

Apprentissage de réseaux de neurones à activations binaires avec garanties statistiques

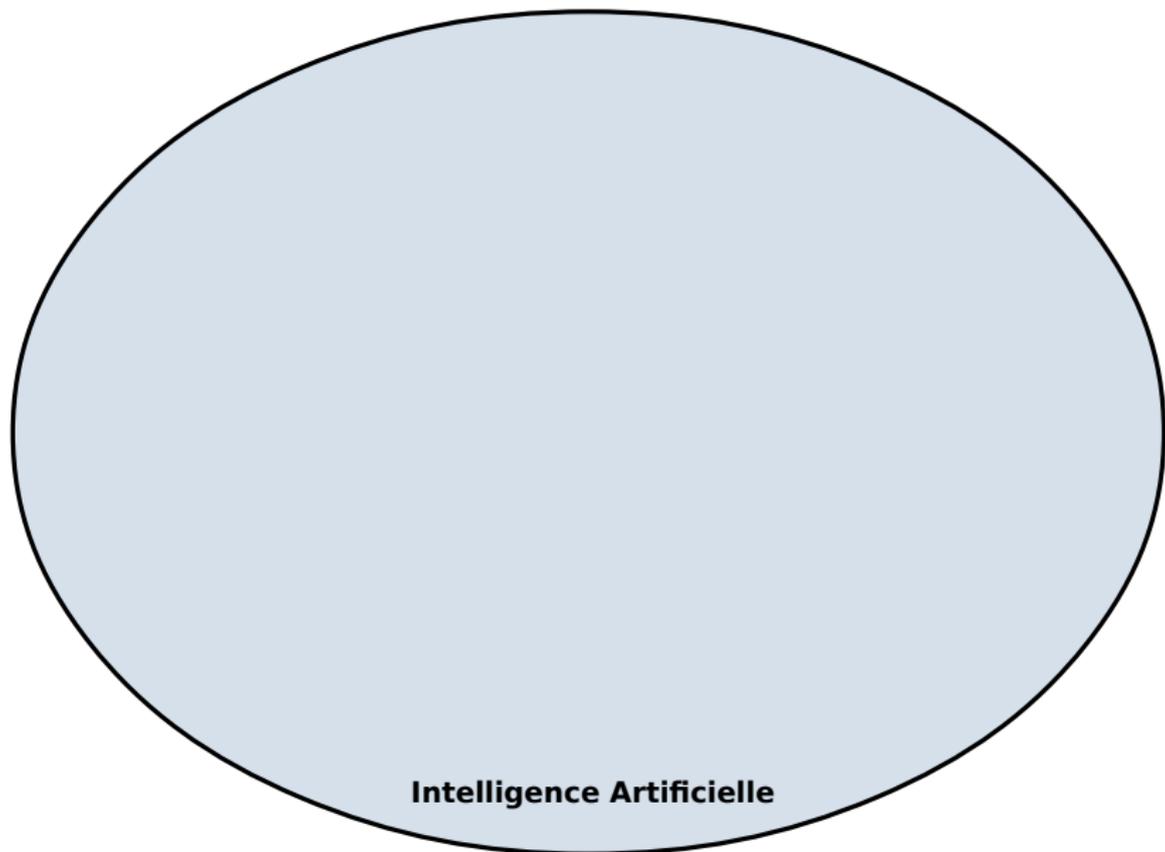
Gaël Letarte

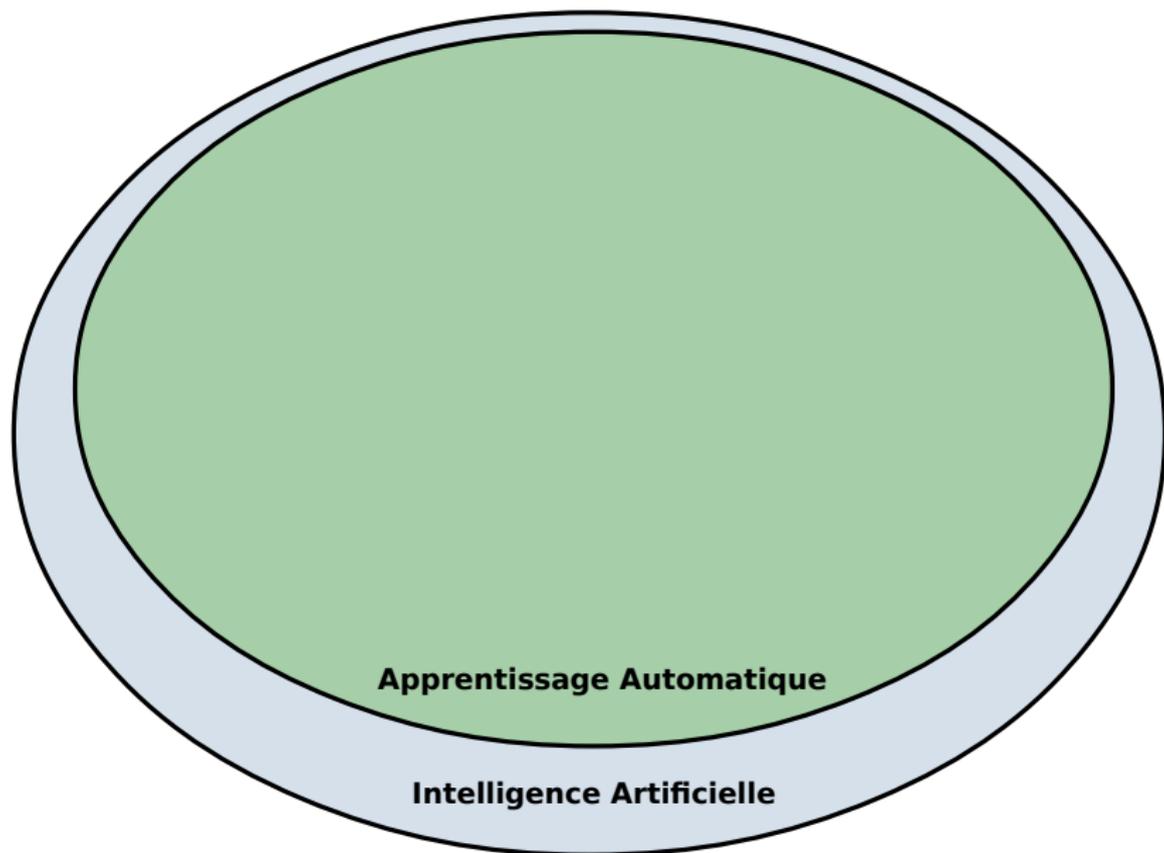
En collaboration avec Pascal Germain, Benjamin Guedj et François Laviolette

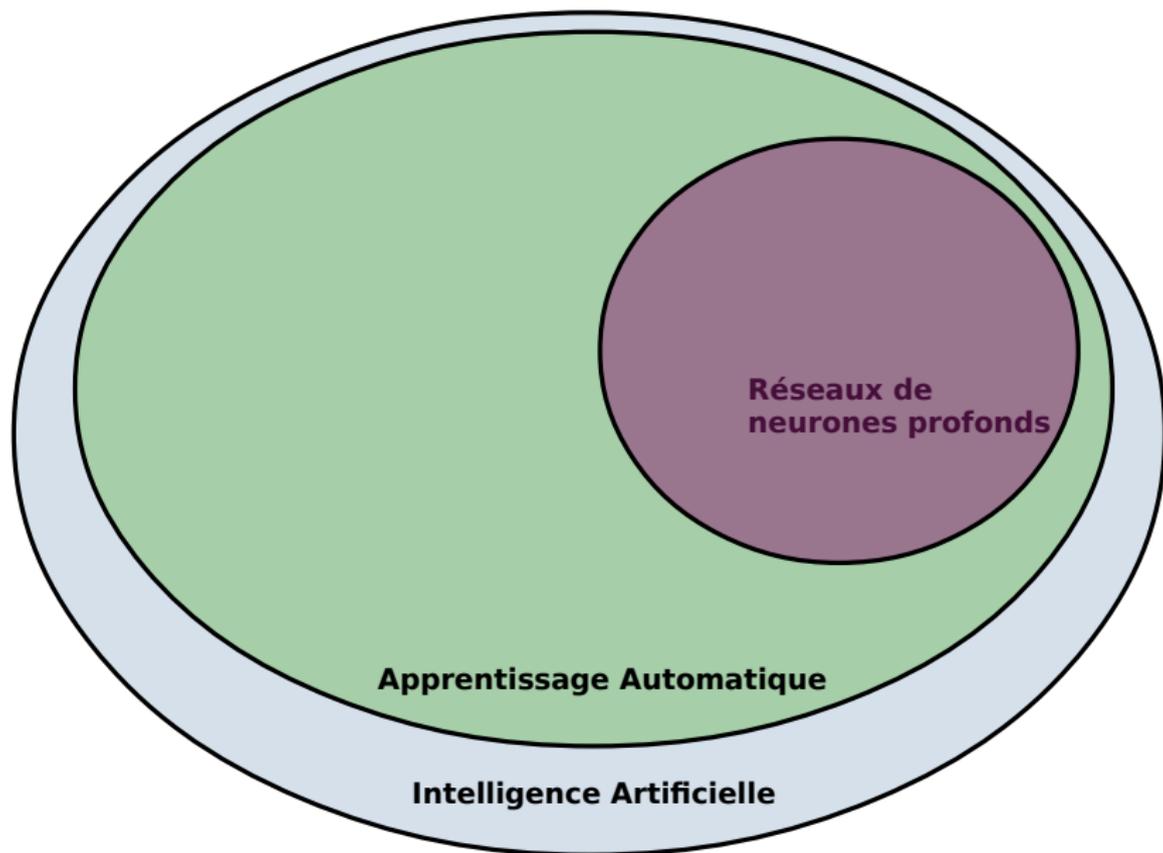
Département d'informatique et de génie logiciel, Université Laval

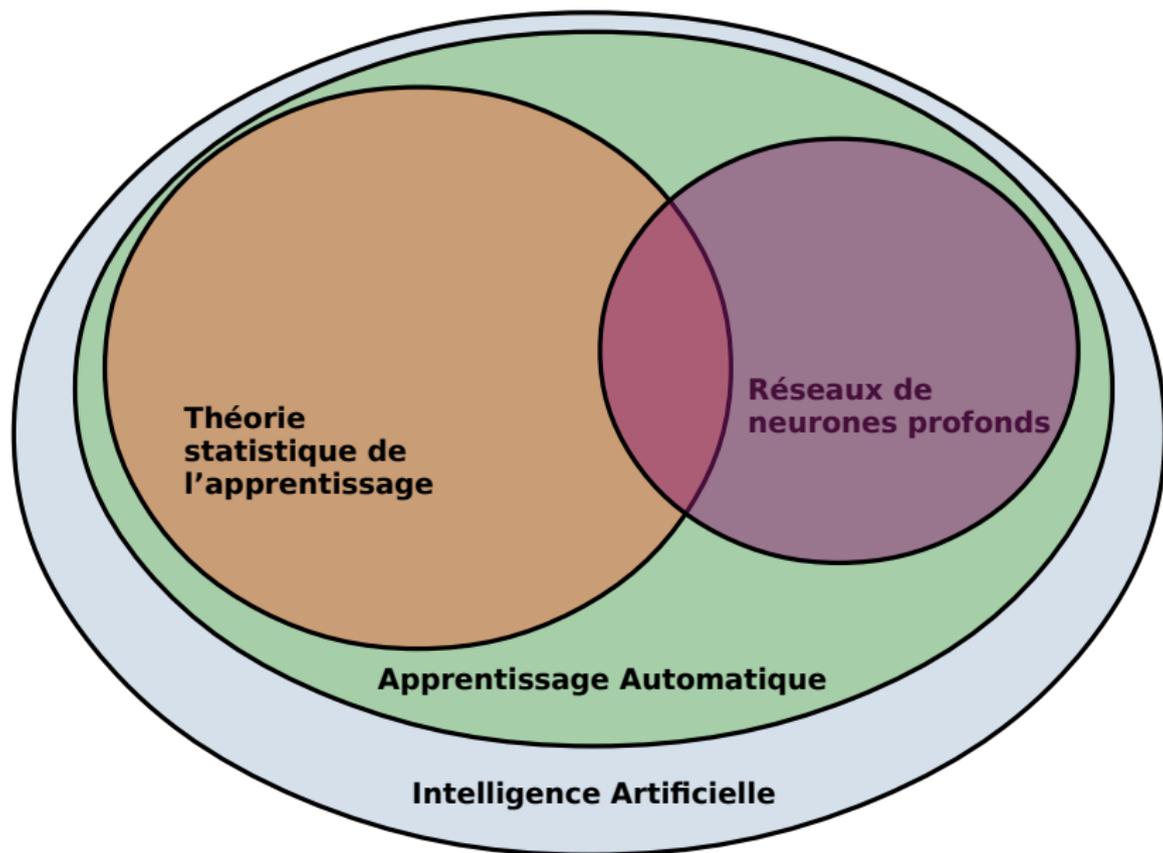
17 Janvier, 2020

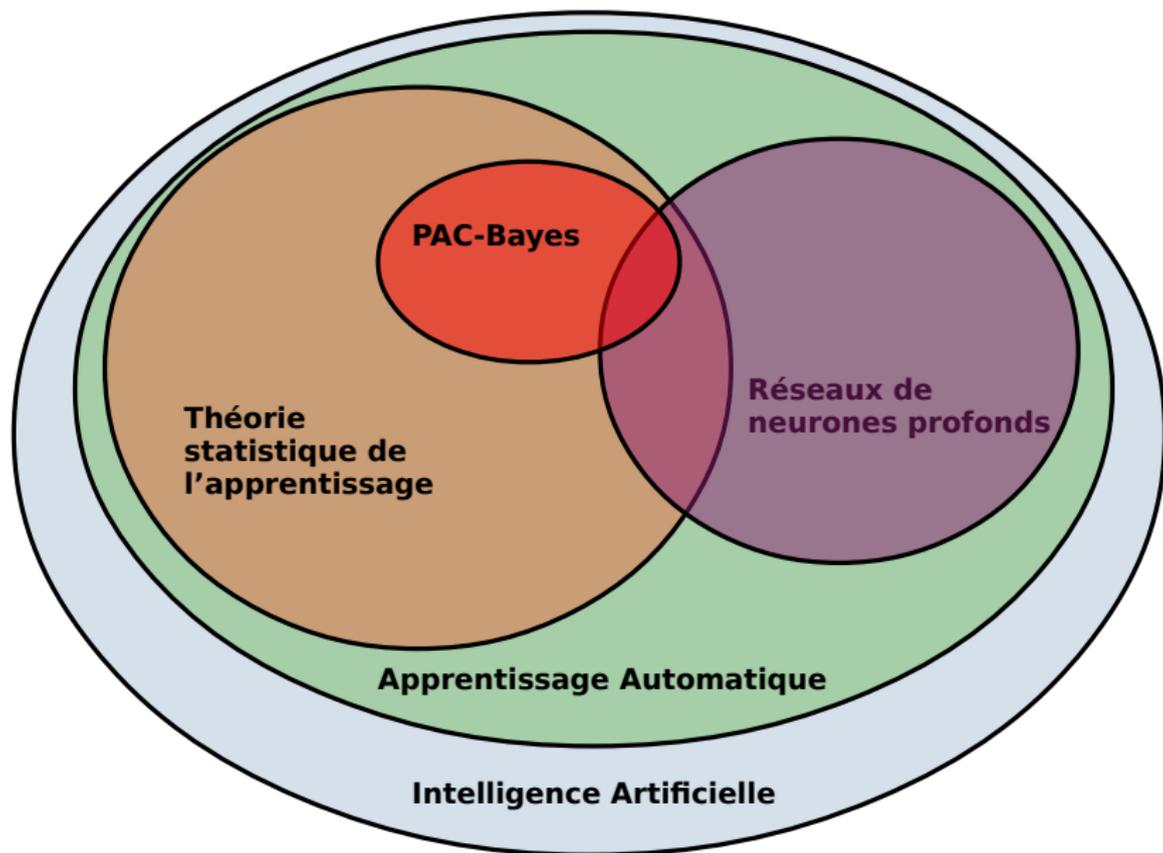


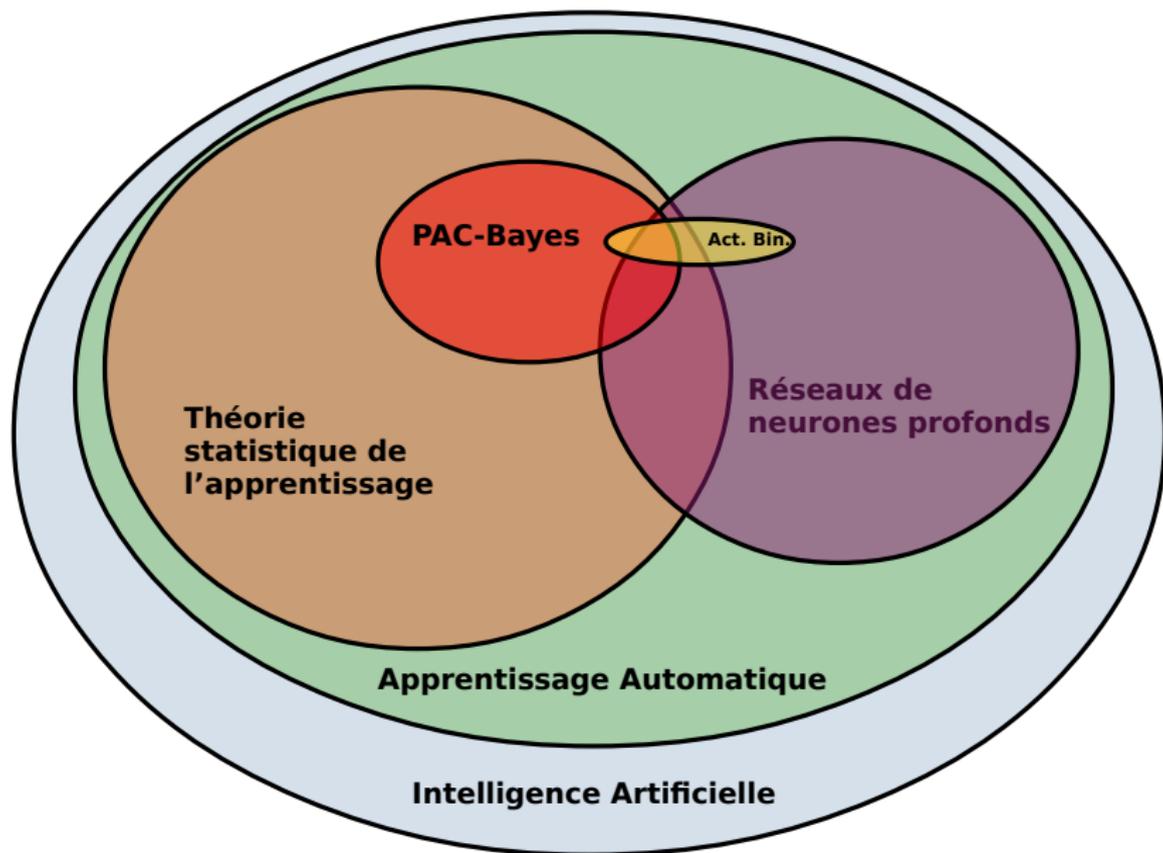


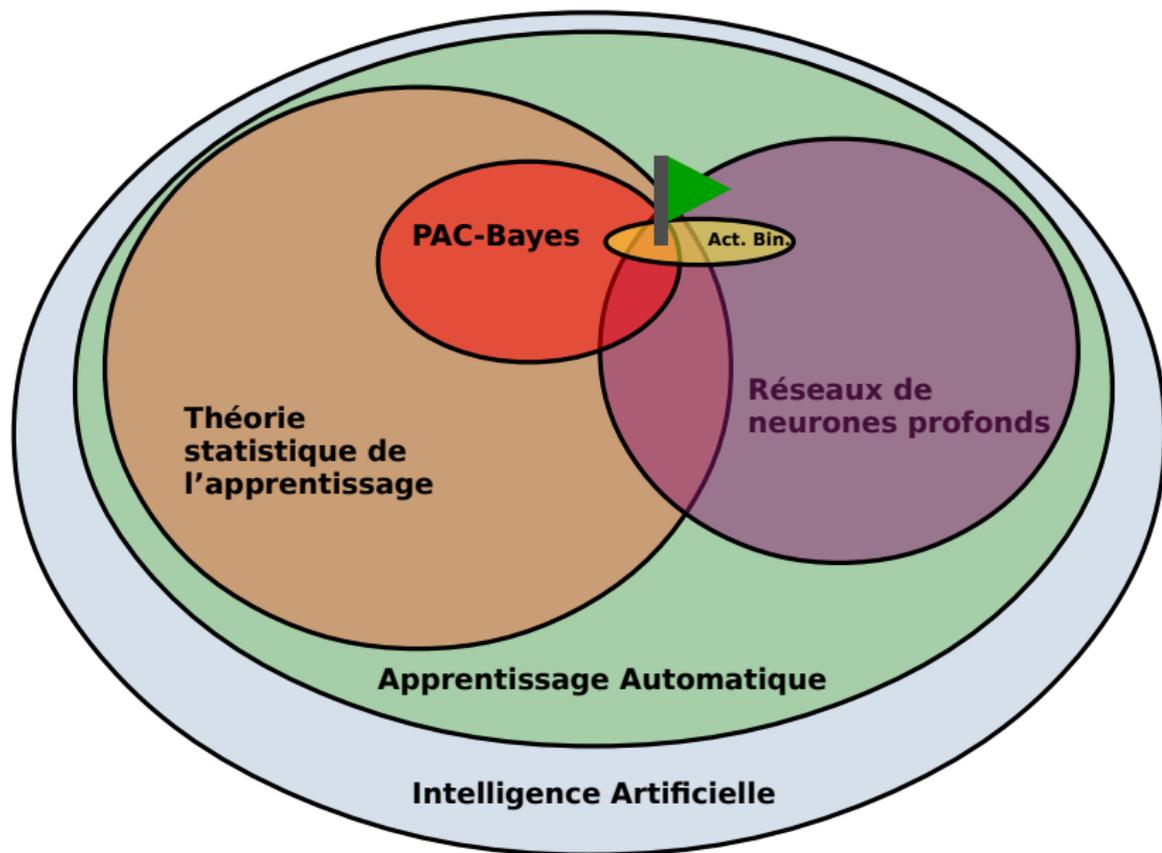


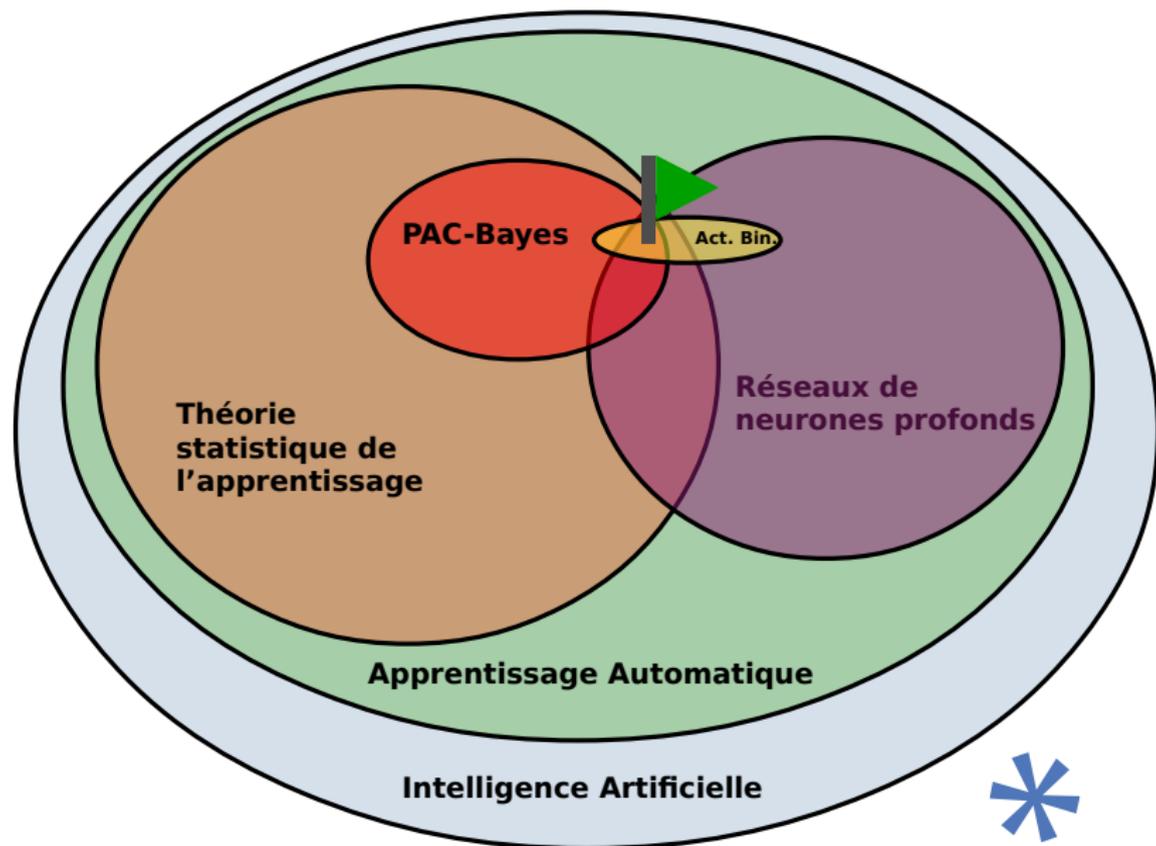


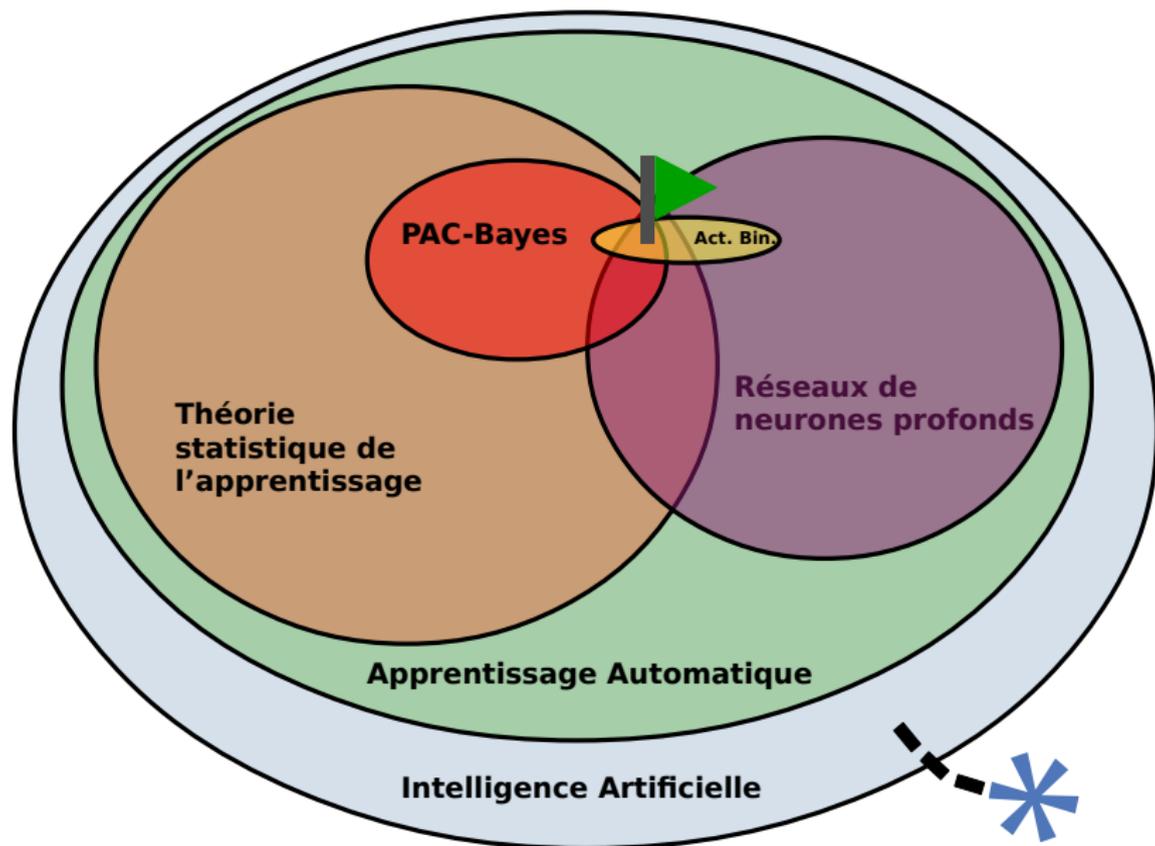




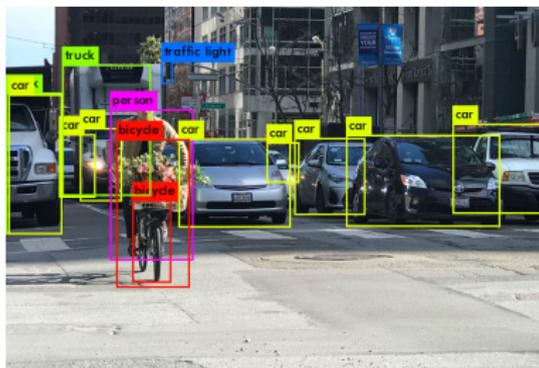








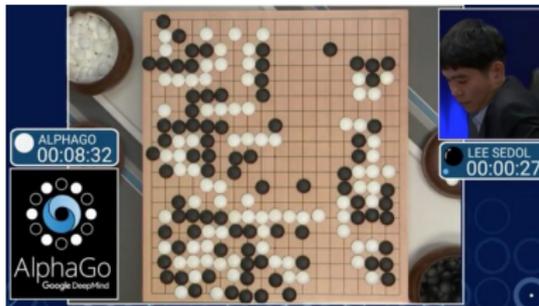
Des succès!



Crédit: https://medium.com/@jonathan_hui/



Crédit: www.google.com



Crédit: <https://hicomm.bg/>



Crédit: <https://entertainment.ie/>

Des échecs...

The New York Times

Apple Card Investigated After Gender Discrimination Complaints

A prominent software developer said on Twitter that the credit card was “sexist” against women applying for credit.

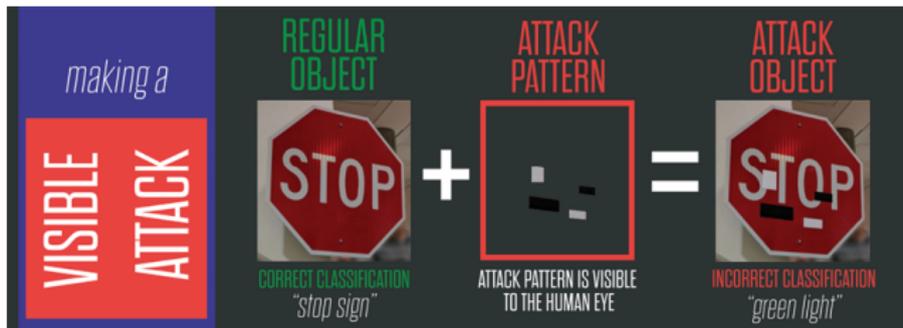


EXCLUSIVE

STAT+

IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show

By CASEY ROSS @caseyross and IKE SWETLITZ / JULY 25, 2018



Crédit: Comiter, Marcus. "Attacking Artificial Intelligence: AI's Security Vulnerability and What Policymakers Can Do About It."

Des garanties?

The New York Times

Apple Card Investigated After Gender Discrimination Complaints

A prominent software developer said on Twitter that the credit card was "sexist".

Équité

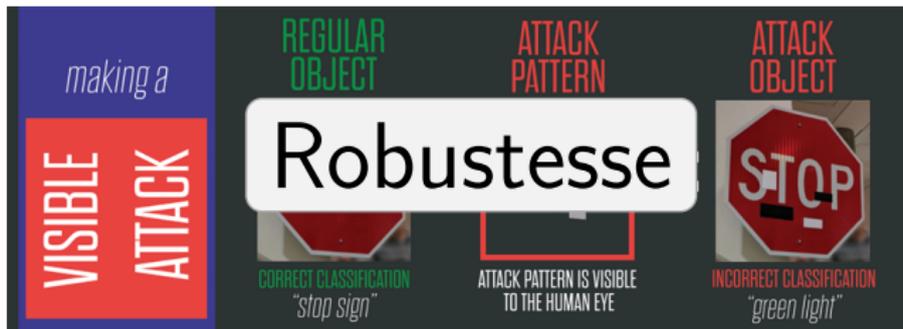
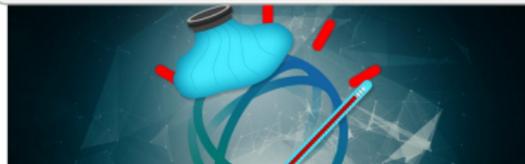


EXCLUSIVE

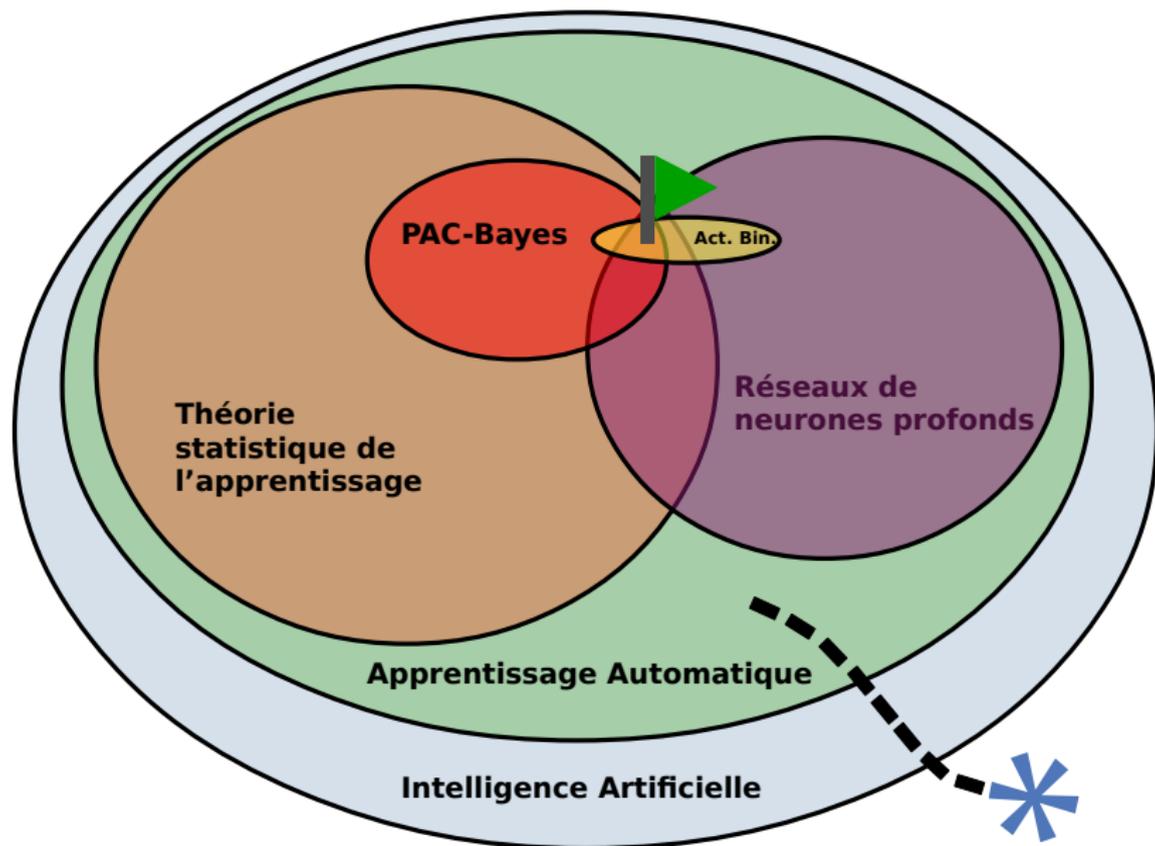
STAT+

IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal

Généralisation



Crédit: Comiter, Marcus. "Attacking Artificial Intelligence: AI's Security Vulnerability and What Policymakers Can Do About It."



Apprentissage automatique: Exemple

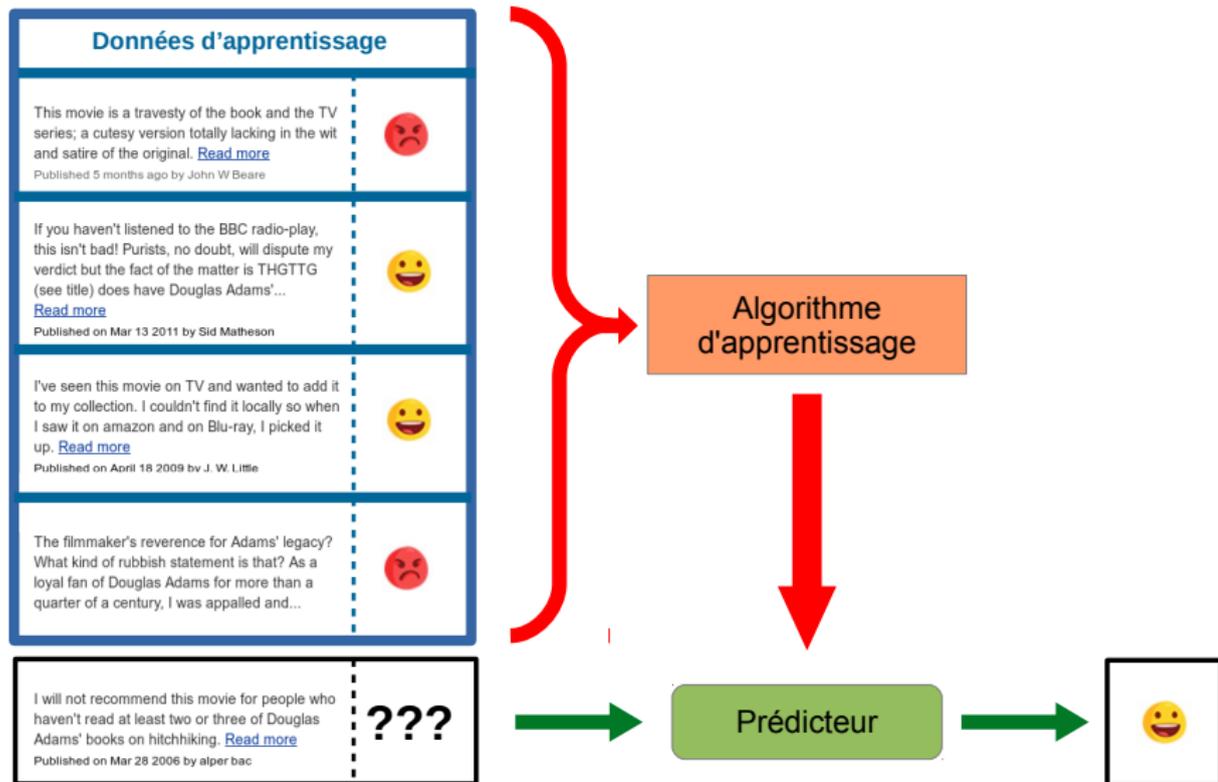
Données d'apprentissage	
<p>This movie is a travesty of the book and the TV series; a cutesy version totally lacking in the wit and satire of the original. Read more</p> <p>Published 5 months ago by John W Beare</p>	
<p>If you haven't listened to the BBC radio-play, this isn't bad! Purists, no doubt, will dispute my verdict but the fact of the matter is THGTTG (see title) does have Douglas Adams'... Read more</p> <p>Published on Mar 13 2011 by Sid Matheson</p>	
<p>I've seen this movie on TV and wanted to add it to my collection. I couldn't find it locally so when I saw it on amazon and on Blu-ray, I picked it up. Read more</p> <p>Published on April 18 2009 by J. W. Little</p>	
<p>The filmmaker's reverence for Adams' legacy? What kind of rubbish statement is that? As a loyal fan of Douglas Adams for more than a quarter of a century, I was appalled and...</p>	
<p>I will not recommend this movie for people who haven't read at least two or three of Douglas Adams' books on hitchhiking. Read more</p> <p>Published on Mar 28 2006 by alper bac</p>	???



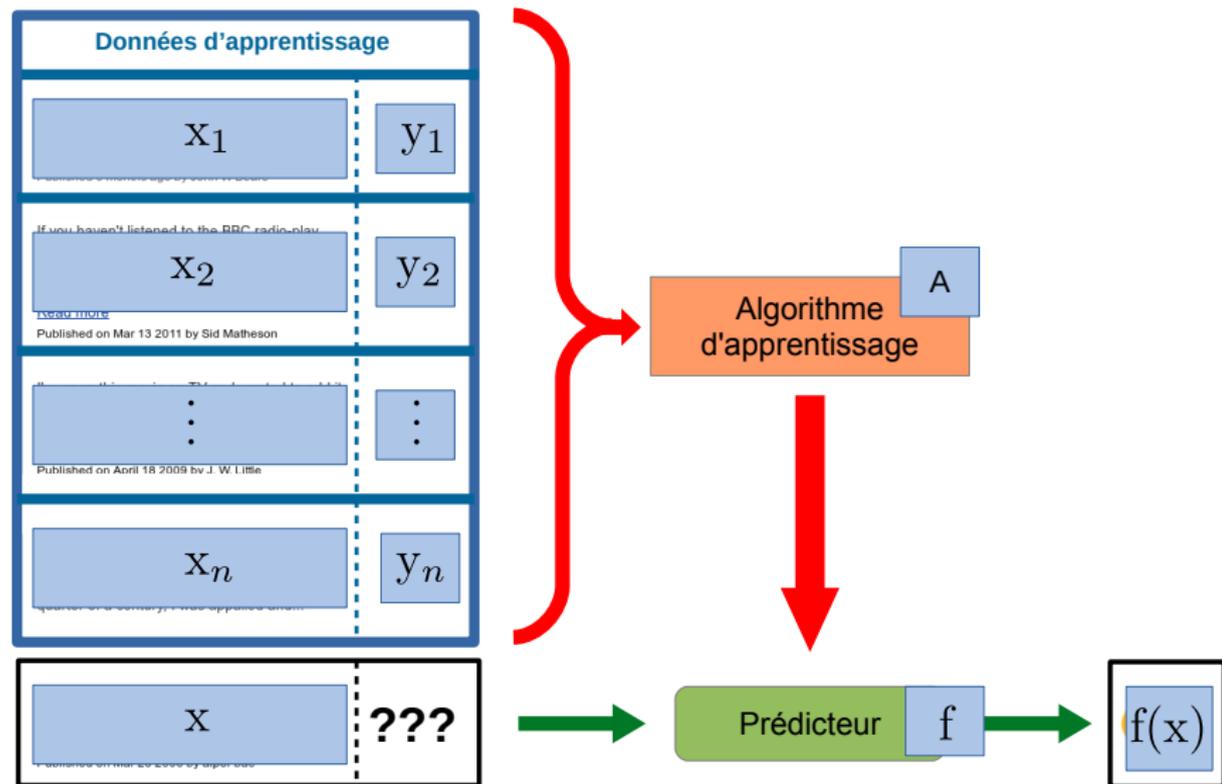
Prédicteur



Apprentissage automatique: Exemple



Abstraction mathématique



Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Ensemble d'apprentissage

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

Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Ensemble d'apprentissage

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

Prédicteur

$$f : \mathcal{X} \rightarrow \mathcal{Y}, \quad f \in \mathcal{F}$$

Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Ensemble d'apprentissage

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

Prédicteur

$$f : \mathcal{X} \rightarrow \mathcal{Y}, \quad f \in \mathcal{F}$$

Algorithme d'apprentissage

$$A(S) \rightarrow f$$

Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Ensemble d'apprentissage

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

Prédicteur

$$f : \mathcal{X} \rightarrow \mathcal{Y}, \quad f \in \mathcal{F}$$

Algorithme d'apprentissage

$$A(S) \rightarrow f$$

Fonction de perte

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

Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Ensemble d'apprentissage

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

Prédicteur

$$f : \mathcal{X} \rightarrow \mathcal{Y}, \quad f \in \mathcal{F}$$

Algorithme d'apprentissage

$$A(S) \rightarrow f$$

Fonction de perte

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

Exemple: perte zéro-un

$$\ell_{01}(f(x), y) = \begin{cases} 0 & \text{si } f(x) = y \\ 1 & \text{sinon.} \end{cases}$$

Définitions

Observation d'apprentissage

Une observation $(x, y) \in \mathcal{X} \times \mathcal{Y}$ est une paire **description-étiquette**.

Ensemble d'apprentissage

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

Prédicteur

$$f : \mathcal{X} \rightarrow \mathcal{Y}, \quad f \in \mathcal{F}$$

Algorithme d'apprentissage

$$A(S) \rightarrow f$$

Fonction de perte

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

Exemple: perte zéro-un

$$\ell_{01}(f(x), y) = \begin{cases} 0 & \text{si } f(x) = y \\ 1 & \text{sinon.} \end{cases}$$

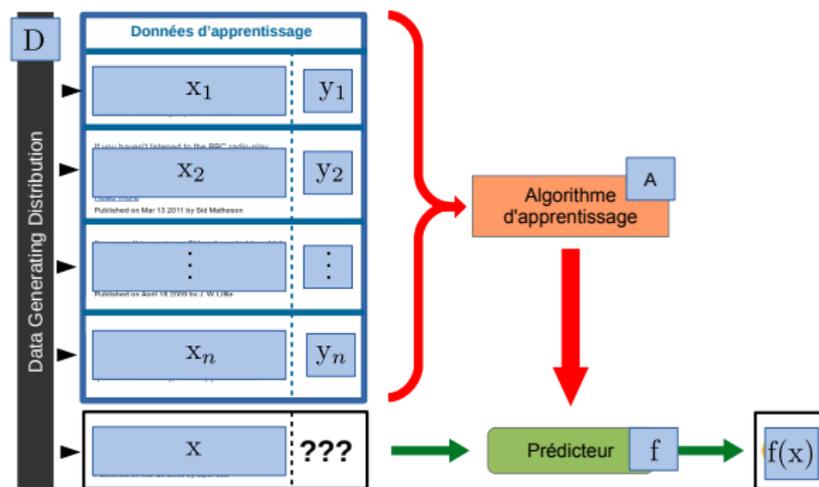
Perte empirique

$$\hat{\mathcal{L}}_S(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i)$$

Défi de généralisation

Hypothèse fréquentiste

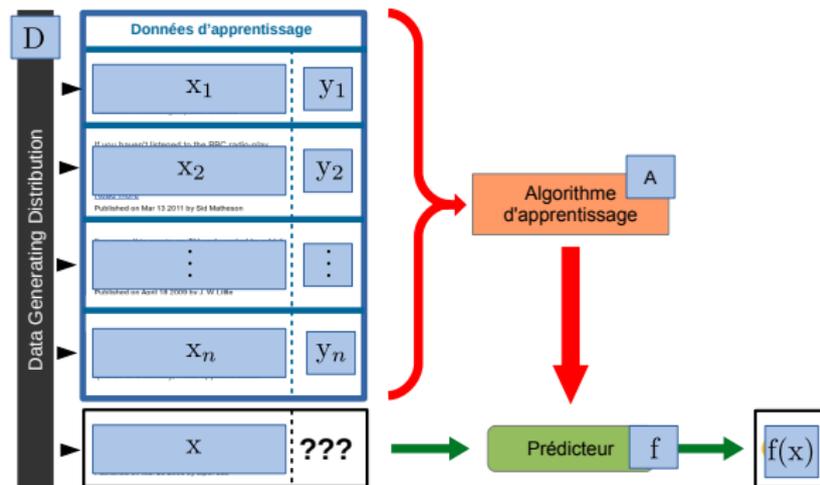
$$S = \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \} \sim D^n$$



Défi de généralisation

Hypothèse fréquentiste

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



Objectif

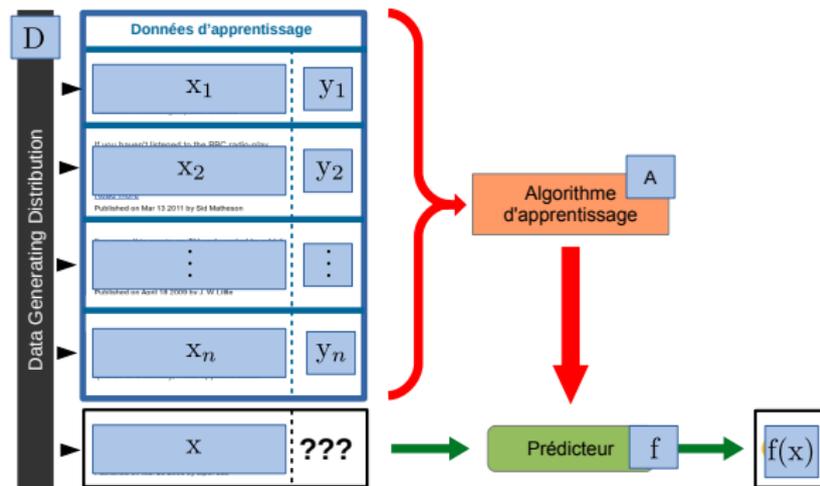
Minimiser la **perte en généralisation**

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

Défi de généralisation

Hypothèse fréquentiste

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



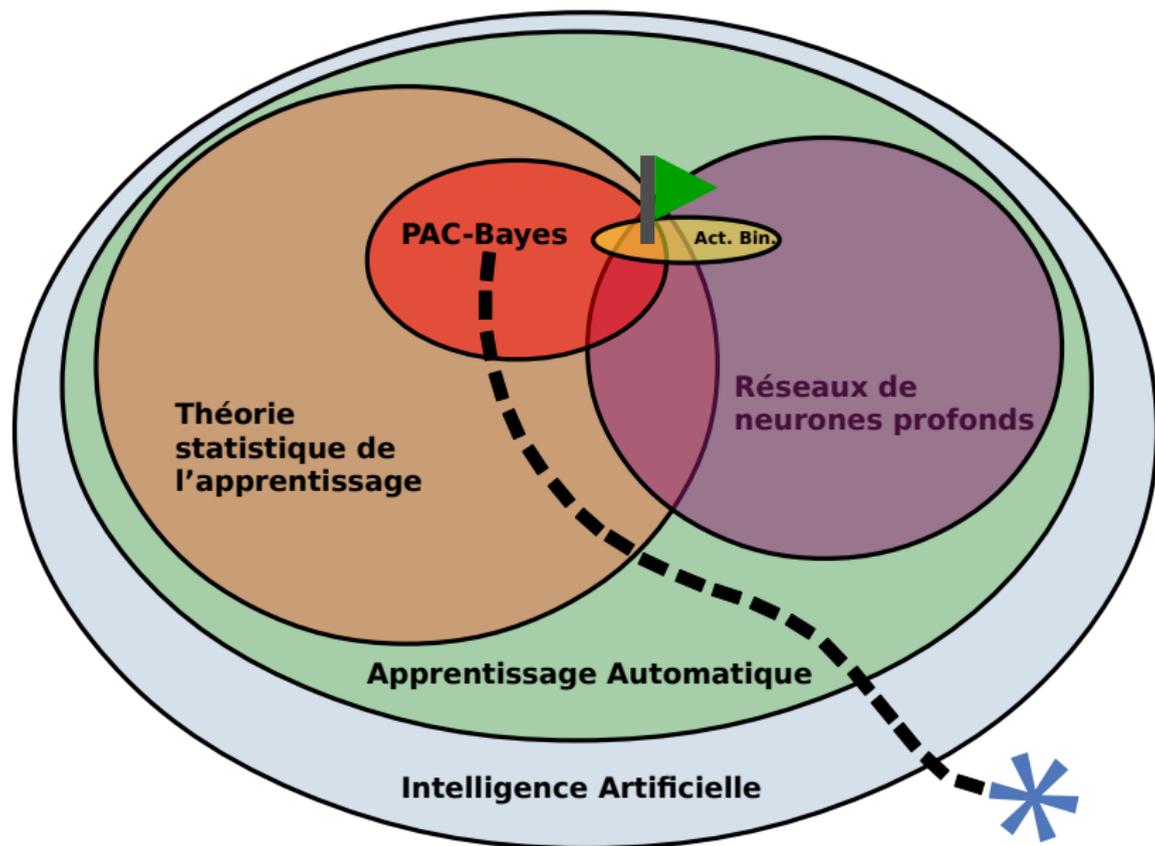
Objectif

Minimiser la **perte en généralisation**

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

On a accès *seulement* à la **perte empirique** sur S

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



Théorie PAC-Bayésienne (McAllester 1999)

La théorie PAC-Bayésienne est une théorie fréquentiste qui formule des garanties PAC sur des prédicteurs apparentés aux méthodes bayésiennes.

Théorie PAC-Bayésienne (McAllester 1999)

La théorie PAC-Bayésienne est une **théorie fréquentiste** qui formule des garanties PAC sur des prédicteurs apparentés aux méthodes bayésiennes.

Hypothèse fréquentiste

Chaque observation est générée *i.i.d.* par une **distribution de données D** .

Théorie PAC-Bayésienne (McAllester 1999)

La théorie PAC-Bayésienne est une théorie fréquentiste qui formule des **garanties PAC** sur des prédicteurs apparentés aux méthodes bayésiennes.

Hypothèse fréquentiste

Chaque observation est générée *i.i.d.* par une **distribution de données D** .

Bornes PAC (probablement approximativement correctes)

Avec probabilité $1-\delta$, la perte \mathcal{L}_D du prédicteur f est inférieure à ϵ :

$$\Pr\left(\mathcal{L}_D(f) \leq \epsilon(\hat{\mathcal{L}}_S(f), n, \delta, \dots)\right) \geq 1-\delta$$

Théorie PAC-Bayésienne (McAllester 1999)

La théorie PAC-Bayésienne est une théorie fréquentiste qui formule des **garanties PAC** sur des prédicteurs apparentés aux méthodes bayésiennes.

Hypothèse fréquentiste

Chaque observation est générée *i.i.d.* par une **distribution de données** D .

Bornes PAC (probablement approximativement correctes)

Avec probabilité $1-\delta$, la **perte** \mathcal{L}_D du **prédicteur** f est inférieure à ϵ :

$$\Pr \left(\mathcal{L}_D(f) \leq \epsilon(\hat{\mathcal{L}}_S(f), n, \delta, \dots) \right) \geq 1-\delta$$

Théorie PAC-Bayésienne (McAllester 1999)

La théorie PAC-Bayésienne est une théorie fréquentiste qui formule des **garanties PAC** sur des prédicteurs apparentés aux méthodes bayésiennes.

Hypothèse fréquentiste

Chaque observation est générée *i.i.d.* par une **distribution de données D** .

Bornes PAC (probablement approximativement correctes)

Avec probabilité $1-\delta$, la **perte \mathcal{L}_D** du **prédicteur f** est **inférieure à ϵ** :

$$\Pr\left(\mathcal{L}_D(f) \leq \epsilon(\hat{\mathcal{L}}_S(f), n, \delta, \dots)\right) \geq 1-\delta$$

Théorie PAC-Bayésienne (McAllester 1999)

La théorie PAC-Bayésienne est une théorie fréquentiste qui formule des garanties PAC sur des prédicteurs **apparentés aux méthodes bayésiennes**.

Hypothèse fréquentiste

Chaque observation est générée *i.i.d.* par une **distribution de données D** .

Bornes PAC (probablement approximativement correctes)

Avec probabilité $1-\delta$, la **perte \mathcal{L}_D** du **prédicteur f** est **inférieure à ϵ** :

$$\Pr \left(\mathcal{L}_D(f) \leq \epsilon(\hat{\mathcal{L}}_S(f), n, \delta, \dots) \right) \geq 1-\delta$$

Inspiration bayésienne

- Considère une **agrégation** de prédicteurs $\mathbf{E}_{f \sim Q} \mathcal{L}_D(f)$ (distribution Q)
- Incorpore des **connaissances a priori** sur le problème (distribution P)

Théorème PAC-Bayésien

Bornes PAC (probablement approximativement correctes)

Avec probabilité $1-\delta$, la perte \mathcal{L}_D du prédicteur f est inférieure à ϵ :

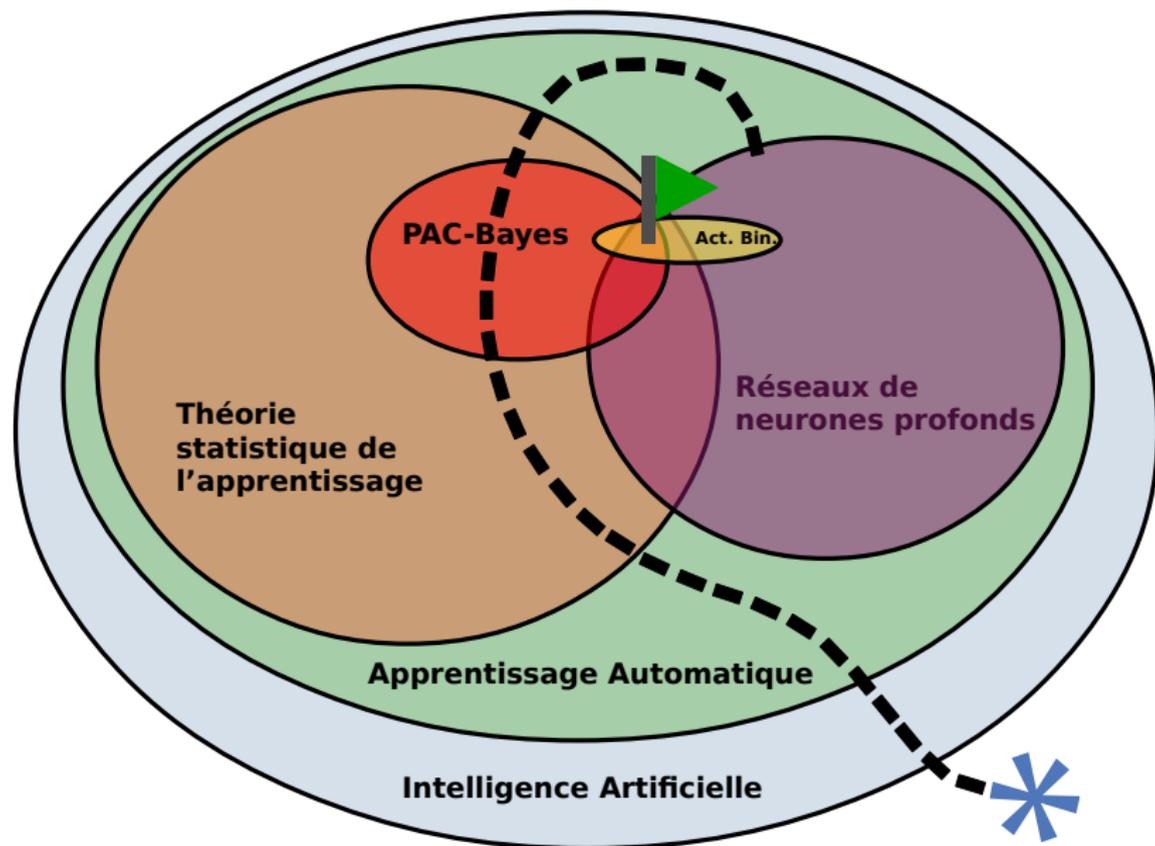
$$\Pr\left(\mathcal{L}_D(f) \leq \epsilon(\hat{\mathcal{L}}_S(f), n, \delta, \dots)\right) \geq 1-\delta$$

Borne PAC-Bayésienne

$$\Pr_{S \sim D^n} \left(\forall Q \text{ sur } \mathcal{F}, \forall C > 0 : \mathbf{E}_{f \sim Q} \mathcal{L}_D(f) \leq \epsilon \left(\mathbf{E}_{f \sim Q} \hat{\mathcal{L}}_S(f), n, \delta, P, Q, C \right) \right) \geq 1-\delta$$

$$\epsilon(\dots) = \frac{1}{1-e^{-C}} \left(1 - \exp \left(-C \mathbf{E}_{f \sim Q} \hat{\mathcal{L}}_S(f) - \frac{1}{n} \left[\text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right] \right) \right),$$

où $\text{KL}(Q \| P) = \mathbf{E}_{h \sim Q} \ln \frac{Q(h)}{P(h)}$ est la **divergence de Kullback-Leibler**.



Réseau de neurones

Problème de classification binaire:

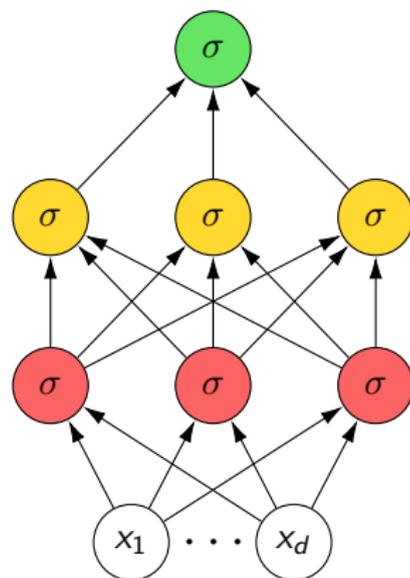
- $\mathbf{x} \in \mathbb{R}^d$
- $y \in \{-1, 1\}$

Architecture:

- L couches *pleinement connectées*
- La k -ième couche possède d_k neurones
- $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ est une *fonction d'activation*

Paramètres:

- $\mathbf{W}_k \in \mathbb{R}^{d_k \times d_{k-1}}$ sont des matrices de poids.
- $\theta = \text{vec}(\{\mathbf{W}_k\}_{k=1}^L) \in \mathbb{R}^D$



Réseau de neurones

Problème de classification binaire:

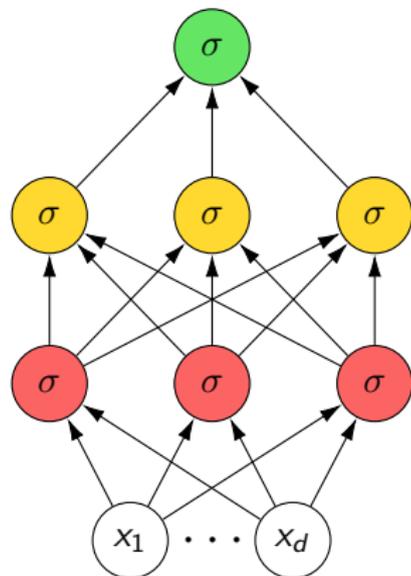
- $\mathbf{x} \in \mathbb{R}^d$
- $y \in \{-1, 1\}$

Architecture:

- L couches *pleinement connectées*
- La k -ième couche possède d_k neurones
- $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ est une *fonction d'activation*

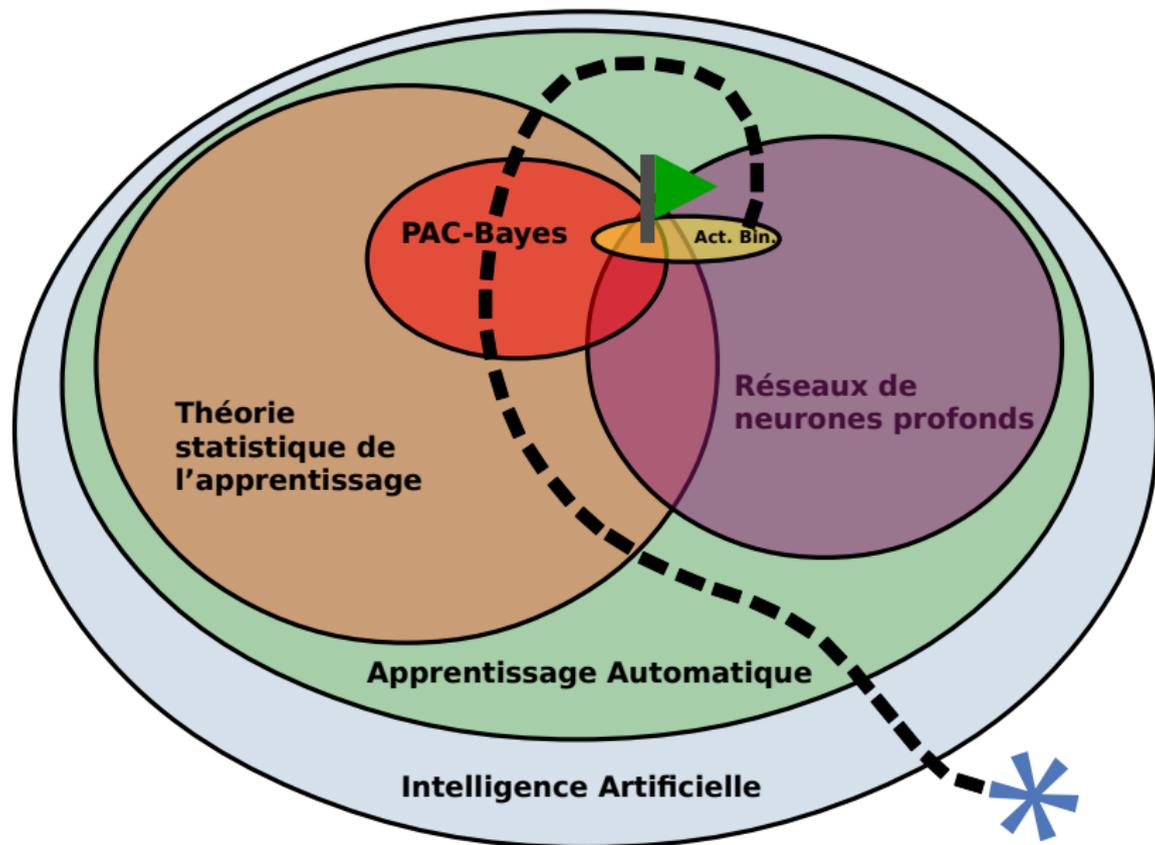
Paramètres:

- $\mathbf{W}_k \in \mathbb{R}^{d_k \times d_{k-1}}$ sont des matrices de poids.
- $\theta = \text{vec}(\{\mathbf{W}_k\}_{k=1}^L) \in \mathbb{R}^D$



Prédiction

$$f_{\theta}(\mathbf{x}) = \sigma(\mathbf{w}_L \sigma(\mathbf{W}_{L-1} \sigma(\dots \sigma(\mathbf{W}_1 \mathbf{x})))) .$$



Réseau de neurones à *activations binaires*

Problème de classification binaire:

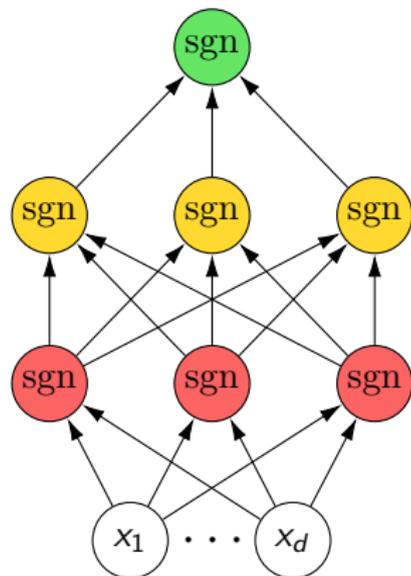
- $\mathbf{x} \in \mathbb{R}^d$
- $y \in \{-1, 1\}$

Architecture:

- L couches *pleinement connectées*
- La k -ième couche possède d_k neurones
- $\text{sgn}(a) = 1$ si $a > 0$ et $\text{sgn}(a) = -1$ sinon

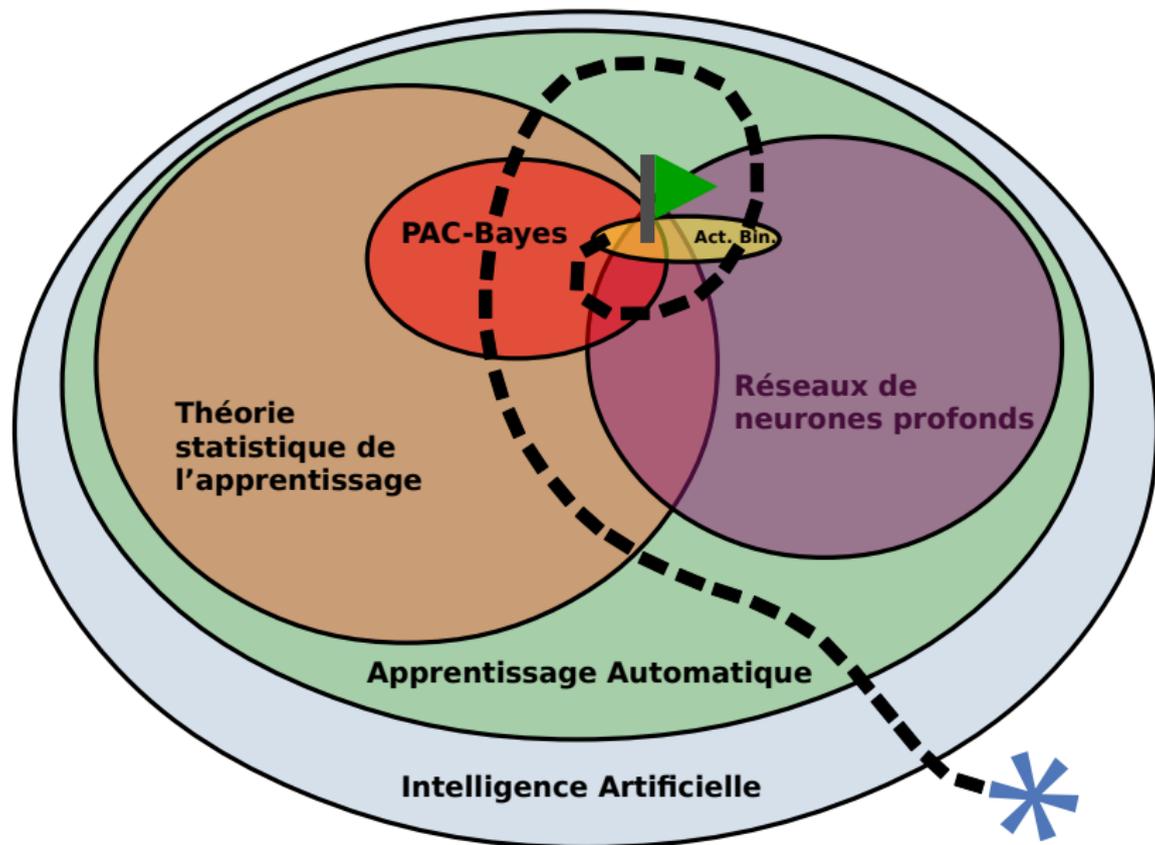
Paramètres:

- $\mathbf{W}_k \in \mathbb{R}^{d_k \times d_{k-1}}$ sont des matrices de poids.
- $\theta = \text{vec}(\{\mathbf{W}_k\}_{k=1}^L) \in \mathbb{R}^D$



Prédiction

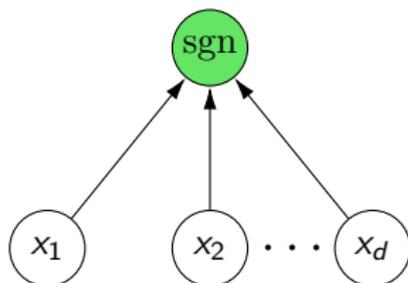
$$f_{\theta}(\mathbf{x}) = \text{sgn}(\mathbf{w}_L \text{sgn}(\mathbf{W}_{L-1} \text{sgn}(\dots \text{sgn}(\mathbf{W}_1 \mathbf{x})))) .$$



Une couche

“PAC-Bayesian Learning of Linear Classifiers” (Germain et al., 2009)

$$f_{\mathbf{w}}(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x}), \text{ avec } \mathbf{w} \in \mathbb{R}^d.$$



Une couche

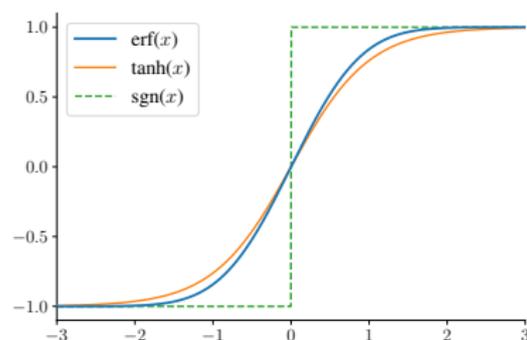
“PAC-Bayesian Learning of Linear Classifiers” (Germain et al., 2009)

$$f_{\mathbf{w}}(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x}), \text{ avec } \mathbf{w} \in \mathbb{R}^d.$$

Analyse PAC-Bayésienne:

- Espace des prédicteurs $\mathcal{F}_d = \{f_{\mathbf{v}} | \mathbf{v} \in \mathbb{R}^d\}$
- Prior gaussien $P_{\mathbf{w}_0} = \mathcal{N}(\mathbf{w}_0, I_d)$ sur \mathcal{F}_d
- Posterior gaussien $Q_{\mathbf{w}} = \mathcal{N}(\mathbf{w}, I_d)$ sur \mathcal{F}_d
- Prédicteur

$$F_{\mathbf{w}}(\mathbf{x}) = \mathbf{E}_{\mathbf{v} \sim Q_{\mathbf{w}}} f_{\mathbf{v}}(\mathbf{x}) = \text{erf}\left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{2}\|\mathbf{x}\|}\right)$$



Une couche

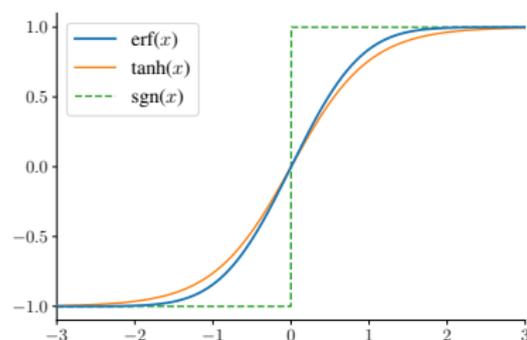
“PAC-Bayesian Learning of Linear Classifiers” (Germain et al., 2009)

$$f_{\mathbf{w}}(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x}), \text{ avec } \mathbf{w} \in \mathbb{R}^d.$$

Analyse PAC-Bayésienne:

- Espace des prédicteurs $\mathcal{F}_d = \{f_{\mathbf{v}} | \mathbf{v} \in \mathbb{R}^d\}$
- Prior gaussien $P_{\mathbf{w}_0} = \mathcal{N}(\mathbf{w}_0, I_d)$ sur \mathcal{F}_d
- Posterior gaussien $Q_{\mathbf{w}} = \mathcal{N}(\mathbf{w}, I_d)$ sur \mathcal{F}_d
- Prédicteur

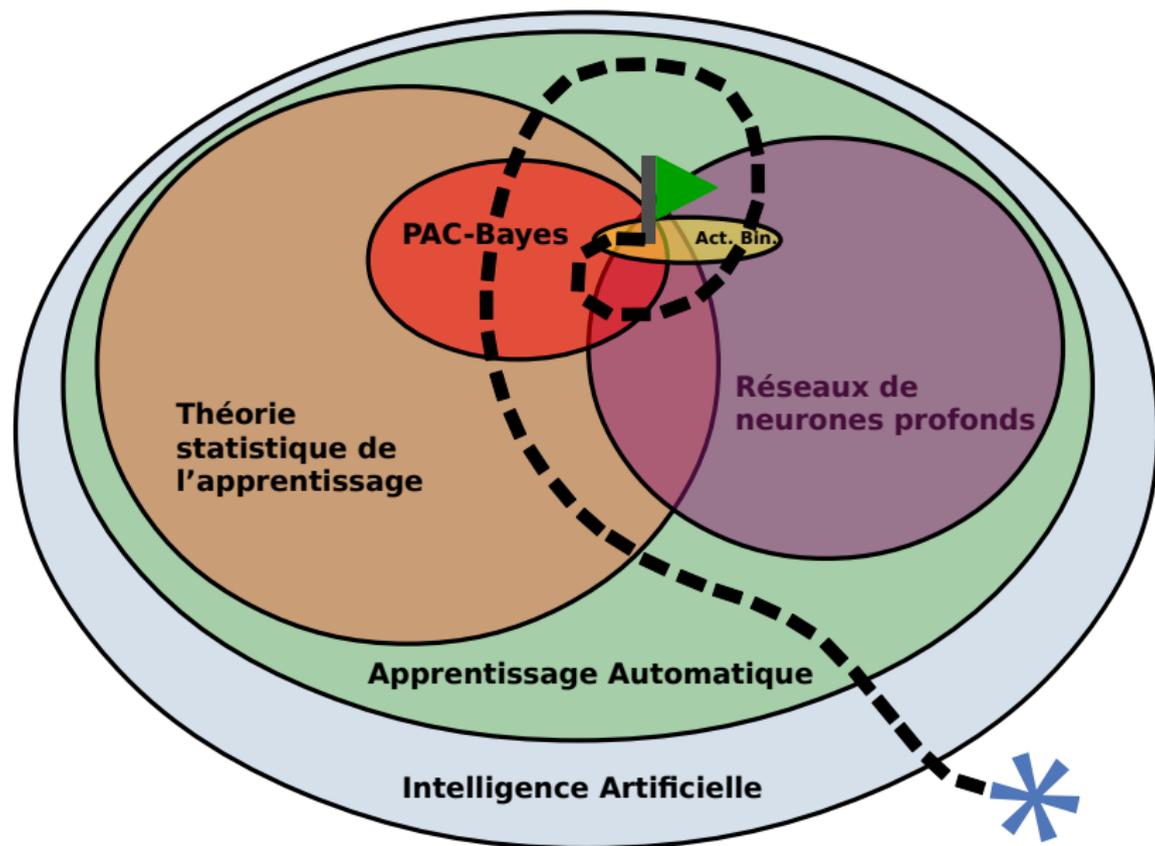
$$F_{\mathbf{w}}(\mathbf{x}) = \mathbf{E}_{\mathbf{v} \sim Q_{\mathbf{w}}} f_{\mathbf{v}}(\mathbf{x}) = \text{erf}\left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{2}\|\mathbf{x}\|}\right)$$



Minimisation de la borne

(perte linéaire $\ell(y, y') \stackrel{\text{def}}{=} \frac{1}{2}(1 - yy')$)

$$C n \hat{\mathcal{L}}_S(F_{\mathbf{w}}) + \text{KL}(Q_{\mathbf{w}} \| P_{\mathbf{w}_0}) = C \frac{1}{2} \sum_{i=1}^n \text{erf}\left(-y_i \frac{\mathbf{w} \cdot \mathbf{x}_i}{\sqrt{2}\|\mathbf{x}_i\|}\right) + \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|^2.$$



Deux couches

Posterior $Q_\theta = \mathcal{N}(\theta, I_D)$ sur l'ensemble des réseaux $\mathcal{F}_D = \{f_{\tilde{\theta}} \mid \tilde{\theta} \in \mathbb{R}^D\}$, avec

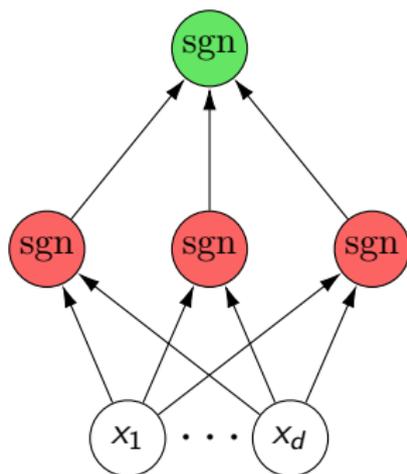
$$f_{\tilde{\theta}}(\mathbf{x}) = \text{sgn}(\mathbf{w}_2 \cdot \text{sgn}(\mathbf{W}_1 \mathbf{x})) .$$

$$F_\theta(\mathbf{x}) = \mathbf{E}_{\tilde{\theta} \sim Q_\theta} f_{\tilde{\theta}}(\mathbf{x})$$

$$= \int_{\mathbb{R}^{d_1 \times d_0}} Q_1(\mathbf{V}_1) \int_{\mathbb{R}^{d_1}} Q_2(\mathbf{v}_2) \text{sgn}(\mathbf{v}_2 \cdot \text{sgn}(\mathbf{V}_1 \mathbf{x})) d\mathbf{v}_2 d\mathbf{V}_1$$

$$= \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} \text{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right) \int_{\mathbb{R}^{d_1 \times d_0}} \mathbb{1}[\mathbf{s} = \text{sgn}(\mathbf{V}_1 \mathbf{x})] Q_1(\mathbf{V}_1) d\mathbf{V}_1$$

$$= \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} \underbrace{\text{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right)}_{F_{\mathbf{w}_2}(\mathbf{s})} \underbrace{\prod_{i=1}^{d_1} \left[\frac{1}{2} + \frac{s_i}{2} \text{erf}\left(\frac{\mathbf{w}_1^i \cdot \mathbf{x}}{\sqrt{2} \|\mathbf{x}\|}\right) \right]}_{\Pr(\mathbf{s}|\mathbf{x}, \mathbf{W}_1)} .$$



Ingrédients de la borne PAC-Bayes

Prédicteur

$$F_{\theta}(\mathbf{x}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} f_{\tilde{\theta}}(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} \operatorname{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right) \prod_{i=1}^{d_1} \left[\frac{1}{2} + \frac{s_i}{2} \operatorname{erf}\left(\frac{\mathbf{w}_1^i \cdot \mathbf{x}}{\sqrt{2} \|\mathbf{x}\|}\right) \right].$$

Ingrédients de la borne PAC-Bayes

Prédicteur

$$F_{\theta}(\mathbf{x}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} f_{\tilde{\theta}}(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} \operatorname{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right) \prod_{i=1}^{d_1} \left[\frac{1}{2} + \frac{s_i}{2} \operatorname{erf}\left(\frac{\mathbf{w}_1^i \cdot \mathbf{x}}{\sqrt{2} \|\mathbf{x}\|}\right) \right].$$

Perte empirique

$$\hat{\mathcal{L}}_S(F_{\theta}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} \hat{\mathcal{L}}_S(f_{\tilde{\theta}}) = \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} - \frac{1}{2} y_i F_{\theta}(\mathbf{x}_i) \right].$$

Ingrédients de la borne PAC-Bayes

Prédicteur

$$F_{\theta}(\mathbf{x}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} f_{\tilde{\theta}}(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} \operatorname{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right) \prod_{i=1}^{d_1} \left[\frac{1}{2} + \frac{s_i}{2} \operatorname{erf}\left(\frac{\mathbf{w}_1^i \cdot \mathbf{x}}{\sqrt{2} \|\mathbf{x}\|}\right) \right].$$

Perte empirique

$$\hat{\mathcal{L}}_S(F_{\theta}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} \hat{\mathcal{L}}_S(f_{\tilde{\theta}}) = \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} - \frac{1}{2} y_i F_{\theta}(\mathbf{x}_i) \right].$$

Terme de complexité

$$\operatorname{KL}(Q_{\theta} \| P_{\theta_0}) = \frac{1}{2} \|\theta - \theta_0\|^2.$$

Ingédients de la borne PAC-Bayes

Prédicteur

$$F_{\theta}(\mathbf{x}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} f_{\tilde{\theta}}(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} \operatorname{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right) \prod_{i=1}^{d_1} \left[\frac{1}{2} + \frac{s_i}{2} \operatorname{erf}\left(\frac{\mathbf{w}_1^i \cdot \mathbf{x}}{\sqrt{2} \|\mathbf{x}\|}\right) \right].$$

Perte empirique

$$\hat{\mathcal{L}}_S(F_{\theta}) = \mathbf{E}_{\tilde{\theta} \sim Q_{\theta}} \hat{\mathcal{L}}_S(f_{\tilde{\theta}}) = \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} - \frac{1}{2} y_i F_{\theta}(\mathbf{x}_i) \right].$$

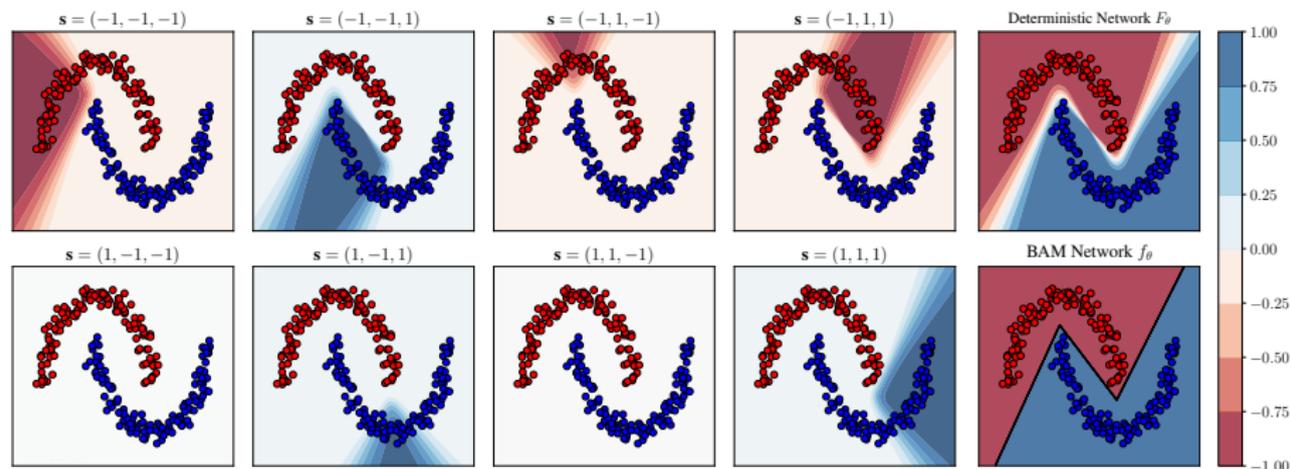
Terme de complexité

$$\operatorname{KL}(Q_{\theta} \| P_{\theta_0}) = \frac{1}{2} \|\theta - \theta_0\|^2.$$

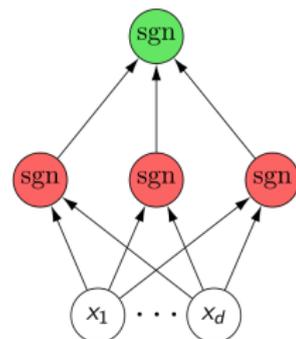
Garantie de généralisation

$$\frac{1}{1-e^{-c}} \left(1 - \exp \left(-c \hat{\mathcal{L}}_S(F_{\theta}) - \frac{1}{n} [\operatorname{KL}(Q_{\theta} \| P_{\mu}) + \ln \frac{2\sqrt{n}}{\delta}] \right) \right)$$

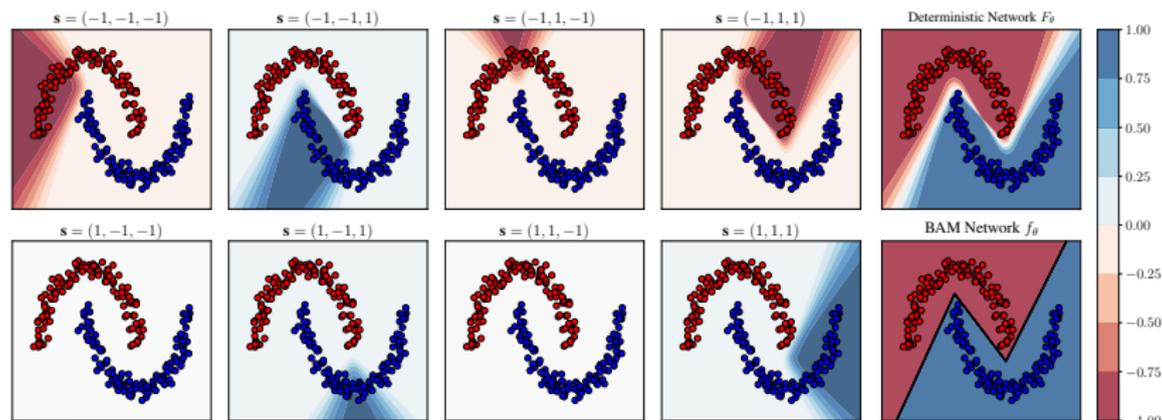
Visualisation



$$F_\theta(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1, 1\}^{d_1}} \underbrace{\operatorname{erf}\left(\frac{\mathbf{w}_2 \cdot \mathbf{s}}{\sqrt{2d_1}}\right)}_{F_{w_2}(\mathbf{s})} \underbrace{\prod_{i=1}^{d_1} \left[\frac{1}{2} + \frac{s_i}{2} \operatorname{erf}\left(\frac{\mathbf{w}_1^i \cdot \mathbf{x}}{\sqrt{2} \|\mathbf{x}\|}\right) \right]}_{\Pr(\mathbf{s}|\mathbf{x}, \mathbf{W}_1)}.$$

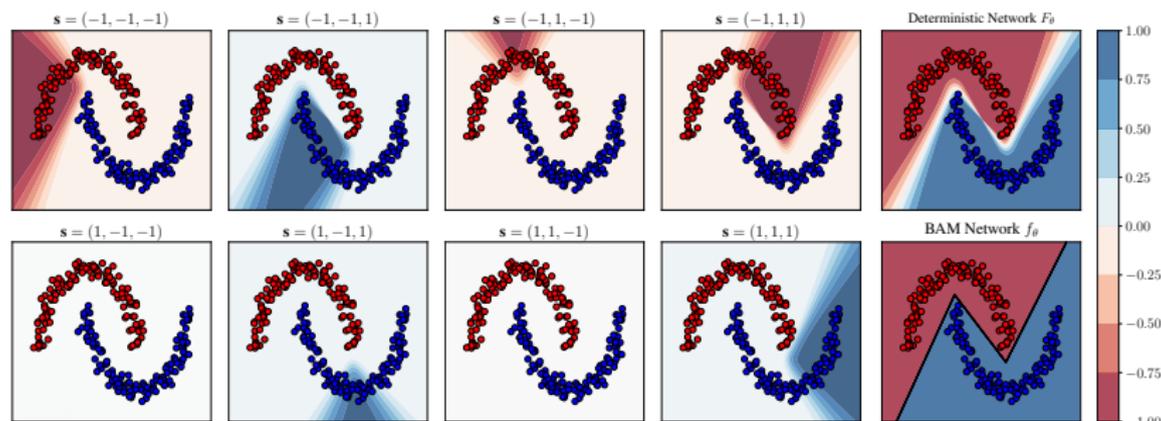


Approximation stochastique



$$F_\theta(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1, 1\}^{d_1}} F_{w_2}(\mathbf{s}) \Pr(\mathbf{s} | \mathbf{x}, \mathbf{W}_1)$$

Approximation stochastique



$$F_{\theta}(\mathbf{x}) = \sum_{\mathbf{s} \in \{-1,1\}^{d_1}} F_{\mathbf{w}_2}(\mathbf{s}) \Pr(\mathbf{s}|\mathbf{x}, \mathbf{W}_1)$$

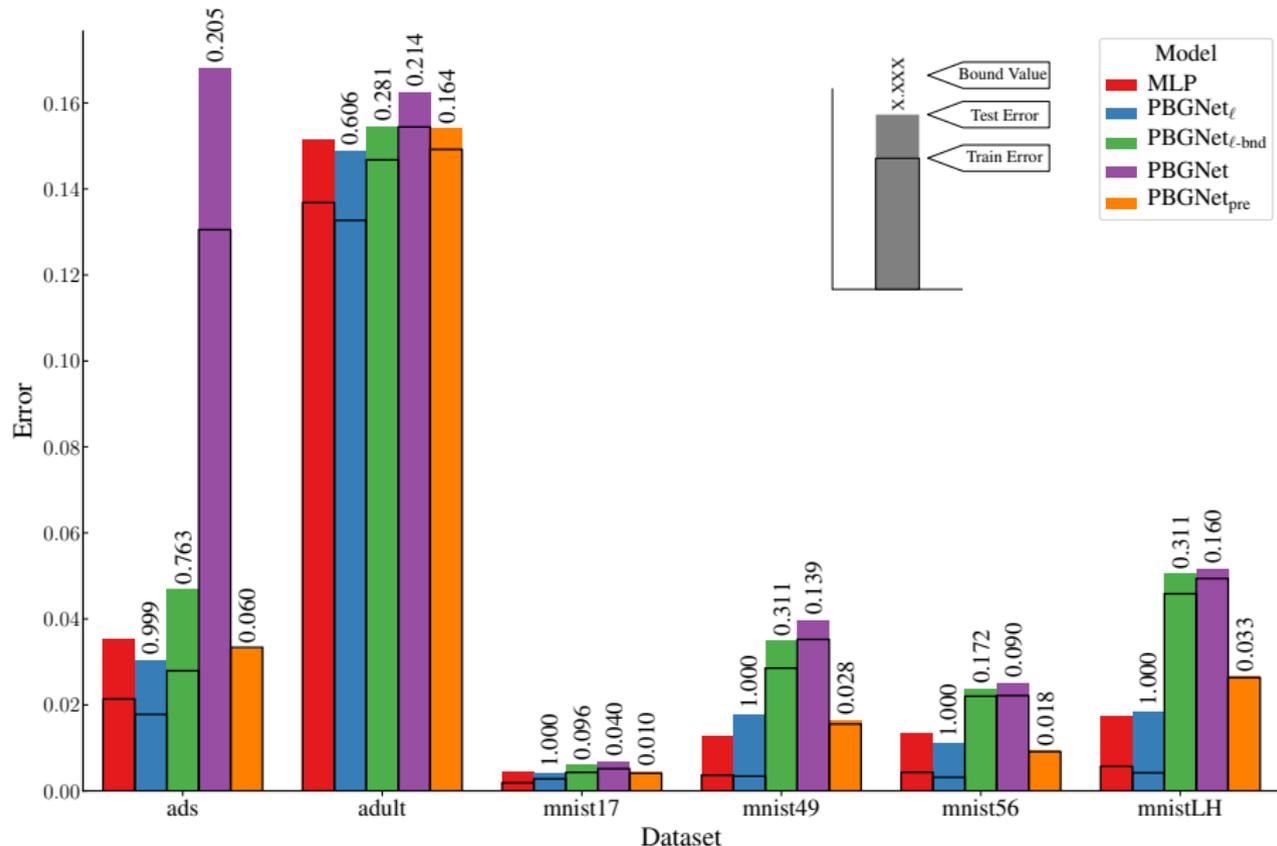
Échantillonnage de Monte Carlo

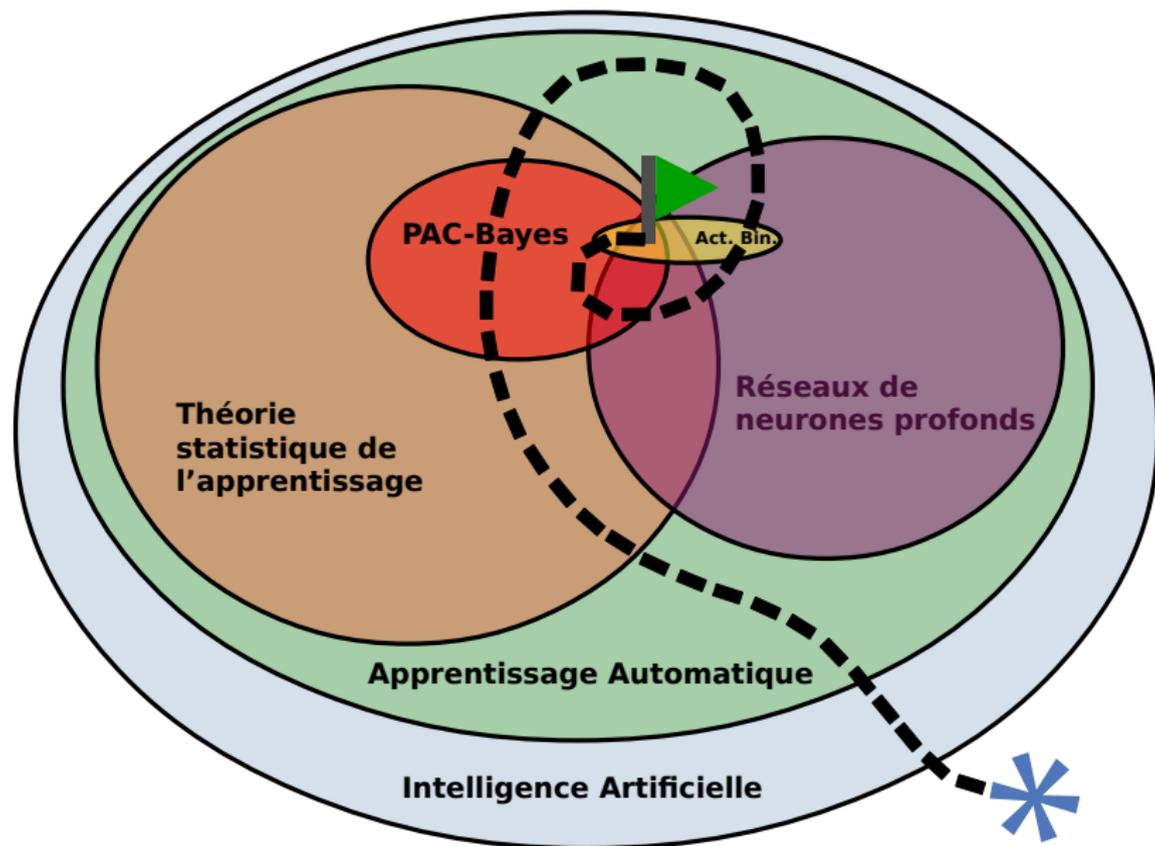
On génère T vecteurs binaires aléatoires $\{\mathbf{s}^t\}_{t=1}^T$ selon $\Pr(\mathbf{s}|\mathbf{x}, \mathbf{W}_1)$

Prédiction.

$$F_{\theta}(\mathbf{x}) \approx \frac{1}{T} \sum_{t=1}^T F_{\mathbf{w}_2}(\mathbf{s}^t).$$

Expérimentation: Résultats





Étendre la méthode:

- Multiclasse
- Réseaux à convolutions
- Apprentissage par transfert

Analyser les propriétés:

- Robustesse
- Interprétabilité

Étendre la méthode:

- Multiclasse
- Réseaux à convolutions
- Apprentissage par transfert

Analyser les propriétés:

- Robustesse
- Interprétabilité

Opportunités de stages!