

Generalization of the PAC-Bayesian Theory and Applications to Semi-Supervised Learning

Pascal Germain

INRIA Paris (SIERRA Team)

Modal Seminar

INRIA Lille

January 24, 2017

*Dans la vie, l'essentiel est de porter
sur tout des jugements a priori.*

— Boris Vian

1 Introduction

2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A General PAC-Bayesian Theorem
- Bounding the Majority Vote Risk

3 Semi-Supervised Learning and Variations

- Semi-Supervised Learning
- Transductive Learning
- Domain Adaptation

4 Conclusion

1 Introduction

2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A General PAC-Bayesian Theorem
- Bounding the Majority Vote Risk

3 Semi-Supervised Learning and Variations

- Semi-Supervised Learning
- Transductive Learning
- Domain Adaptation

4 Conclusion

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

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

Predictors (or hypothesis)

$$h : \mathcal{X} \rightarrow \mathcal{Y}, \quad h \in \mathcal{H}$$

Learning algorithm

$$A(S) \rightarrow h$$

Loss function

$$\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$

Empirical loss

$$\widehat{\mathcal{L}}_S^\ell(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, x_i, y_i)$$

Generalization loss

$$\mathcal{L}_D^\ell(h) = \mathbf{E}_{(x,y) \sim D} \ell(h, x, y)$$

PAC-Bayesian Theory

Initiated by McAllester (1999), the PAC-Bayesian theory gives **PAC** generalization guarantees to “**Bayesian** like” algorithms.

PAC guarantees (Probably Approximately Correct)

With probability at least “ $1-\delta$ ”, the loss of predictor h is less than “ ε ”

$$\Pr_{S \sim D^n} \left(\mathcal{L}_D^\ell(h) \leq \varepsilon(\widehat{\mathcal{L}}_S^\ell(h), n, \delta, \dots) \right) \geq 1-\delta$$

Bayesian flavor

Given:

- A **prior** distribution P on \mathcal{H} .
- A **posterior** distribution Q on \mathcal{H} .

$$\Pr_{S \sim D^n} \left(\mathbb{E}_{h \sim Q} \mathcal{L}_D^\ell(h) \leq \varepsilon \left(\mathbb{E}_{h \sim Q} \widehat{\mathcal{L}}_S^\ell(h), n, \delta, P, \dots \right) \right) \geq 1-\delta$$

A Classical PAC-Bayesian Theorem

PAC-Bayesian theorem (adapted from McAllester 1999, 2003)

For any distribution D on $\mathcal{X} \times \mathcal{Y}$, for any set of predictors \mathcal{H} , for any loss $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$, for any distribution P on \mathcal{H} , for any $\delta \in (0, 1]$, we have,

$$\Pr_{S \sim D^n} \left(\forall Q \text{ on } \mathcal{H} : \mathbf{E}_{h \sim Q} \mathcal{L}_D^\ell(h) \leq \mathbf{E}_{h \sim Q} \widehat{\mathcal{L}}_S^\ell(h) + \sqrt{\frac{1}{2n} \left[\text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right]} \right) \geq 1 - \delta,$$

where $\text{KL}(Q \| P) = \mathbf{E}_{h \sim Q} \ln \frac{Q(h)}{P(h)}$ is the **Kullback-Leibler divergence**.

Training bound

- Gives generalization guarantees **not based on testing sample**.

Valid for all posterior Q on \mathcal{H}

- Inspiration for conceiving **new learning algorithms**.

1 Introduction

2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A General PAC-Bayesian Theorem
- Bounding the Majority Vote Risk

3 Semi-Supervised Learning and Variations

- Semi-Supervised Learning
- Transductive Learning
- Domain Adaptation

4 Conclusion

- 1 Introduction
- 2 PAC-Bayesian Theory
 - Majority Vote Classifiers
 - A General PAC-Bayesian Theorem
 - Bounding the Majority Vote Risk
- 3 Semi-Supervised Learning and Variations
 - Semi-Supervised Learning
 - Transductive Learning
 - Domain Adaptation
- 4 Conclusion

Majority Vote Classifiers

Consider a binary classification problem, where $\mathcal{Y} = \{-1, +1\}$ and the set \mathcal{H} contains **binary voters** $h: \mathcal{X} \rightarrow \{-1, +1\}$

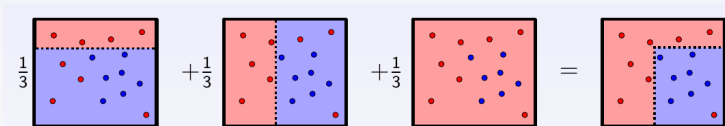
Weighted majority vote

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

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

Many learning algorithms output majority vote classifiers

AdaBoost, Random Forests, Bagging, ...



A Surrogate Loss

Majority vote risk

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

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

Gibbs Risk / Linear Loss

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

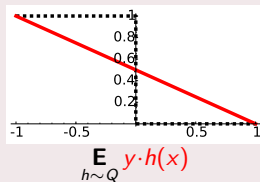
$$\begin{aligned} R_D(G_Q) &= \mathbf{E}_{(x,y) \sim D} \mathbf{E}_{h \sim Q} \mathbf{I} [h(x) \neq y] \\ &= \mathbf{E}_{h \sim Q} \mathcal{L}_D^{\ell_{01}}(h), \end{aligned}$$

where $\ell_{01}(h, x, y) = \mathbf{I}[h(x) \neq y]$.

Factor two

It is well-known that

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

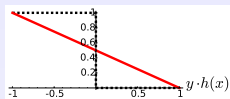


From the *Factor 2* to the \mathcal{C} -bound

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

Factor 2 bound

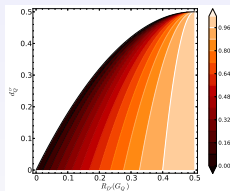
$$\begin{aligned} R_D(B_Q) &= \Pr_{(x,y) \sim D} (1 - y \cdot h(x) \geq 1) \\ &\leq \mathbf{E}_{(x,y) \sim D} (1 - y \cdot h(x)) = 2 R_D(G_Q). \end{aligned}$$



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)

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



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

$$d_Q^D = \mathbf{E}_{(x, \cdot) \sim D} \mathbf{E}_{h_i \sim Q} \mathbf{E}_{h_j \sim Q} \mathbb{I}[h_i(x) \neq h_j(x)] = \frac{1}{2} \left(1 - \mathbf{E}_{(x, \cdot) \sim D} \left[\mathbf{E}_{h \sim Q} h(x) \right]^2 \right).$$

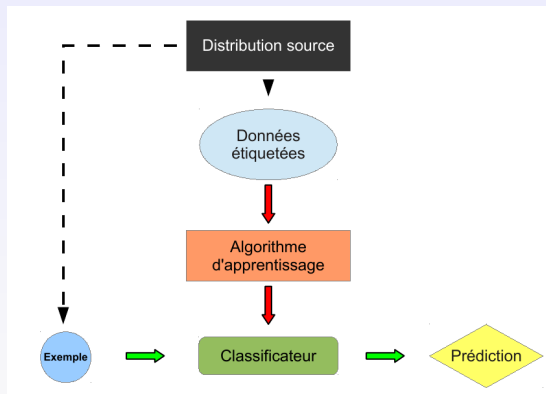
- 1 Introduction
- 2 PAC-Bayesian Theory
 - Majority Vote Classifiers
 - **A General PAC-Bayesian Theorem**
 - Bounding the Majority Vote Risk
- 3 Semi-Supervised Learning and Variations
 - Semi-Supervised Learning
 - Transductive Learning
 - Domain Adaptation
- 4 Conclusion

I.I.D. Assumption

Assumption

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

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



A General PAC-Bayesian Theorem

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

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

General theorem (Bégin et al. 2014, 2016; Germain 2015)

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 Δ -function, we have, with probability at least $1 - \delta$ over the choice of $S \sim D^n$,

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

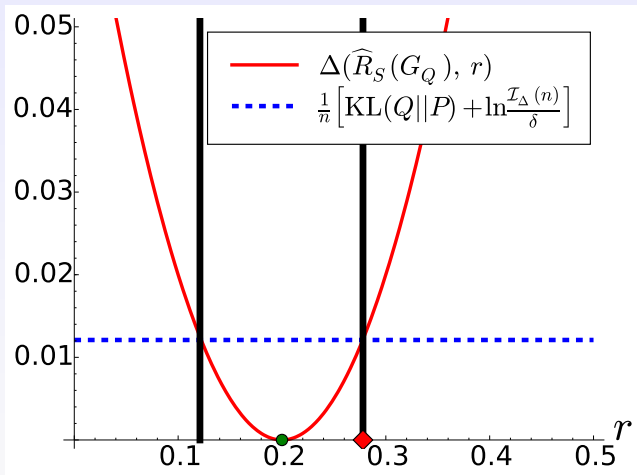
where

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

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** $\mathcal{L}_D^{\ell_{01}}(h)$:

$$\begin{aligned} \Pr_{S \sim D^n} \left(\widehat{\mathcal{L}}_S^{\ell_{01}}(h) = \frac{k}{n} \right) &= \binom{n}{k} \left(\mathcal{L}_D^{\ell_{01}}(h) \right)^k \left(1 - \mathcal{L}_D^{\ell_{01}}(h) \right)^{n-k} \\ &= \mathbf{Bin} \left(k; n, \mathcal{L}_D^{\ell_{01}}(h) \right) \end{aligned}$$

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.

$$\begin{aligned}
 & n \cdot \Delta \left(\mathbf{E}_{h \sim Q} \widehat{\mathcal{L}}_S^\ell(h), \mathbf{E}_{h \sim Q} \mathcal{L}_D^\ell(h) \right) \\
 \text{Jensen's Inequality} & \leq \mathbf{E}_{h \sim Q} n \cdot \Delta \left(\widehat{\mathcal{L}}_S^\ell(h), \mathcal{L}_D^\ell(h) \right) \\
 \text{Change of measure} & \leq \text{KL}(Q \| P) + \ln \mathbf{E}_{h \sim P} e^{n \Delta \left(\widehat{\mathcal{L}}_S^\ell(h), \mathcal{L}_D^\ell(h) \right)} \\
 \text{Markov's Inequality} & \leq_{1-\delta} \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim D^n} \mathbf{E}_{h \sim P} e^{n \Delta \left(\widehat{\mathcal{L}}_{S'}^\ell(h), \mathcal{L}_D^\ell(h) \right)} \\
 \text{Expectation swap} & = \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim D^n} e^{n \Delta \left(\widehat{\mathcal{L}}_{S'}^\ell(h), \mathcal{L}_D^\ell(h) \right)} \\
 \text{Binomial law} & = \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k=0}^n \text{Bin}(k; n, \mathcal{L}_D^\ell(h)) e^{n \Delta \left(\frac{k}{n}, \mathcal{L}_D^\ell(h) \right)} \\
 \text{Supremum over risk} & \leq \text{KL}(Q \| P) + \ln \frac{1}{\delta} \sup_{r \in [0,1]} \left[\sum_{k=0}^n \text{Bin}(k; n, r) e^{n \Delta \left(\frac{k}{n}, r \right)} \right] \\
 & = \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathcal{I}_\Delta(n). \quad \square
 \end{aligned}$$

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 $S \sim D^n$, for all Q on \mathcal{H} :

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

$$(b) \quad R_D(G_Q) \leq \hat{R}_S(G_Q) + \sqrt{\frac{1}{2n} \left[\text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right]}, \quad (\text{McAllester 1999, 2003})$$

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

$$(d) \quad R_D(G_Q) \leq \hat{R}_S(G_Q) + \frac{1}{\lambda} \left[\text{KL}(Q \| P) + \ln \frac{1}{\delta} + f(\lambda, n) \right]. \quad (\text{Alquier et al. 2015})$$

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

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

$$\Delta_\lambda(q, p) = \frac{\lambda}{n} (p - q).$$

- 1 Introduction
- 2 PAC-Bayesian Theory
 - Majority Vote Classifiers
 - A General PAC-Bayesian Theorem
 - **Bounding the Majority Vote Risk**
- 3 Semi-Supervised Learning and Variations
 - Semi-Supervised Learning
 - Transductive Learning
 - Domain Adaptation
- 4 Conclusion

Bounding the Expected Disagreement

Expected disagreement

$$d_Q^D = \mathbf{E}_{(x, \cdot) \sim D} \mathbf{E}_{h_i \sim Q} \mathbf{E}_{h_j \sim Q} \mathbf{I} \left[h_i(x) \neq h_j(x) \right] = \mathbf{E}_{(x, \cdot) \sim D} \mathbf{E}_{h_{ij} \sim Q^2} \ell_d(h_{ij}, x, \cdot),$$

$$\text{where } Q^2(h_{ij}) = Q(h_i) Q(h_j) \implies \text{KL}(Q^2 \| P^2) = 2 \text{KL}(Q \| P).$$

General theorem

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

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

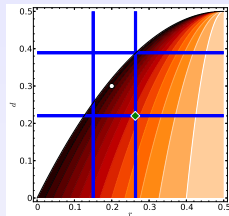
Corollary

$$(a) \quad \text{kl} \left(\widehat{d}_Q^S, d_Q^D \right) \leq \frac{1}{n} \left[2 \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right].$$

The \mathcal{C} -bound

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

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



PAC-Bayes \mathcal{C} -bound 1

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

$$\forall Q \text{ on } \mathcal{H} : R_D(B_Q) \leq 1 - \frac{(1 - 2 \cdot \bar{r})^2}{1 - 2 \cdot \underline{d}},$$

with

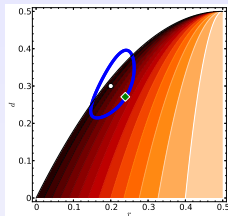
$$\bar{r} = \max_{r \in [0, \frac{1}{2}]} \left\{ \Delta(\hat{R}_S(G_Q), r) \leq \frac{1}{n} \left[2\text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta/2} \right] \right\},$$

$$\underline{d} = \min_{d \in [0, \frac{1}{2}]} \left\{ \Delta(\hat{d}_Q^S, d) \leq \frac{1}{n} \left[2\text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta/2} \right] \right\}.$$

The \mathcal{C} -bound

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

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



PAC-Bayes \mathcal{C} -bound 2

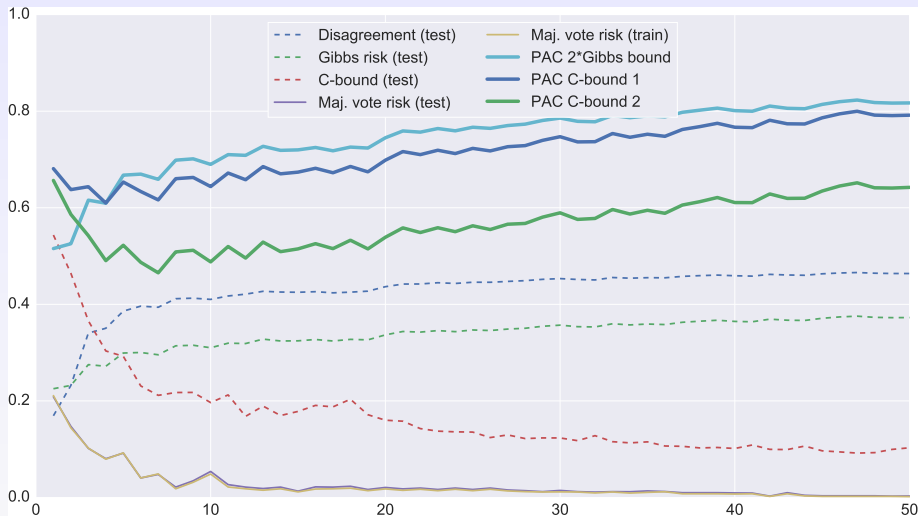
[...] with probability at least $1 - \delta$ over the choice of $S \sim D^n$,

$$\forall Q \text{ on } \mathcal{H} : R_D(B_Q) \leq \sup_{(r,d) \in \mathcal{A}} \left\{ 1 - \frac{(1 - 2 \cdot r)^2}{1 - 2 \cdot d} \right\},$$

with

$$\mathcal{A} = \left\{ (r, d) \in [0, \frac{1}{2}] \mid \Delta_2 \left((\hat{R}_S(G_Q), r), (\hat{d}_Q^S, d) \right) \leq \frac{1}{n} \left[2\text{KL}(Q \| P) + \ln \frac{\mathcal{I}_{\Delta_2}(n)}{\delta} \right] \right\}.$$

Bounds Values (adaboost iterates)



1 Introduction

2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A General PAC-Bayesian Theorem
- Bounding the Majority Vote Risk

3 Semi-Supervised Learning and Variations

- Semi-Supervised Learning
- Transductive Learning
- Domain Adaptation

4 Conclusion

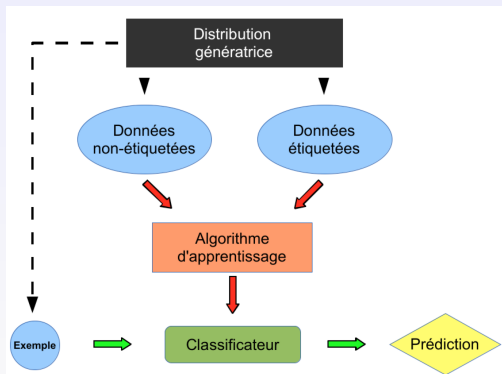
- 1 Introduction
- 2 PAC-Bayesian Theory
 - Majority Vote Classifiers
 - A General PAC-Bayesian Theorem
 - Bounding the Majority Vote Risk
- 3 Semi-Supervised Learning and Variations
 - Semi-Supervised Learning
 - Transductive Learning
 - Domain Adaptation
- 4 Conclusion

Semi-Supervised Learning

Assumption

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

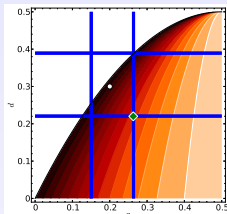
$$\begin{aligned} S &= \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim D^n \\ U &= \{ (x_{n+1}, \cdot), (x_{n+2}, \cdot), \dots, (x_{n+n'}, \cdot) \} \sim D^{n'} \end{aligned}$$



Semi-Supervised Learning

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

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



PAC \mathcal{C} -bound semi-supervised

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

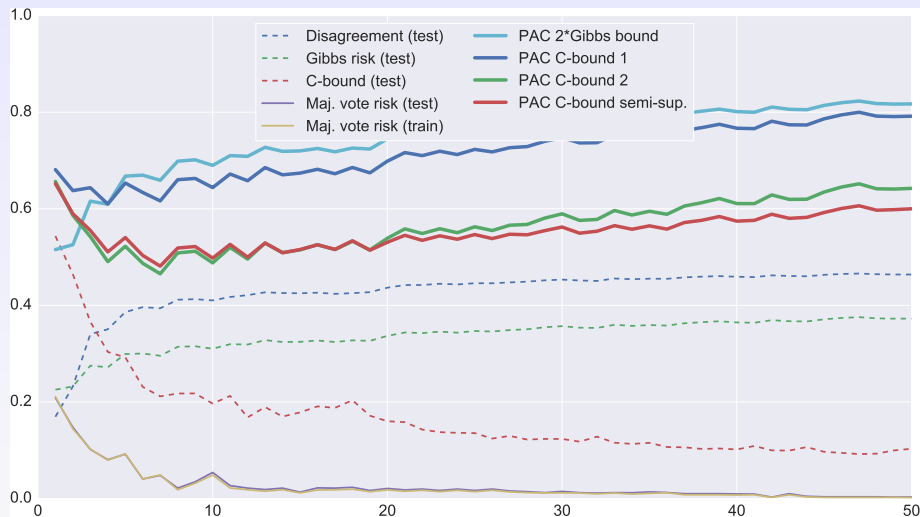
$$\forall Q \text{ on } \mathcal{H} : R_D(B_Q) \leq 1 - \frac{(1 - 2 \cdot \bar{r})^2}{1 - 2 \cdot \underline{d}},$$

with

$$\bar{r} = \max_{r \in [0, \frac{1}{2}]} \left\{ \Delta(\hat{R}_S(G_Q), r) \leq \frac{1}{n} \left[2\text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n)}{\delta/2} \right] \right\},$$

$$\underline{d} = \min_{d \in [0, \frac{1}{2}]} \left\{ \Delta(\hat{d}_Q^{SUU}, d) \leq \frac{1}{n+n'} \left[2\text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(n+n')}{\delta/2} \right] \right\}.$$

Bounds Values (adaboost iterates)



1 Introduction

2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A General PAC-Bayesian Theorem
- Bounding the Majority Vote Risk

3 Semi-Supervised Learning and Variations

- Semi-Supervised Learning
- **Transductive Learning**
- Domain Adaptation

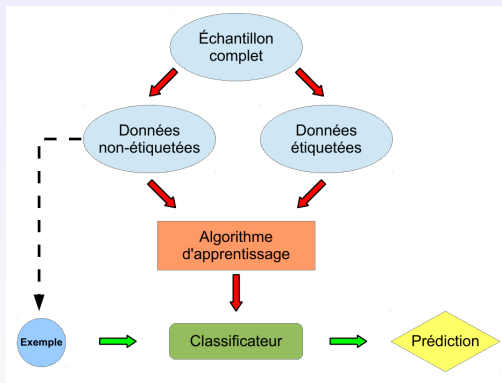
4 Conclusion

Transductive Learning

Assumption

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

$$S = \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \subset Z$$
$$U = \{ (x_{n+1}, \cdot), (x_{n+2}, \cdot), \dots, (x_N, \cdot) \} = Z \setminus S$$



Transductive Learning

Assumption

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

$$\begin{aligned} S &= \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \subset Z \\ U &= \{ (x_{n+1}, \cdot), (x_{n+2}, \cdot), \dots, (x_N, \cdot) \} = Z \setminus S \end{aligned}$$

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

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

Theorem

(Bégin et al. 2014)

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

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

where

$$\mathcal{T}_\Delta(n, N) = \max_{K=0 \dots N} \left[\sum_{k=\max[0, K+n-N]}^{\min[n, K]} \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}} e^{n\Delta(\frac{k}{n}, \frac{K}{N})} \right].$$

Theorem

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

Proof.

$$n \cdot \Delta \left(\mathbf{E}_{h \sim Q} \widehat{\mathcal{L}}_S^\ell(h), \mathbf{E}_{h \sim Q} \widehat{\mathcal{L}}_Z^\ell(h) \right)$$

Jensen's inequality

$$\leq \mathbf{E}_{h \sim Q} n \cdot \Delta \left(\widehat{\mathcal{L}}_S^\ell(h), \widehat{\mathcal{L}}_Z^\ell(h) \right)$$

Change of measure

$$\leq \text{KL}(Q \| P) + \ln \mathbf{E}_{h \sim P} e^{n \Delta \left(\widehat{\mathcal{L}}_S^\ell(h), \widehat{\mathcal{L}}_Z^\ell(h) \right)}$$

Markov's inequality

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

Expectations swap

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

Hypergeometric law

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_k \frac{\binom{N \cdot \widehat{\mathcal{L}}_Z^\ell(h)}{k} \binom{N - N \cdot \widehat{\mathcal{L}}_Z^\ell(h)}{n-k}}{\binom{N}{n}} e^{n \cdot \Delta \left(\frac{k}{n}, \widehat{\mathcal{L}}_Z^\ell(h) \right)}$$

Supremum over risk

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

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

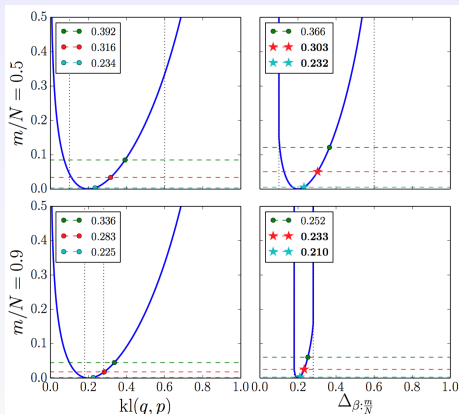
□

A New Transductive Bound

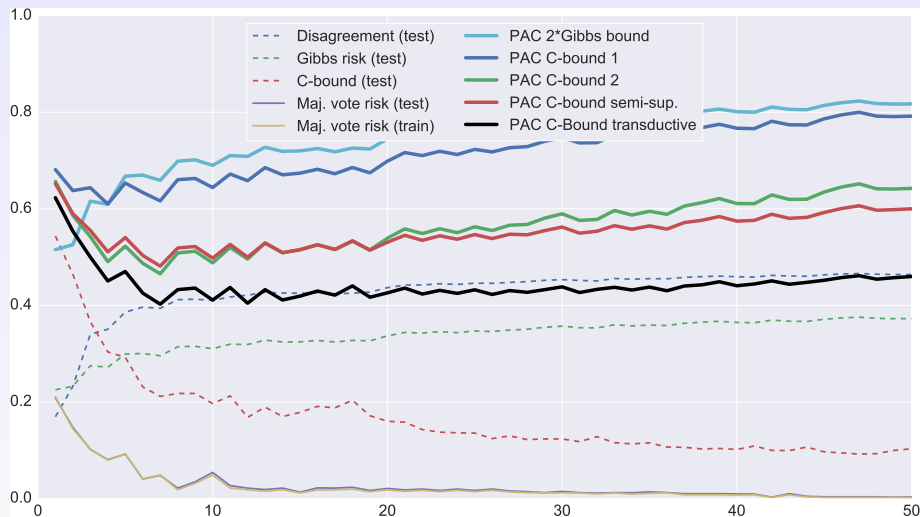
A new Δ -function

(Bégin et al. 2014)

$$\Delta_{\beta}(q, p) = \text{kl}(q, p) + \frac{1-\beta}{\beta} \text{kl}\left(\frac{p-\beta q}{1-\beta}, p\right).$$



Bounds Values (adaboost iterates)



1 Introduction

2 PAC-Bayesian Theory

- Majority Vote Classifiers
- A General PAC-Bayesian Theorem
- Bounding the Majority Vote Risk

3 Semi-Supervised Learning and Variations

- Semi-Supervised Learning
- Transductive Learning
- Domain Adaptation

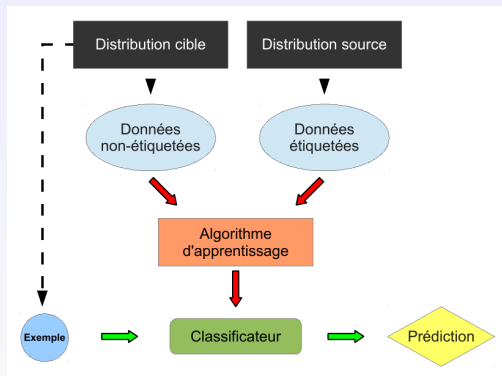
4 Conclusion

Domain Adaptation

Assumption

Source and target examples are generated by different distributions.

$$S = \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim (D_S)^n$$
$$T = \{ (x_1, \cdot), (x_2, \cdot), \dots, (x_n, \cdot) \} \sim (D_T)^{n'}$$



Our Domain Adaptation Setting

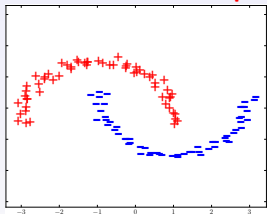
Assumption

Source and target examples are generated by different distributions.

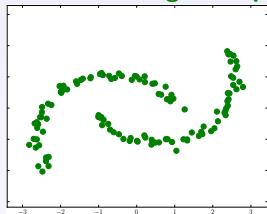
$$\begin{aligned} S &= \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim (D_S)^n \\ T &= \{ (x_1, \cdot), (x_2, \cdot), \dots, (x_n, \cdot) \} \sim (D_T)^{n'} \end{aligned}$$

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

a **labeled source sample**



an **unlabeled target sample**



The goal is to build a classifier with a low **target risk** $R_{D_T}(h)$

Generalization Bound

Observation

$$R_D(G_Q) = \frac{1}{2}d_Q^D + e_Q^D = \frac{1}{2}d_Q^D + \mathbf{E}_{(x,y) \sim D} \mathbf{E}_{h_i \sim Q} \mathbf{E}_{h_j \sim Q} \mathbb{I}[h_i(x) \neq y] \mathbb{I}[h_j(x) \neq y].$$

PAC-Bayesian DA Bound

(Germain, Habrard, et al. 2016)

For any prior P over \mathcal{H} , any $\delta \in (0, 1]$, any real numbers $b > 1$ and $c > 1$, with a probability at least $1 - \delta$ over the choices of $S \sim (D_S)^n$ and $T \sim (D_T)^{n'}$, we have

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

$$R_{D_T}(G_Q) \leq c \times \frac{1}{2} \underbrace{\widehat{d}_Q^T}_{\text{empirical target disagreement}} + b \times \underbrace{\beta_\infty(D_T \| D_S)}_{\text{domain divergence}} \underbrace{\widehat{e}_Q^S}_{\text{empirical source joint error}} + \underbrace{O(\text{KL}(Q \| P) + \ln \frac{1}{\delta})}_{\text{complexity term}}.$$

Linear trade-off between \widehat{d}_Q^T and \widehat{e}_Q^S :

\Rightarrow We consider $\beta_\infty(D_T \| D_S) = \sup_x \frac{D_T(x)}{D_S(x)}$ as a parameter to tune.

Learning algorithm for Linear Classifiers

As many PAC-Bayesian works (since Langford and Shawe-Taylor 2002) :

- We consider the set \mathcal{H} of all linear classifiers $h_{\mathbf{v}}$ in $\mathcal{X} = \mathbb{R}^d$:

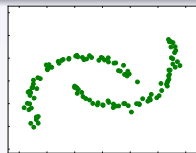
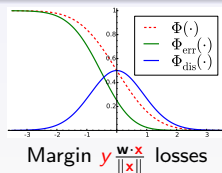
$$h_{\mathbf{v}}(\mathbf{x}) = \text{sgn}(\mathbf{v} \cdot \mathbf{x}).$$

- Let $Q_{\mathbf{w}}$ on \mathcal{H} be a Gaussian distribution centered on \mathbf{w} (with $\Sigma = \mathbf{I}_d$):

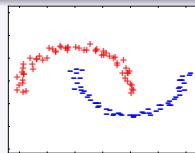
$$h_{\mathbf{w}}(\mathbf{x}) = \text{sgn} \left[\mathbf{E}_{\mathbf{v} \sim Q_{\mathbf{w}}} h_{\mathbf{v}}(\mathbf{x}) \right].$$

Given $T = \{\mathbf{x}_i\}_{i=1}^{n'}$ and $S = \{(\mathbf{x}_j, y_j)\}_{j=1}^n$, find $\mathbf{w} \in \mathbb{R}^d$ that minimizes:

$$\underbrace{C \times \widehat{d}_Q^T}_{\frac{C}{n'} \sum_i \Phi_{\text{dis}}\left(\frac{\mathbf{w} \cdot \mathbf{x}_i}{\|\mathbf{x}_i\|}\right)} + \underbrace{B \times \widehat{e}_Q^S}_{\frac{B}{n} \sum_j \Phi_{\text{err}}\left(y_j \frac{\mathbf{w} \cdot \mathbf{x}_j}{\|\mathbf{x}_j\|}\right)} + \underbrace{\text{KL}(Q_{\mathbf{w}} \| P_0)}_{\frac{1}{2} \|\mathbf{w}\|^2}.$$



Low density region on target

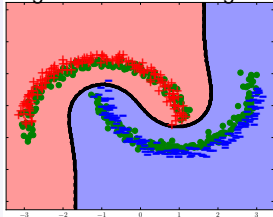


Classification accuracy on source

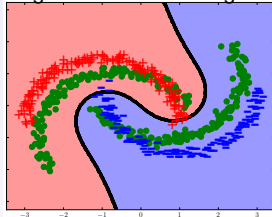
Toy Experiment

- RBF kernel
- $B = 1$
- $C = 1$

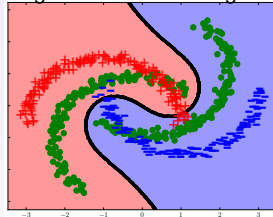
Target rotation of 10 degrees



Target rotation of 30 degrees



Target rotation of 50 degrees



Empirical results on Amazon Dataset

- Linear kernel
- Hyper-parameter selection by reverse cross-validation

	svm	dasvm	coda	PBDA	DALC
books→DVDs	<i>0.179</i>	0.193	0.181	0.183	0.178
books→electro	0.290	<i>0.226</i>	0.232	0.263	0.212
books→kitchen	0.251	0.179	0.215	0.229	<i>0.194</i>
DVDs→books	0.203	0.202	0.217	<i>0.197</i>	0.186
DVDs→electro	0.269	0.186	<i>0.214</i>	0.241	0.245
DVDs→kitchen	0.232	0.183	<i>0.181</i>	0.186	0.175
electro→books	0.287	0.305	0.275	0.232	<i>0.240</i>
electro→DVDs	0.267	0.214	0.239	<i>0.221</i>	0.256
electro→kitchen	<i>0.129</i>	0.149	0.134	0.141	0.123
kitchen→books	0.267	0.259	<i>0.247</i>	<i>0.247</i>	0.236
kitchen→DVDs	0.253	0.198	0.238	0.233	<i>0.225</i>
kitchen→electro	0.149	0.157	0.153	0.129	<i>0.131</i>
Average	0.231	<i>0.204</i>	0.210	0.208	0.200

- 1 Introduction
- 2 PAC-Bayesian Theory
 - Majority Vote Classifiers
 - A General PAC-Bayesian Theorem
 - Bounding the Majority Vote Risk
- 3 Semi-Supervised Learning and Variations
 - Semi-Supervised Learning
 - Transductive Learning
 - Domain Adaptation
- 4 Conclusion

An original PAC-Bayesian approach

- General theorem from which we recover existing results;
- Modular proof, easy to adapt to various frameworks (*i.e.*, transductive learning).

The virtue of disagreement

- Second-order statistic from unlabeled sample (useful in semi-supervised setting and variants);
- Allows to reduce the value of PAC-Bayesian bounds (\mathcal{C} -bound);
- Also used in a domain adaptation learning algorithm.

Some (self-)pointers

Our 74 pages journal paper (JMLR)

- Risk Bounds for the Majority Vote: From a PAC-Bayesian Analysis to a Learning Algorithm (Germain, Lacasse, et al. 2015)

My PhD thesis (in french)

- Généralisations de la théorie PAC-bayésienne pour l'apprentissage inductif, l'apprentissage transductif et l'adaptation de domaine (Germain 2015)
<http://www.di.ens.fr/~germain/publis/these.pdf>

Recent papers

- ICML: A New PAC-Bayesian Perspective on Domain Adaptation (Germain, Habrard, et al. 2016)
- AISTATS: PAC-Bayesian Bounds Based on the Rényi Divergence (Bégin et al. 2016)
- NIPS: PAC-Bayesian Theory Meets Bayesian Inference (Germain, Bach, et al. 2016)

Even some code on my GitHub....

- <https://github.com/pgermain>

References I

- Alquier, Pierre, James Ridgway, and Nicolas Chopin (2015). “On the properties of variational approximations of Gibbs posteriors”. In: *ArXiv e-prints*. url: <http://arxiv.org/abs/1506.04091>.
- Bégin, Luc, Pascal Germain, François Laviolette, and Jean-François Roy (2014). “PAC-Bayesian Theory for Transductive Learning”. In: *AISTATS*.
- (2016). “PAC-Bayesian Bounds based on the Rényi Divergence”. In: *AISTATS*.
- Catoni, Olivier (2007). *PAC-Bayesian supervised classification: the thermodynamics of statistical learning*. Vol. 56. Inst. of Mathematical Statistic.
- Germain, Pascal (2015). “Généralisations de la théorie PAC-bayésienne pour l'apprentissage inductif, l'apprentissage transductif et l'adaptation de domaine.” PhD thesis. Université Laval. url: <http://www.theses.ulaval.ca/2015/31774/>.
- Germain, Pascal, Francis Bach, Alexandre Lacoste, and Simon Lacoste-Julien (2016). “PAC-Bayesian Theory Meets Bayesian Inference”. In: *ArXiv e-prints*. url: <http://arxiv.org/abs/1605.08636>.
- Germain, Pascal, Amaury Habrard, François Laviolette, and Emilie Morvant (2016). “A New PAC-Bayesian Perspective on Domain Adaptation”. In: *ICML*. url: <http://arxiv.org/abs/1506.04573>.
- Germain, Pascal, Alexandre Lacasse, François Laviolette, Mario Marchand, and Jean-François Roy (2015). “Risk Bounds for the Majority Vote: From a PAC-Bayesian Analysis to a Learning Algorithm”. In: *JMLR 16*.

References II

- Langford, John and Matthias Seeger (2001). *Bounds for averaging classifiers*. Tech. rep. Carnegie Mellon, Department of Computer Science.
- Langford, John and John Shawe-Taylor (2002). "PAC-Bayes & Margins". In: *NIPS*.
- McAllester, David (1999). "Some PAC-Bayesian Theorems". In: *Machine Learning* 37.3.
- (2003). "PAC-Bayesian Stochastic Model selection". In: *Machine Learning* 51.1.