quad \hat{R}_Z(B_Q) \leq 1 - \frac{(1 - 2 \times \bar{r})^2}{1 - 2 \times d_Q^Z} \quad (\mathcal{C}\text{-bound})$$

where

$$\bar{r} := \hat{R}_S(G_Q) + \sqrt{\frac{1-\frac{n}{N}}{2n} \left[ \text{KL}(Q\|P) + \ln \frac{3 \ln(n) \sqrt{n(1-\frac{n}{N})}}{\delta} \right]},$$

$$d_Q^Z = \frac{1}{2} \left( 1 - \sum_{i=1}^N \left[ \mathbf{E}_{h \sim Q} h(x_i) \right]^2 \right).$$

# Empirical Comparison

Majority votes of *decision stumps* obtained with *AdaBoost*.

Dataset	N	n/N	$R_S(B_Q)$	Factor 2	C-bound
car	1728	0.1	0.105	1.092	-
		0.5	0.115	0.830	<b>0.819</b>
letter_AB	1555	0.1	0.000	<b>0.914</b>	0.961
		0.5	0.000	0.797	<b>0.626</b>
mushroom	8124	0.1	0.000	<b>0.964</b>	0.966
		0.5	0.000	0.875	<b>0.546</b>
nursery	12959	0.1	0.009	0.798	<b>0.692</b>
		0.5	0.010	0.711	<b>0.379</b>
optdigits	3823	0.1	0.000	1.055	-
		0.5	0.026	0.917	<b>0.793</b>
pageblock	5473	0.1	0.048	<b>0.979</b>	0.992
		0.5	0.057	0.894	<b>0.697</b>
pendigits	7494	0.1	0.023	<b>0.989</b>	0.997
		0.5	0.041	0.912	<b>0.706</b>
segment	2310	0.1	0.000	1.101	-
		0.5	0.014	0.920	<b>0.834</b>
spambase	4601	0.1	0.115	1.096	-
		0.5	0.137	0.973	<b>0.961</b>

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# A New Change of Measure

## Kullback-Leibler Change of Measure Inequality

For any  $P$  and  $Q$  on  $\mathcal{H}$ , and for any  $\phi : \mathcal{H} \rightarrow \mathbb{R}$ , we have

$$\mathbf{E}_{h \sim Q} \phi(h) \leq \text{KL}(Q \| P) + \ln \left( \mathbf{E}_{h \sim P} e^{\phi(h)} \right).$$

## Rényi Change of Measure Inequality

For any  $P$  and  $Q$  on  $\mathcal{H}$ , any  $\phi : \mathcal{H} \rightarrow \mathbb{R}$ , and for any  $\alpha > 1$ , we have

$$\frac{\alpha}{\alpha-1} \ln \mathbf{E}_{h \sim Q} \phi(h) \leq D_{\alpha}(Q \| P) + \ln \left( \mathbf{E}_{h \sim P} \phi(h)^{\frac{\alpha}{\alpha-1}} \right),$$

with  $D_{\alpha}(Q \| P) \stackrel{\text{def}}{=} \frac{1}{\alpha-1} \ln \left[ \mathbf{E}_{h \sim P} \left( \frac{Q(h)}{P(h)} \right)^{\alpha} \right].$

# Rényi-Based General Theorem

## General theorem

[...] for any  $\alpha > 1$ , with probability at least  $1 - \delta$  over the choice of  $S \sim D^n$ ,

$$\forall Q \text{ on } \mathcal{H}: \quad \ln \Delta \left( \widehat{R}_S(G_Q), R_D(G_Q) \right) \leq \frac{1}{\alpha'} \left[ D_\alpha(Q \| P) + \ln \frac{\mathcal{I}_\Delta^R(n, \alpha')}{\delta} \right],$$

with

$$\mathcal{I}_\Delta^R(n, \alpha') \stackrel{\text{def}}{=} \sup_{r \in [0, 1]} \left[ \sum_{k=0}^n \text{Bin}(k; n, r) \Delta\left(\frac{k}{n}, r\right)^{\alpha'} \right],$$

and  $\alpha' := \frac{\alpha}{\alpha - 1} > 1$ .

# General theorem (Rényi-Based)

$$\Pr_{S \sim D^n} \left( \forall Q \text{ on } \mathcal{H} : \ln \Delta \left( \widehat{R}_S(G_Q), R_D(G_Q) \right) \leq \frac{1}{\alpha'} \left[ D_\alpha(Q \| P) + \ln \frac{\mathcal{I}_\Delta^R(n, \alpha')}{\delta} \right] \right) \geq 1 - \delta.$$

**Proof.**

$$\alpha' := \frac{\alpha}{\alpha - 1}$$

$$\alpha' \cdot \ln \Delta \left( \mathbf{E}_{h \sim Q} \widehat{R}_S(h), \mathbf{E}_{h \sim Q} R_D(h) \right)$$

Jensen's Inequality

$$\leq \alpha' \cdot \ln \mathbf{E}_{h \sim Q} \Delta \left( \widehat{R}_S(h), R_D(h) \right)$$

Change of measure

$$\leq D_\alpha(Q \| P) + \ln \mathbf{E}_{h \sim P} \Delta \left( \widehat{R}_S(h), R_D(h) \right)^{\alpha'}$$

Markov's Inequality

$$\leq_{1-\delta} D_\alpha(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim D^n} \mathbf{E}_{h \sim P} \Delta \left( R_{S'}(h), R_D(h) \right)^{\alpha'}$$

Expectation swap

$$= D_\alpha(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim D^n} \Delta \left( R_{S'}(h), R_D(h) \right)^{\alpha'}$$

Binomial law

$$= D_\alpha(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k=0}^n \text{Bin}(k; n, R_D(h)) \Delta \left( \frac{k}{n}, R_D(h) \right)^{\alpha'}$$

Supremum over risk

$$\leq D_\alpha(Q \| P) + \ln \frac{1}{\delta} \sup_{r \in [0,1]} \left[ \sum_{k=0}^n \text{Bin}(k; n, r) \Delta \left( \frac{k}{n}, r \right)^{\alpha'} \right]$$

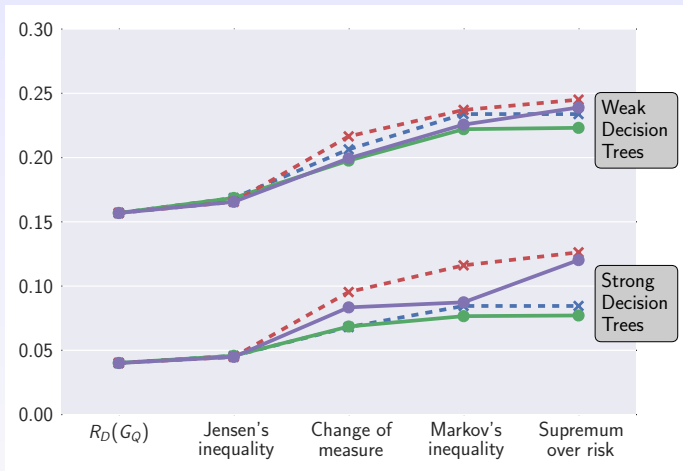
$$= D_\alpha(Q \| P) + \ln \frac{1}{\delta} \mathcal{I}_\Delta^R(n, \alpha').$$

□



# Empirical Study

Majority votes of 500 decision trees on *Mushroom* dataset



✕  $\text{KL}(Q\|P)$  and  $\Delta := 2(q-p)^2$

●  $D_\alpha(Q\|P)$  and  $\Delta := 2(q-p)^2$

✕  $\text{KL}(Q\|P)$  and  $\Delta := \text{kl}(q, p)$

●  $D_\alpha(Q\|P)$  and  $\Delta := \text{kl}(q, p)$

# Plan

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

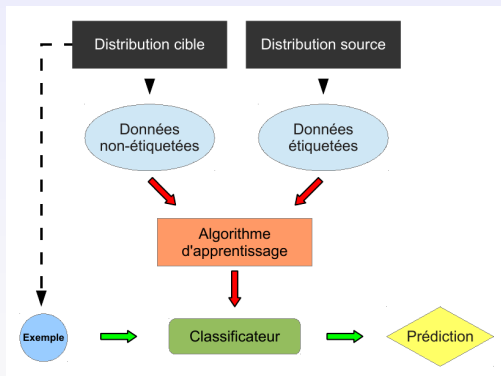
## 4 Conclusion and future works

# Domain Adaptation

## Assumption

Source and target examples are generated by different distributions

$$\begin{aligned} S &= \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim (D_S)^n \\ T &= \{ (x_1, \cdot), (x_2, \cdot), \dots, (x_n, \cdot) \} \sim (D_T)^n \end{aligned}$$



# Our Domain Adaptation Setting

## Binary classification tasks

- Input space :  $\mathbb{R}^d$
- Labels :  $\{-1, 1\}$

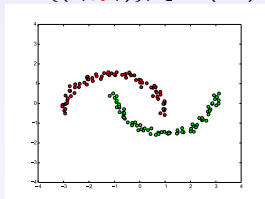
## Two different data distributions

- Source domain :  $D_S$
- Target domain :  $D_T$

A **domain adaptation** learning algorithm is provided with

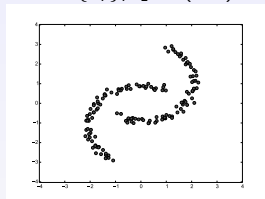
a **labeled source sample**

$$S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^n \sim (D_S)^n,$$



an **unlabeled target sample**

$$T = \{\mathbf{x}_i^t\}_{i=1}^n \sim (D_T)^n.$$



The goal is to build a classifier  $h : \mathbb{R}^d \rightarrow \{-1, 1\}$  with a low **target risk**

$$R_{D_T}(h) \stackrel{\text{def}}{=} \Pr_{(\mathbf{x}^t, y^t) \sim D_T} [h(\mathbf{x}^t) \neq y^t].$$

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

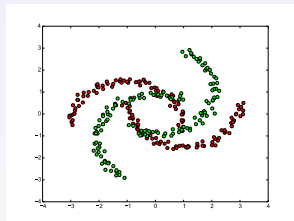
# Divergence between source and target domains

## Definition (Ben David et al., 2006)

Given two domain distributions  $D_S$  and  $D_T$ , and a **hypothesis class**  $\mathcal{H}$ , the  **$\mathcal{H}$ -divergence** between  $D_S$  and  $D_T$  is

$$d_{\mathcal{H}}(D_S, D_T) \stackrel{\text{def}}{=} 2 \sup_{h \in \mathcal{H}} \left| \Pr_{\mathbf{x}^s \sim D_S} [h(\mathbf{x}^s) = 1] + \Pr_{\mathbf{x}^t \sim D_T} [h(\mathbf{x}^t) = -1] - 1 \right|.$$

The  **$\mathcal{H}$ -divergence** measures the ability of an hypothesis class  $\mathcal{H}$  to **discriminate** between source  $D_S$  and target  $D_T$  distributions.



# Bound on the target risk

## Theorem (Ben David et al., 2006)

Let  $\mathcal{H}$  be a hypothesis class of VC-dimension  $d$ . With probability  $1 - \delta$  over the choice of samples  $S \sim (\textcolor{red}{D}_S)^n$  and  $T \sim (\textcolor{green}{D}_T)^n$ , for every  $h \in \mathcal{H}$  :

$$\textcolor{red}{R}_{D_T}(h) \leq \hat{R}_S(h) + \frac{4}{n} \sqrt{d \log \frac{2en}{d} + \log \frac{4}{\delta}} + \hat{d}_{\mathcal{H}}(\textcolor{red}{S}, \textcolor{green}{T}) + \frac{4}{n^2} \sqrt{d \log \frac{2n}{d} + \log \frac{4}{\delta}} + \beta$$

with  $\beta \geq \inf_{h^* \in \mathcal{H}} [\textcolor{red}{R}_{D_S}(h^*) + R_{D_T}(h^*)]$ .

Empirical risk on the **source sample** :

$$\hat{R}_S(h) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h(\textcolor{red}{\mathbf{x}}_i^S) \neq \textcolor{red}{y}_i^S].$$

Empirical  $\mathcal{H}$ -divergence :

$$\hat{d}_{\mathcal{H}}(\textcolor{red}{S}, \textcolor{green}{T}) \stackrel{\text{def}}{=} 2 \max_{h \in \mathcal{H}} \left[ \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h(\textcolor{red}{\mathbf{x}}_i^S) = 1] + \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h(\textcolor{green}{\mathbf{x}}_i^T) = -1] - 1 \right].$$

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works



# Nouvelle borne pour l'adaptation de domaine

$\mathcal{H}\Delta\mathcal{H}$ -distance (Ben-David et al., 2006, 2010)

$$d_{\mathcal{H}\Delta\mathcal{H}}(D_S, D_T) \stackrel{\text{def}}{=} 2 \sup_{h, h' \in \mathcal{H}} \left| \mathbf{E}_{(x^S, \cdot) \sim D_S} \mathbb{I}[h(x^S) \neq h'(x^S)] - \mathbf{E}_{(x^T, \cdot) \sim D_T} \mathbb{I}[h(x^T) \neq h'(x^T)] \right|$$

Distributions disagreement

$$\text{dis}_Q(D_S, D_T) \stackrel{\text{def}}{=} \left| d_Q^{D_T} - d_Q^{D_S} \right|$$

Theorem

[...] with probability  $1-\delta$  over the choice of  $S \times T \sim (D_S \times D_T)^n$ , we have

$\forall Q$  on  $\mathcal{H}$  :

$$R_{D_T}(G_Q) \leq c' \cdot \widehat{R}_S(G_Q) + a' \cdot \widehat{\text{dis}}_Q(S, T) + \left( \frac{c'}{c} + \frac{2a'}{a} \right) \frac{\text{KL}(Q \| P) + \ln \frac{3}{\delta}}{n} + \lambda_Q^* + a' - 1$$

where  $a' \stackrel{\text{def}}{=} \frac{2a}{1-e^{-2a}}$  et  $c' \stackrel{\text{def}}{=} \frac{c}{1-e^{-c}}$ .

# A New Domain Adaptation Algorithm

For linear classifiers

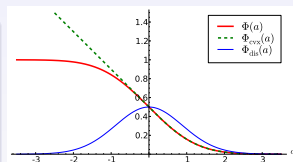
## PAC-Bayes specialization to linear classifier (Langford and Shawe-Taylor, 2002)

- Linear classifier :  $h_{\mathbf{w}}(\mathbf{x}) = \text{sgn}[\mathbf{w} \cdot \mathbf{x}]$
- Voters :  $\mathcal{H} = \{h_{\mathbf{v}} \mid \mathbf{v} \in \mathbb{R}^d\}$
- Prior  $P_0$  : isotropic Gaussian centered on  $\mathbf{0}$
- Posterior  $Q_{\mathbf{w}}$  : isotropic Gaussian centered on  $\mathbf{w}$
- $h_{\mathbf{w}}(\mathbf{x}) = B_{Q_{\mathbf{w}}}(\mathbf{x})$
- $R_D(Q_{\mathbf{w}}) = \mathbf{E}_{(\mathbf{x}, y) \sim D} \Phi\left(y \frac{\mathbf{w} \cdot \mathbf{x}}{\|\mathbf{x}\|}\right)$
- $\text{KL}(Q_{\mathbf{w}} \| P_0) = \frac{1}{2} \|\mathbf{w}\|^2$

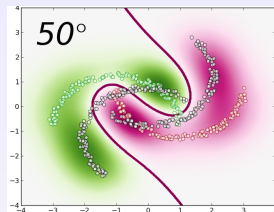
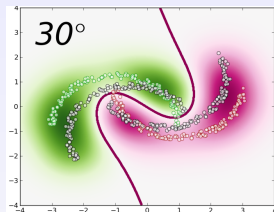
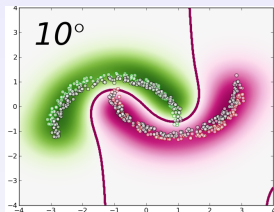
## PBDA

Minimize :

$$C \sum_{i=1}^n \Phi_c\left(y_i^S \frac{\mathbf{w} \cdot \mathbf{x}_i^S}{\|\mathbf{x}_i^S\|}\right) + A \left| \sum_{i=1}^n \Phi_d\left(\frac{\mathbf{w} \cdot \mathbf{x}_i^S}{\|\mathbf{x}_i^S\|}\right) - \Phi_d\left(\frac{\mathbf{w} \cdot \mathbf{x}_i^T}{\|\mathbf{x}_i^T\|}\right) \right| + \frac{\|\mathbf{w}\|^2}{2}$$

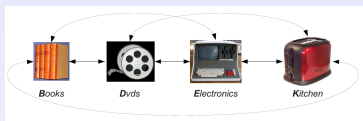


# Intertwining Moons Toy Dataset



# Amazon Reviews Dataset

**Input** : product review (bag of words) — **Output** : positive or negative rating.



source → target	PBGD	SVM	DASVM	CODA	PBDA
books→dvds	<b>0.174</b>	0.179	0.193	0.181	0.183
books→electronics	0.275	0.290	<b>0.226</b>	0.232	0.263
books→kitchen	0.236	0.251	<b>0.179</b>	0.215	0.229
dvds→books	<b>0.192</b>	0.203	0.202	0.217	0.197
dvds→electronics	0.256	0.269	<b>0.186</b>	0.214	0.241
dvds→kitchen	0.211	0.232	0.183	<b>0.181</b>	0.186
electronics→books	0.268	0.287	0.305	0.275	<b>0.232</b>
electronics→dvds	0.245	0.267	<b>0.214</b>	0.239	0.221
electronics→kitchen	<b>0.127</b>	0.129	0.149	0.134	0.141
kitchen→books	0.255	0.267	0.259	<b>0.247</b>	<b>0.247</b>
kitchen→dvds	0.244	0.253	<b>0.198</b>	0.238	0.233
kitchen→electronics	0.235	0.149	0.157	0.153	<b>0.129</b>
<b>Mean</b>	0.226	0.231	<b>0.204</b>	0.210	0.208

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

# A New Perspective on Domain Adaptation

## Theorem

For all pairs  $D_S$  and  $D_T$  on  $\mathcal{X} \times \mathcal{Y}$ , for all set  $\mathcal{H}$  of voters, and for all  $q > 0$ ,

$$\forall Q \text{ sur } \mathcal{H}, \quad R_{D_T}(G_Q) \leq \frac{1}{2} d_Q^{D_T} + \beta_q(D_T \| D_S) \times \left[ e_Q^{D_S} \right]^{1 - \frac{1}{q}}.$$

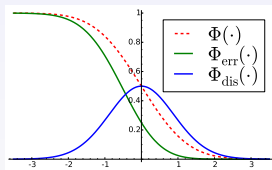
where  $\beta_q(D_T \| D_S) = \left[ \mathbf{E}_{(x,y) \sim D_S} \left( \frac{D_T(x,y)}{D_S(x,y)} \right)^q \right]^{\frac{1}{q}}$ ,

and  $e_Q^{D_S} \stackrel{\text{def}}{=} \mathbf{E}_{(x,y) \sim D} \mathbf{E}_{h_1 \sim Q} \mathbf{E}_{h_2 \sim Q} \mathbb{I}[h_1(x) \neq y] \times \mathbb{I}[h_2(x) \neq y]$ .

## DALC

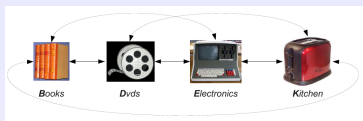
Minimize :

$$A \sum_{i=1}^{n_t} \Phi_d \left( \frac{\mathbf{w} \cdot \mathbf{x}_i^T}{\|\mathbf{x}_i^T\|} \right) + B \sum_{i=1}^{n_s} \Phi_{\text{err}} \left( y_i^S \frac{\mathbf{w} \cdot \mathbf{x}_i^S}{\|\mathbf{x}_i^S\|} \right) + \frac{\|\mathbf{w}\|^2}{2}$$



# Amazon Reviews Dataset

**Input** : product review (bag of words) — **Output** : positive or negative rating.



source → target	PBGD	SVM	DASVM	CODA	PBDA	DALC
books→dvds	<b>0.174</b>	0.179	0.193	0.181	0.183	0.178
books→electronics	0.275	0.290	0.226	0.232	0.263	<b>0.212</b>
books→kitchen	0.236	0.251	<b>0.179</b>	0.215	0.229	0.194
dvds→books	0.192	0.203	0.202	0.217	0.197	<b>0.186</b>
dvds→electronics	0.256	0.269	<b>0.186</b>	0.214	0.241	0.245
dvds→kitchen	0.211	0.232	0.183	0.181	0.186	<b>0.175</b>
electronics→books	0.268	0.287	0.305	0.275	<b>0.232</b>	0.240
electronics→dvds	0.245	0.267	<b>0.214</b>	0.239	0.221	0.256
electronics→kitchen	0.127	0.129	0.149	0.134	0.141	<b>0.123</b>
kitchen→books	0.255	0.267	0.259	0.247	0.247	<b>0.236</b>
kitchen→dvds	0.244	0.253	<b>0.198</b>	0.238	0.233	0.225
kitchen→electronics	0.235	0.149	0.157	0.153	<b>0.129</b>	0.131
<b>Mean</b>	0.226	0.231	0.204	0.210	0.208	<b>0.200</b>

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works



# Bound on the target risk

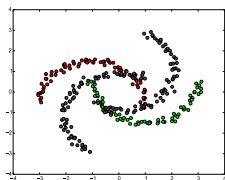
## Theorem (Ben David et al., 2006)

Let  $\mathcal{H}$  be a hypothesis class of VC-dimension  $d$ . With probability  $1 - \delta$  over the choice of samples  $S \sim (\textcolor{red}{D}_S)^n$  and  $T \sim (\textcolor{green}{D}_T)^n$ , for every  $h \in \mathcal{H}$  :

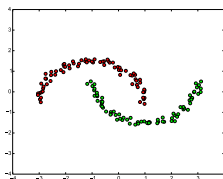
$$\textcolor{red}{R}_{D_T}(h) \leq \hat{R}_S(h) + \frac{4}{n} \sqrt{d \log \frac{2en}{d} + \log \frac{4}{\delta}} + \hat{d}_{\mathcal{H}}(\textcolor{red}{S}, \textcolor{green}{T}) + \frac{4}{n^2} \sqrt{d \log \frac{2n}{d} + \log \frac{4}{\delta}} + \beta$$

with  $\beta \geq \inf_{h^* \in \mathcal{H}} [\textcolor{red}{R}_{D_S}(h^*) + \textcolor{red}{R}_{D_T}(h^*)]$ .

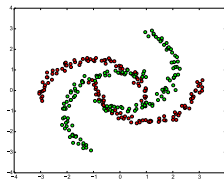
**Target risk  $\textcolor{red}{R}_{D_T}(h)$  is low**  
if, given  $\textcolor{red}{S}$  and  $\textcolor{green}{T}$ ,



**$\hat{R}_S(h)$  is small,**  
i.e.,  $h \in \mathcal{H}$  is good on



**and  $\hat{d}_{\mathcal{H}}(\textcolor{red}{S}, \textcolor{green}{T})$  is small,**  
i.e., all  $h' \in \mathcal{H}$  are bad on



# Domain-Adversarial Neural Network (DANN)

## Empirical $\mathcal{H}$ -divergence

$$\hat{d}_{\mathcal{H}}(\mathcal{S}, \mathcal{T}) \stackrel{\text{def}}{=} 2 \max_{h \in \mathcal{H}} \left[ \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h(\mathbf{x}_i^{\mathcal{S}}) = 1] + \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h(\mathbf{x}_i^{\mathcal{T}}) = -1] - 1 \right].$$

We estimate the  $\mathcal{H}$ -divergence by a logistic regressor that model the probability that a given input (either  $\mathbf{x}^{\mathcal{S}}$  or  $\mathbf{x}^{\mathcal{T}}$ ) is from the source domain :

$$o(\mathbf{h}(\mathbf{x})) \stackrel{\text{def}}{=} \text{sigm}(d + \mathbf{w}^{\top} \mathbf{h}(\mathbf{x})).$$

**Given a representation output by the hidden layer  $\mathbf{h}(\cdot)$  :**

$$\hat{d}_{\mathcal{H}}(\mathbf{h}(\mathcal{S}), \mathbf{h}(\mathcal{T})) \approx 2 \max_{\mathbf{w}, d} \left[ \frac{1}{n} \sum_{i=1}^n \log(o(\mathbf{h}(\mathbf{x}_i^{\mathcal{S}}))) + \frac{1}{n} \sum_{i=1}^n \log(1 - o(\mathbf{h}(\mathbf{x}_i^{\mathcal{T}}))) - 1 \right].$$

# Domain-Adversarial Neural Network (DANN)

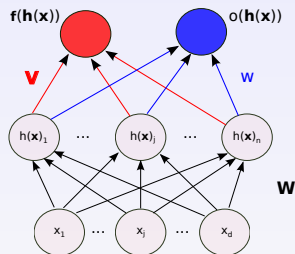
$$\min_{\mathbf{w}, \mathbf{v}, \mathbf{b}, \mathbf{c}} \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n -\log(f_{y_i^s}(\mathbf{x}_i^s))}_{\text{source loss}} + \lambda \underbrace{\max_{\mathbf{w}, d} \left( \frac{1}{n} \sum_{i=1}^n \log(o(\mathbf{h}(\mathbf{x}_i^s))) + \frac{1}{n} \sum_{i=1}^n \log(1 - o(\mathbf{h}(\mathbf{x}_i^t))) \right)}_{\text{adaptation regularizer}} \right],$$

where  $\lambda > 0$  weights the domain adaptation regularization term.

Given a **source sample**  $S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^m \sim (D_S)^m$ ,  
and a **target sample**  $T = \{(\mathbf{x}_i^t)\}_{i=1}^m \sim (D_T)^m$ ,

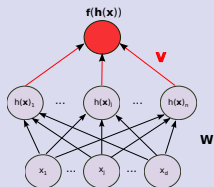
1. Pick a  $\mathbf{x}^s \in S$  and  $\mathbf{x}^t \in T$
2. Update  $\mathbf{v}$  towards  $\mathbf{f}(\mathbf{h}(\mathbf{x}^s)) = y^s$
3. Update  $\mathbf{W}$  towards  $\mathbf{f}(\mathbf{h}(\mathbf{x}^s)) = y^s$
4. Update  $\mathbf{w}$  towards  $o(\mathbf{h}(\mathbf{x}^s)) = 1$  and  $o(\mathbf{h}(\mathbf{x}^t)) = -1$
5. Update  $\mathbf{W}$  towards  $o(\mathbf{h}(\mathbf{x}^s)) = -1$  and  $o(\mathbf{h}(\mathbf{x}^t)) = 1$

**DANN finds a representation  $\mathbf{h}(\cdot)$  that are good on  $S$ ;  
but **unable to discriminate** between  $S$  and  $T$ .**

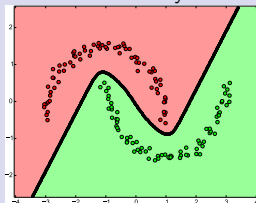


# Toy Dataset

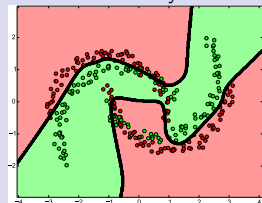
## Standard Neural Network (NN)



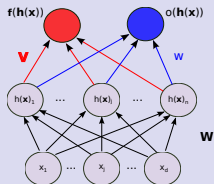
Trained to classify source



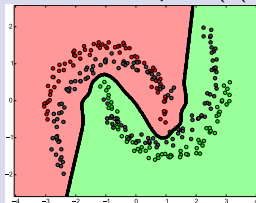
Trained to classify domains



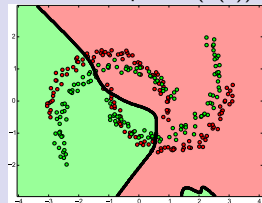
## Domain-Adversarial Neural Networks (DANN)



Classification output :  $f(h(x))$



Domain output :  $o(h(x))$



# Amazon Reviews

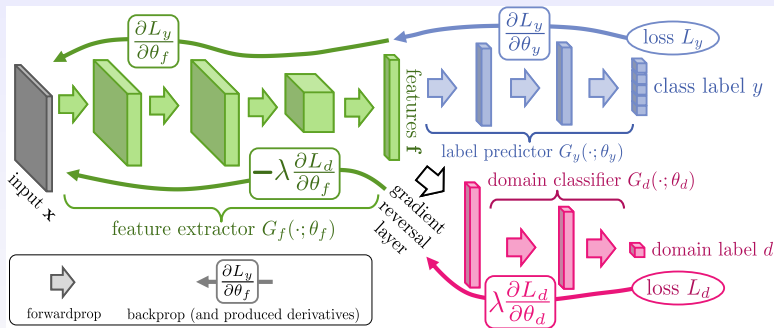
**Input** : product review (bag of words) — **Output** : positive or negative rating.

Dataset	DANN	NN
books → dvd	0.201	<b>0.199</b>
books → electronics	<b>0.246</b>	0.251
books → kitchen	<b>0.230</b>	0.235
dvd → books	<b>0.247</b>	0.261
dvd → electronics	<b>0.247</b>	0.256
dvd → kitchen	0.227	0.227
electronics → books	<b>0.280</b>	0.281
electronics → dvd	<b>0.273</b>	0.277
electronics → kitchen	<b>0.148</b>	0.149
kitchen → books	<b>0.283</b>	0.288
kitchen → dvd	0.261	0.261
kitchen → electronics	0.161	0.161

# Deeper and deeper...

To appear in JMLR : **Domain-Adversarial Neural Networks.**

by Ganin, Ustinova, Ajakan, Germain, Larochelle, Laviolette, Marchand and Lempitsky



# Plan

## 1 Basic Definitions

## 2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A Classical PAC-Bayesian Theorem
- A General PAC-Bayesian Theorem
- Transductive Learning
- Rényi-Based Theorem

## 3 Domain Adaptation Algorithms

- Ben-David et al.'s Domain Divergence
- A First PAC-Bayesian Algorithm
- A Second PAC-Bayesian Algorithm
- A Neural Network / Representation Learning Algorithm

## 4 Conclusion and future works

## An original PAC-Bayesian approach

- General theorem from which we recover existing results ;
- Modular proof, easy to adapt to various frameworks.

## Domain adaptation algorithms

- Two algorithms for linear classifiers derived from PAC-Bayesian bounds ;
- One *representation learning* Network inspired by the seminal work of Ben-David et al.



- Explore the relationships between PAC-Bayesian and *truly* Bayesian approaches ;
- Speed-up domain adaptation algorithms with stochastic gradient ;
- Go beyond simple binary classification setting ;
- Apply PAC-Bayes to your problems !

*If you only have a hammer,  
you tend to see every problem as a nail.*

— Abraham Maslow, 1966