

Apprentissage Machine (supervisé)

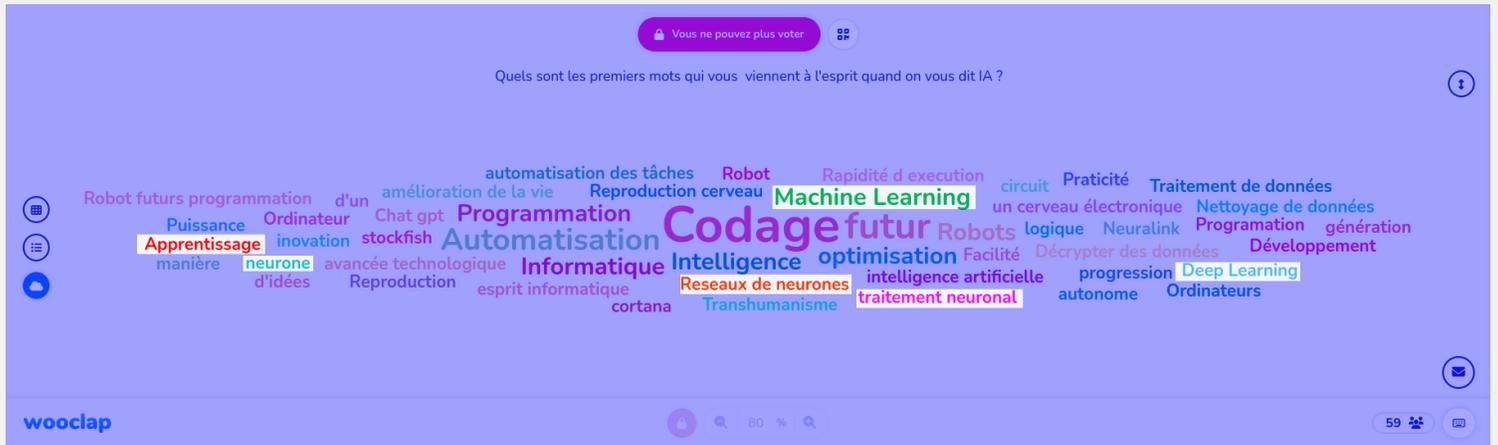
Florent Capelli
24 Février 2024

1

Introduction

2

Qu'est-ce que c'est l'IA pour vous ?



Vos réponses au Wooclap !

De la spécification au programme

Écrire un programme revient souvent à :

- Identifier le problème
 - Entrée ?
 - Sortie ?
- Spécification :
- Produire un programme respectant la spécification

De la spécification au programme

Écrire un programme revient souvent à :

- Identifier le problème
 - Entrée ? **Un tableau T**
 - Sortie ?
- Spécification :

- Produire un programme respectant la spécification

De la spécification au programme

Écrire un programme revient souvent à :

- Identifier le problème
 - Entrée ? **Un tableau T**
 - Sortie ? **Un tableau U représentant T trié**
- Spécification :

- Produire un programme respectant la spécification

De la spécification au programme

Écrire un programme revient souvent à :

- Identifier le problème
 - Entrée ? **Un tableau T**
 - Sortie ? **Un tableau U représentant T trié**
- Spécification :
 - $U[i] \leq U[i+1]$
 - **Il existe bijection $\pi : [n] \rightarrow [n], T[i] = U[\pi[i]]$**
 - Produire un programme respectant la spécification

4.3

De la spécification au programme

Écrire un programme revient souvent à :

- Identifier le problème
 - Entrée ? **Un tableau T**
 - Sortie ? **Un tableau U représentant T trié**
- Spécification :
 - $U[i] \leq U[i+1]$
 - **Il existe bijection $\pi : [n] \rightarrow [n], T[i] = U[\pi[i]]$**
 - Produire un programme respectant la spécification

```
def selection_sort(t):
    u=t.copy()
    for i in range(len(u)):
        m=i
        for j in range(i,len(u)):
            if u[j] < u[m]:
                m=j
        u[m], u[i] = u[i], u[m]
    return u
```

4.4

Quelle spécification ?

Quelle spécification ?

- Une banque a-t-elle **intérêt** à accorder un crédit à X ?

FIRST BANK OF WIKI
1425 JAMES ST. PO BOX 4000
VICTORIA BC V8X 3X4 1-800-555-5555

CHEQUING ACCOUNT STATEMENT
Page : 1 of 1

JOHN JONES
1643 DUNDAS ST W APT 27
TORONTO ON M6K 1V2

Date	Description	Ref.	Withdrawals	Deposits	Balance
2003-10-08	Previous balance				0.55
2003-10-14	Payroll Deposit - HOTEL			694.81	695.36
2003-10-14	Web Bill Payment - MASTERCARD	9685	200.00		495.36
2003-10-16	ATM Withdrawal - INTERAC	3990	21.25		474.11
2003-10-16	Fees - Interac		1.50		472.61
2003-10-20	Interac Purchase - ELECTRONICS	1975	2.99		469.62
2003-10-21	Web Bill Payment - AMEX	3314	300.00		169.62
2003-10-22	ATM Withdrawal - FIRST BANK	0064	100.00		69.62
2003-10-23	Interac Purchase - SUPERMARKET	1559	29.08		40.54
2003-10-24	Interac Refund - ELECTRONICS	1975		2.99	43.53
2003-10-27	Telephone Bill Payment - VISA	2475	6.77		36.76
2003-10-28	Payroll Deposit - HOTEL			694.81	731.57
2003-10-30	Web Funds Transfer - From SAVINGS	2620		50.00	781.57
2003-11-03	Pre-Auth. Payment - INSURANCE		33.55		748.02
2003-11-03	Cheque No. - 409		100.00		648.02
2003-11-06	Mortgage Payment		710.49		-62.47
2003-11-07	Fees - Overdraft		5.00		-67.47
2003-11-08	Fees - Monthly		5.00		-72.47
*** Totals ***			1,515.63	1,442.61	

Étant donné un dossier de client, quel est la probabilité que le client ne rembourse pas son crédit ?

Quelle spécification ?

- Une banque a-t-elle **intérêt** à accorder un crédit à X ?
- Reconnaître un chat sur des photos.



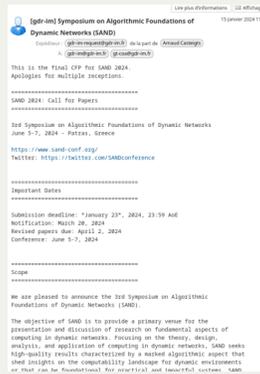
Étant donnée une image, y a-t-il un chat sur cette image ?

Quelle spécification ?

- Une banque a-t-elle **intérêt** à accorder un crédit à X ?
- Reconnaître un chat sur des photos.
- Détecter des spams



Spam!



Spam!

Étant donné un email, dois-je l'envoyer dans le dossier spam ?

Quelle spécification ?

- Une banque a-t-elle **intêret** à accorder un crédit à X ?
- Reconnaître un chat sur des photos.
- Détecter des spams
- Recommander des films qu'un utilisateur aimerait



Algorithmes de recommandation

Connaissant les films regardés précédemment par l'utilisateur, trouver un film qu'il va aimer.

Quelle spécification ?

- Une banque a-t-elle **intérêt** à accorder un crédit à X ?
- Reconnaître un chat sur des photos.
- Détecter des spams
- Recommander des films qu'un utilisateur aimerait

On ne sait pas toujours expliquer ou évaluer ce qu'on veut en sortie. On veut quelque chose "qui marche".

Expliquer par l'exemple

On veut **quelque chose** “qui marche” :

- On a des exemples : les **données**
 - “Bob 35 ans cadre dynamique n’a remboursé que 35% de son crédit au terme”.
 - “Cette image contient un chat et celle-ci non”.
 - “Ce message est un spam”.
 - “Alice a regardé et aimé *Harry Potter* et *Better Call Saul*”.

6

Expliquer par l'exemple

On veut **quelque chose** “qui marche” :

- On a des exemples : les **données**
 - “Bob 35 ans cadre dynamique n’a remboursé que 35% de son crédit au terme”.
 - “Cette image contient un chat et celle-ci non”.
 - “Ce message est un spam”.
 - “Alice a regardé et aimé *Harry Potter* et *Better Call Saul*”.
- On veut créer un programme qui **infère** des règles depuis les exemples :
 - doit marcher sur les exemples (la plupart au moins).
 - doit **généraliser** sur des données non connues.

6.1

Classification

Approche particulièrement adaptée à deux types de programme :

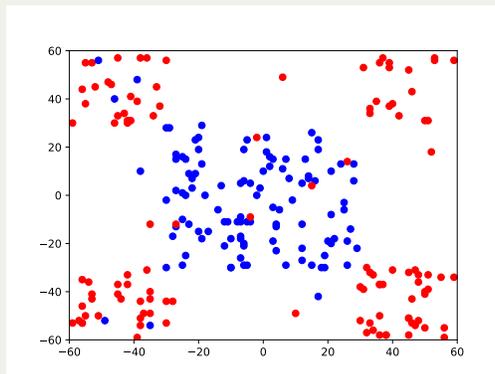
- **Classification** : on veut classer des objets \mathcal{O} dans des classes \mathcal{C}
 - Spam/non spam,
 - Chat/pas chat.

7

Classification

Approche particulièrement adaptée à deux types de programme :

- **Classification** : on veut classer des objets \mathcal{O} dans des classes \mathcal{C}
 - Spam/non spam,
 - Chat/pas chat.

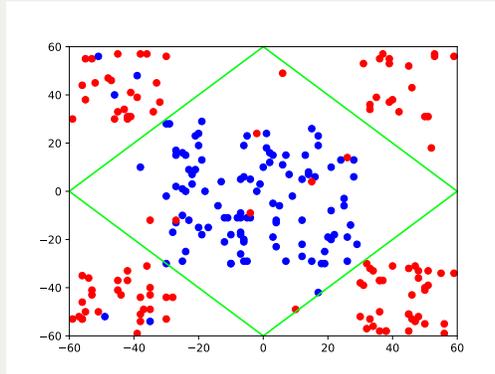


7.1

Classification

Approche particulièrement adaptée à deux types de programme :

- **Classification** : on veut classer des objets \mathcal{O} dans des classes \mathcal{C}
 - Spam/non spam,
 - Chat/pas chat.

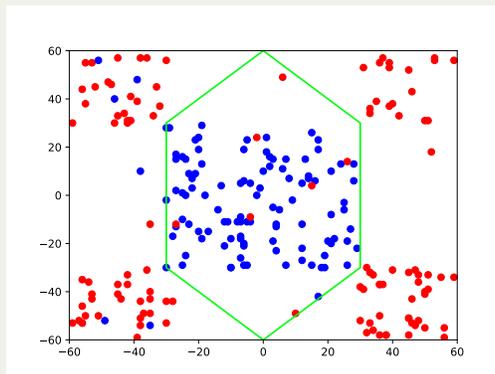


7.2

Classification

Approche particulièrement adaptée à deux types de programme :

- **Classification** : on veut classer des objets \mathcal{O} dans des classes \mathcal{C}
 - Spam/non spam,
 - Chat/pas chat.

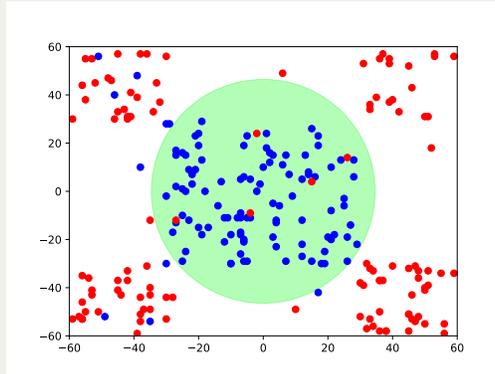


7.3

Classification

Approche particulièrement adaptée à deux types de programme :

- **Classification** : on veut classer des objets \mathcal{O} dans des classes \mathcal{C}
 - Spam/non spam,
 - Chat/pas chat.



7.4

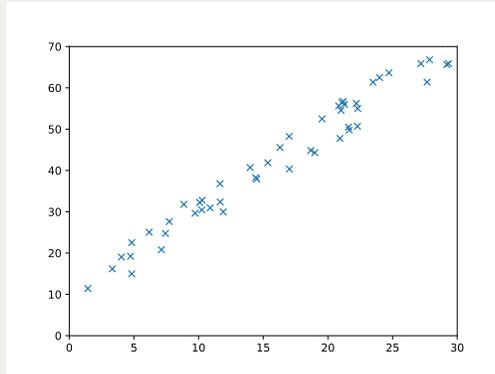
Régression

- **Régression** : on veut estimer une quantité (ou plusieurs) $f(o)$ pour $o \in \mathcal{O}$
 - Pourcentage du crédit remboursé.
 - Note donnée à film par l'utilisateur.

8

Régression

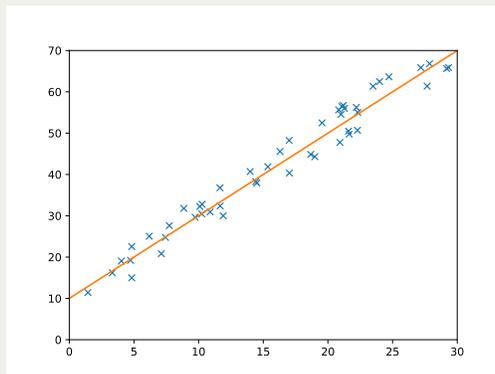
- **Régression** : on veut estimer une quantité (ou plusieurs) $f(o)$ pour $o \in \mathcal{O}$
 - Pourcentage du crédit remboursé.
 - Note donnée à film par l'utilisateur.



8.1

Régression

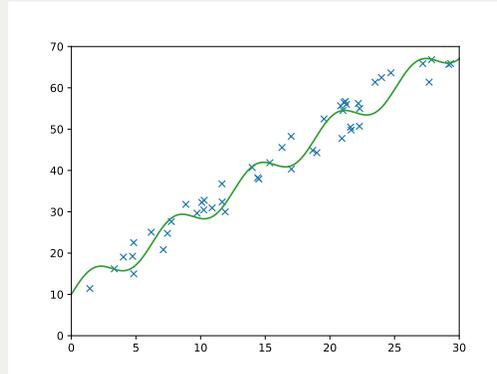
- **Régression** : on veut estimer une quantité (ou plusieurs) $f(o)$ pour $o \in \mathcal{O}$
 - Pourcentage du crédit remboursé.
 - Note donnée à film par l'utilisateur.



8.2

Régression

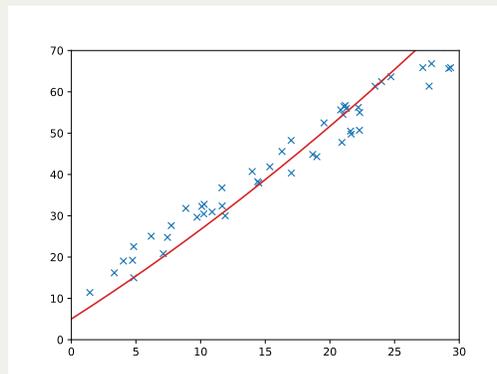
- **Régression** : on veut estimer une quantité (ou plusieurs) $f(o)$ pour $o \in \mathcal{O}$
 - Pourcentage du crédit remboursé.
 - Note donnée à film par l'utilisateur.



8.3

Régression

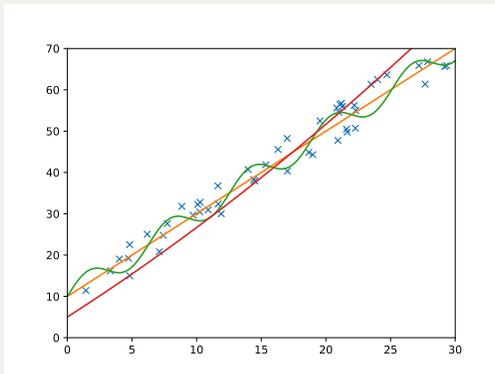
- **Régression** : on veut estimer une quantité (ou plusieurs) $f(o)$ pour $o \in \mathcal{O}$
 - Pourcentage du crédit remboursé.
 - Note donnée à film par l'utilisateur.



8.4

Régression

- **Régression** : on veut estimer une quantité (ou plusieurs) $f(o)$ pour $o \in \mathcal{O}$
 - Pourcentage du crédit remboursé.
 - Note donnée à film par l'utilisateur.



8.5

Classification comme régression

On peut voir la classification comme une régression :

- $f: \mathcal{O} \rightarrow \mathcal{C}$ vue comme $f_c: \mathcal{O} \rightarrow [0,1]$ où $f_c(o)$ est la “probabilité que o soit classé dans c ”
- **Exemple:**
 - $f_{spam}(\text{“Un héritage en attente”}) = 0.9, f_{message}(\text{“Un héritage en attente”}) = 0.1$
 - $f_{spam}(\text{“Invitation à la conférence...”}) = 0.6, f_{message}(\text{“Invitation à la conférence...”}) = 0.4$

9

Cadre théorique

10

Apprentissage supervisé

- On apprend depuis des données **étiquetées**
- On se limite aujourd'hui aux cas : régression et classification

D'autres types d'apprentissage :



"Étant donnée les pixels de l'écran, jouer à Pong"

Apprentissage par renforcement



Génère une image depuis du texte

Apprentissage non
supervisé

11

Cadre théorique idéal (classification)

- On a des objets \mathcal{O}
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{C}$ entre les objets et des *classes* qu'on **ne connaît pas**
- On connaît $f(o_1), \dots, f(o_n)$ pour certains objets $o_i \in \mathcal{O}$

Cadre théorique idéal (classification)

- On a des objets \mathcal{O}
 - **\mathcal{O} est l'ensemble de tous les emails possibles.**
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{C}$ entre les objets et des *classes* qu'on **ne connaît pas**
- On connaît $f(o_1), \dots, f(o_n)$ pour certains objets $o_i \in \mathcal{O}$

Cadre théorique idéal (classification)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les emails possibles.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{C}$ entre les objets et des *classes* qu'on **ne connaît pas**
 - $\mathcal{C} = \{spam, message\}$
 - $f(o) = spam$ seulement si je considère que o est un spam.
- On connaît $f(o_1), \dots, f(o_n)$ pour certains objets $o_i \in \mathcal{O}$

12.2

Cadre théorique idéal (classification)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les emails possibles.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{C}$ entre les objets et des *classes* qu'on **ne connaît pas**
 - $\mathcal{C} = \{spam, message\}$
 - $f(o) = spam$ seulement si je considère que o est un spam.
- On connaît $f(o_1), \dots, f(o_n)$ pour certains objets $o_i \in \mathcal{O}$
 - $f(\text{"Devenir riche en 5 minutes..."}) = spam$
 - $f(\text{"Journées du laboratoire 2024..."}) = message$

12.3

Cadre théorique idéal (classification)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les emails possibles.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{C}$ entre les objets et des *classes* qu'on **ne connaît pas**
 - $\mathcal{C} = \{\textit{spam}, \textit{message}\}$
 - $f(o) = \textit{spam}$ seulement si je considère que o est un spam.
- On connaît $f(o_1), \dots, f(o_n)$ pour certains objets $o_i \in \mathcal{O}$
 - $f(\textit{"Devenir riche en 5 minutes..."}) = \textit{spam}$
 - $f(\textit{"Journées du laboratoire 2024..."}) = \textit{message}$

Objectif: Trouver f^* telle que f^* classe *presque* comme f .

12.4

Cadre théorique idéal (régression)

- On a des objets \mathcal{O}
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{U}$ qu'on **ne connaît pas**
- On connaît $(o_1, f(o_1)), \dots, (o_n, f(o_n))$ pour certains objets $o_i \in \mathcal{O}$

13

Cadre théorique idéal (régression)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les clients possibles d'une banque.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{U}$ qu'on ne connaît pas

- On connaît $(o_1, f(o_1)), \dots, (o_n, f(o_n))$ pour certains objets $o_i \in \mathcal{O}$

13.1

Cadre théorique idéal (régression)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les clients possibles d'une banque.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{U}$ qu'on ne connaît pas
 - $\mathcal{U} = [0,100]$
 - le client o remboursera $f(o)\%$ de l'argent emprunté
- On connaît $(o_1, f(o_1)), \dots, (o_n, f(o_n))$ pour certains objets $o_i \in \mathcal{O}$

13.2

Cadre théorique idéal (régression)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les clients possibles d'une banque.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{U}$ qu'on ne connaît pas
 - $\mathcal{U} = [0,100]$
 - le client o remboursera $f(o)\%$ de l'argent emprunté
- On connaît $(o_1, f(o_1)), \dots, (o_n, f(o_n))$ pour certains objets $o_i \in \mathcal{O}$
 - $f(\text{Alice}) = 100$
 - $f(\text{Bob}) = 35$

13.3

Cadre théorique idéal (régression)

- On a des objets \mathcal{O}
 - \mathcal{O} est l'ensemble de tous les clients possibles d'une banque.
- On suppose qu'il existe $f: \mathcal{O} \rightarrow \mathcal{U}$ qu'on ne connaît pas
 - $\mathcal{U} = [0,100]$
 - le client o remboursera $f(o)\%$ de l'argent emprunté
- On connaît $(o_1, f(o_1)), \dots, (o_n, f(o_n))$ pour certains objets $o_i \in \mathcal{O}$
 - $f(\text{Alice}) = 100$
 - $f(\text{Bob}) = 35$

Objectif : Trouver f^* telle que $f^* \simeq f$.

13.4

Objets et représentations

- Est-ce que le patient a une angine ?
 - tension
 - fièvre
 - couleur de la gorge
 - etc.

Réponse : OUI/NON basée sur ces critères seulement.

Pour un objet $o \in \mathcal{O}$, on ne connaît qu'une représentation de o .

14

Importance de la représentation

- On suppose l'existence d'une fonction $f: \mathcal{O} \rightarrow \mathcal{U}$ sur un ensemble d'objet \mathcal{O} qu'on ne peut pas représenter
 - **Un humain, avec toute sa connaissance, situation personnelle, sociale, économique ...**
- On représente $o \in \mathcal{O}$ à l'aide de certains attributs
 - **Âge, Métier, Revenus ...**

Âge	Métier	Revenus	Crédit initial	Taux	Année	Part remboursée
32	Agent immobilier	45k€	100k€	2.5%	2004	60%
45	Enseignant	30k€	150k€	1.5%	2015	100%
...						
32	Agent immobilier	45k€	100k€	2.5%	2004	90%

15

Importance de la représentation

- On suppose l'existence d'une fonction $f: \mathcal{O} \rightarrow \mathcal{U}$ sur un ensemble d'objet \mathcal{O} qu'on ne peut pas représenter
 - **Un humain, avec toute sa connaissance, situation personnelle, sociale, économique ...**
- On représente $o \in \mathcal{O}$ à l'aide de certains attributs
 - **Âge, Métier, Revenus ...**

Âge	Métier	Revenus	Crédit initial	Taux	Année	Part remboursée
32	Agent immobilier	45k€	100k€	2.5%	2004	60%
45	Enseignant	30k€	150k€	1.5%	2015	100%
...						
32	Agent immobilier	45k€	100k€	2.5%	2004	90%

15.1

Importance de la représentation

- On suppose l'existence d'une fonction $f: \mathcal{O} \rightarrow \mathcal{U}$ sur un ensemble d'objet \mathcal{O} qu'on ne peut pas représenter
 - **Un humain, avec toute sa connaissance, situation personnelle, sociale, économique ...**
- On représente $o \in \mathcal{O}$ à l'aide de certains attributs
 - **Âge, Métier, Revenus ...**

Âge	Métier	Revenus	Crédit initial	Taux	Année	Part remboursée
32	Agent immobilier	45k€	100k€	2.5%	2004	60%
45	Enseignant	30k€	150k€	1.5%	2015	100%
...						
32	Agent immobilier	45k€	100k€	2.5%	2004	90%

Possiblement pas pertinent pour "apprendre" f .
Quid du problème si $r(o)$ =(**"couleur des yeux", "pointure de chaussure", "équipe de foot préférée"**)

15.2

Meilleur cadre théorique

1. Il existe une fonction **idéale**

- $f: \mathcal{O} \rightarrow \mathcal{C}$.

2. On a une **représentation** de \mathcal{O} via des **attributs** $\mathcal{R} = A_1, \dots, A_k$:

- $r: \mathcal{O} \rightarrow \mathcal{R}$

3. On a des **données**

- $r(o_1):f(o_1), \dots, r(o_n):f(o_n)$

4. *Hypothèses:*

- Les données sont des observations d'une variable aléatoire sur $\mathcal{R} \times \mathcal{C}$.
- Quelle distribution ? On suppose une distribution \mathcal{D} sur \mathcal{O} (proba de "rencontrer" $o \in \mathcal{O}$)
- Intuitivement : on observe (x,c) avec proba $Pr_o \simeq \mathcal{G}(x=r(o) \wedge f(o)=c)$

Meilleur cadre théorique

1. Il existe une fonction **idéale**

- $f: \mathcal{O} \rightarrow \mathcal{C}$.

pas forcément toujours vrai ; même ici on peut être incertain

2. On a une **représentation** de \mathcal{O} via des **attributs** $\mathcal{R} = A_1, \dots, A_k$:

- $r: \mathcal{O} \rightarrow \mathcal{R}$

3. On a des **données**

- $r(o_1):f(o_1), \dots, r(o_n):f(o_n)$

4. *Hypothèses:*

- Les données sont des observations d'une variable aléatoire sur $\mathcal{R} \times \mathcal{C}$.
- Quelle distribution ? On suppose une distribution \mathcal{D} sur \mathcal{O} (proba de "rencontrer" $o \in \mathcal{O}$)
- Intuitivement : on observe (x,c) avec proba $Pr_o \simeq \mathcal{G}(x=r(o) \wedge f(o)=c)$

Meilleur cadre théorique

1. Il existe une fonction **idéale**
 - $f: \mathcal{O} \rightarrow \mathcal{C}$.
 - pas forcément toujours vrai ; même ici on peut être incertain**
2. On a une **représentation** de \mathcal{O} via des **attributs** $\mathcal{R} = A_1, \dots, A_k$:
 - $r: \mathcal{O} \rightarrow \mathcal{R}$
3. On a des **données**
 - $r(o_1):f(o_1), \dots, r(o_n):f(o_n)$
4. *Hypothèses*:
 - Les données sont des observations d'une variable aléatoire sur $\mathcal{R} \times \mathcal{C}$.
en pratique, on a des biais d'observation
 - Quelle distribution ? On suppose une distribution \mathcal{D} sur \mathcal{O} (proba de "rencontrer" $o \in \mathcal{O}$)
 - Intuitivement : on observe (x,c) avec proba $Pr_o \simeq \mathcal{G}(x=r(o) \wedge f(o)=c)$

16.2

Meilleur cadre théorique

1. Il existe une fonction **idéale**
 - $f: \mathcal{O} \rightarrow \mathcal{C}$.
 - pas forcément toujours vrai ; même ici on peut être incertain**
2. On a une **représentation** de \mathcal{O} via des **attributs** $\mathcal{R} = A_1, \dots, A_k$:
 - $r: \mathcal{O} \rightarrow \mathcal{R}$
3. On a des **données**
 - $r(o_1):f(o_1), \dots, r(o_n):f(o_n)$
4. *Hypothèses*:
 - Les données sont des observations d'une variable aléatoire sur $\mathcal{R} \times \mathcal{C}$.
en pratique, on a des biais d'observation
 - Quelle distribution ? On suppose une distribution \mathcal{D} sur \mathcal{O} (proba de "rencontrer" $o \in \mathcal{O}$)
très difficile d'estimer vraiment cette distribution
 - Intuitivement : on observe (x,c) avec proba $Pr_o \simeq \mathcal{G}(x=r(o) \wedge f(o)=c)$

16.3

Meilleur cadre théorique

1. Il existe une fonction **idéale**

- $f: \mathcal{O} \rightarrow \mathcal{C}$.

pas forcément toujours vrai ; même ici on peut être incertain

2. On a une **représentation** de \mathcal{O} via des **attributs** $\mathcal{R} = A_1, \dots, A_k$:

- $r: \mathcal{O} \rightarrow \mathcal{R}$

3. On a des **données**

- $r(o_1):f(o_1), \dots, r(o_n):f(o_n)$

4. **Hypothèses:**

- Les données sont des observations d'une variable aléatoire sur $\mathcal{R} \times \mathcal{C}$.

en pratique, on a des biais d'observation

- Quelle distribution ? On suppose une distribution \mathcal{D} sur \mathcal{O} (proba de "rencontrer" $o \in \mathcal{O}$)

très difficile d'estimer vraiment cette distribution

- Intuitivement : on observe (x,c) avec proba $Pr_o \simeq \mathcal{D}(x=r(o) \wedge f(o)=c)$

Attention, ensembles infinis en général, on ne peut les mesurer pas point par point

16.4

À la recherche d'un classifieur

1. On a des **données** $x_1 : y_1, \dots, x_n : y_n$

2. Rappel: $x = r(o) \in \mathcal{R}, y = f(o) \in \mathcal{C}$.

3. On cherche $f^*: \mathcal{R} \rightarrow \mathcal{C}$ qui "**approxime**" f .

Idéalement $f^*(x) = y$ où il existe $o \in \mathcal{O}$ tel que $r(o) = x$ et $f(o) = y$.

17

À la recherche d'un classifieur

1. On a des **données** $x_1 : y_1, \dots, x_n : y_n$
2. Rappel: $x = r(o) \in \mathcal{R}$, $y = f(o) \in \mathcal{C}$.
3. On cherche $f^* : \mathcal{R} \rightarrow \mathcal{C}$ qui “*approxime*” f .

Idéalement $f^*(x) = y$ où il existe $o \in \mathcal{O}$ tel que $r(o) = x$ et $f(o) = y$.

- Plusieurs $y \in \mathcal{C}$ possibles en fonction du o choisit (erreur induite par la représentation!)
- Idéalement : on veut le y le **plus probable**, ie qui maximise $Pr(y|x) := Pr(r(o)=x \wedge f(o)=y)$
-

17.1

À la recherche d'un classifieur

1. On a des **données** $x_1 : y_1, \dots, x_n : y_n$
2. Rappel: $x = r(o) \in \mathcal{R}$, $y = f(o) \in \mathcal{C}$.
3. On cherche $f^* : \mathcal{R} \rightarrow \mathcal{C}$ qui “*approxime*” f .

Idéalement $f^*(x) = y$ où il existe $o \in \mathcal{O}$ tel que $r(o) = x$ et $f(o) = y$.

- Plusieurs $y \in \mathcal{C}$ possibles en fonction du o choisit (erreur induite par la représentation!)
- Idéalement : on veut le y le **plus probable**, ie qui maximise $Pr(y|x) := Pr(r(o)=x \wedge f(o)=y)$
- **Même si on trouve f^* avec exactement ce comportement on fera toujours des erreurs.**

17.2

À la recherche d'un classifieur

1. On a des **données** $x_1 : y_1, \dots, x_n : y_n$
2. Rappel: $x = r(o) \in \mathcal{R}, y = f(o) \in \mathcal{C}$.
3. On cherche $f^* : \mathcal{R} \rightarrow \mathcal{C}$ qui “*approxime*” f .

Idéalement $f^*(x) = y$ où il existe $o \in \mathcal{O}$ tel que $r(o) = x$ et $f(o) = y$.

- Plusieurs $y \in \mathcal{C}$ possibles en fonction du o choisit (erreur induite par la représentation!)
- Idéalement : on veut le y le **plus probable**, ie qui maximise $Pr(y|x) := Pr(r(o)=x \wedge f(o)=y)$
- **Même si on trouve f^* avec exactement ce comportement on fera toujours des erreurs.**

Erreur minimale possible : erreur de Bayes.

17.3

Résumé

- **But:** classifier des objets \mathcal{O} dans des classes \mathcal{C} selon un processus idéal **inconnu** $f: \mathcal{O} \rightarrow \mathcal{C}$
en général on suppose même que f elle-même est probabiliste
- **Entrée:**
 - Observations (données) : $x_1 : y_1, \dots, x_n : y_n$
 - x_i sont des représentations $r(o)$ des objets comme des **listes d'attributs**
 - y_i leur classe (**supervisé**)
- **Sortie idéale:** une **hypothèse** h telle que $h(x) = \operatorname{argmin} Pr(y|x)$
- Même la sortie idéale f^* commet des erreurs : l'erreur de **Bayes**.

On peut définir les mêmes choses pour la regression. Cependant, on essaie désormais de minimiser l'espérance de l'**écart** entre h et f .

Classifier = apprendre une distribution (beaucoup de stat!)

18

De l'importance d'une distribution

On cherche à détecter une maladie :

- $\mathcal{C} = \{\textit{malade}, \textit{sain}\}$.
- $\mathcal{O} = \textit{population}$.
- \mathcal{R} contient : tension, résultat d'une prise de sang complète, fièvre etc.
- **99% des gens sont sains.**

19

De l'importance d'une distribution

On cherche à détecter une maladie :

- $\mathcal{C} = \{\textit{malade}, \textit{sain}\}$.
- $\mathcal{O} = \textit{population}$.
- \mathcal{R} contient : tension, résultat d'une prise de sang complète, fièvre etc.
- **99% des gens sont sains.**

h qui renvoie "sain" quelque soit x est une hypothèse raisonnable.

19.1

Limites et impossibilités

Trop d'inconnues dans le problème :

- **impossible** de **calculer** f^*
- si un oracle vous donne f^* , **impossible** de **vérifier** si l'oracle ment!

Défis :

- **trouver** h qui commet une erreur faible par rapport à f^* à partir des données
- **évaluer** l'erreur de h

Limites et impossibilités

Trop d'inconnues dans le problème :

- **impossible** de **calculer** f^*
- si un oracle vous donne f^* , **impossible** de **vérifier** si l'oracle ment!

Défis :

- **trouver** h qui commet une erreur faible par rapport à f^* à partir des données
Où doit-on chercher h ?
- **évaluer** l'erreur de h

Limites et impossibilités

Trop d'inconnues dans le problème :

- **impossible** de **calculer** f^*
- si un oracle vous donne f^* , **impossible** de **vérifier** si l'oracle ment!

Défis :

- **trouver** h qui commet une erreur faible par rapport à f^* à partir des données
Où doit-on chercher h ?
- **évaluer** l'erreur de h
Définir "l'erreur", et différencier entre erreur apparente et erreur réelle.

Trouver une bonne hypothèse

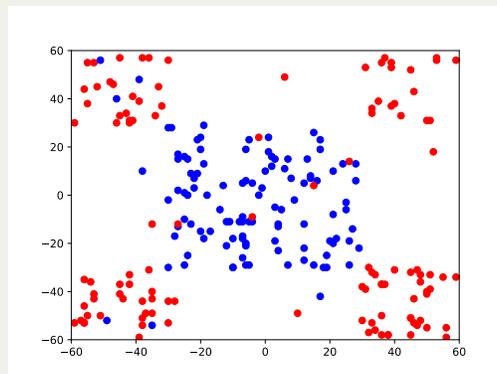
Classification

On cherche h qui sépare les points bleus des points rouges. À quoi doit ressembler h ?

22

Classification

On cherche h qui sépare les points bleus des points rouges. À quoi doit ressembler h ?

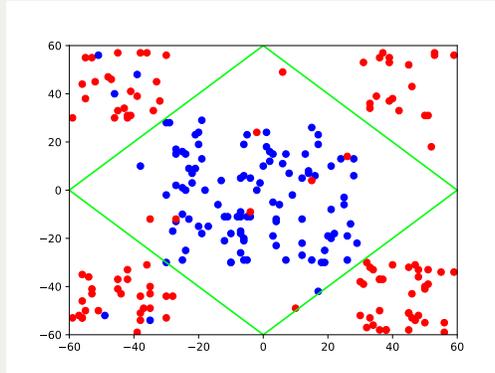


Où tracer les "frontières" ?

22.1

Classification

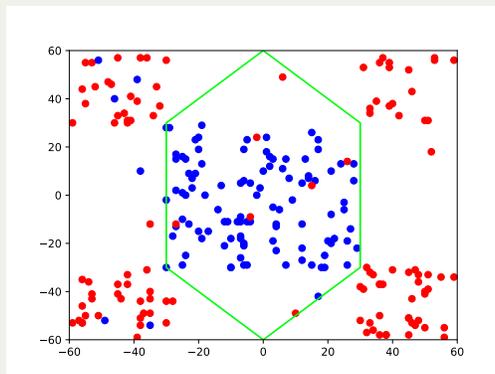
On cherche h qui sépare les points bleus des points rouges. À quoi doit ressembler h ?



4 droites ?

Classification

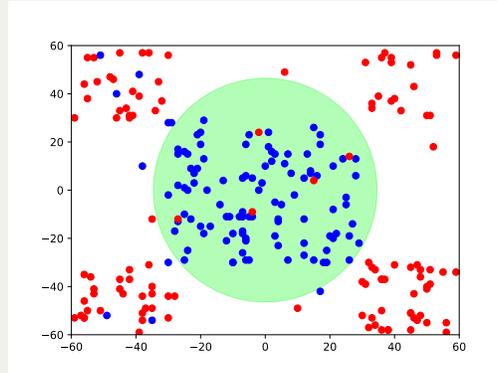
On cherche h qui sépare les points bleus des points rouges. À quoi doit ressembler h ?



6 droites ?

Classification

On cherche h qui sépare les points bleus des points rouges. À quoi doit ressembler h ?

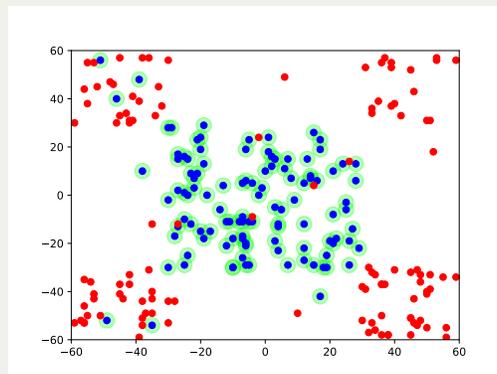


Un cercle ?

22.4

Classification

On cherche h qui sépare les points bleus des points rouges. À quoi doit ressembler h ?



Euh ? Tout plein de petits cercles ?

22.5

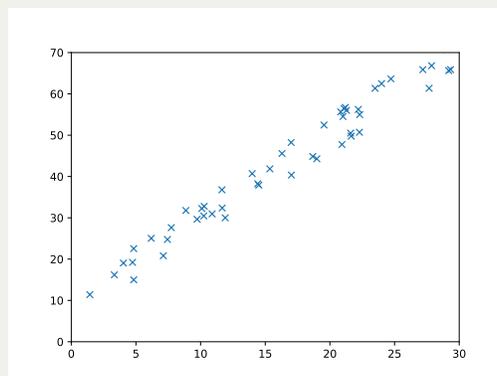
Régression

On veut trouver h qui “approxime” f à partir d'exemple. Quel type de fonction ?

23

Régression

On veut trouver h qui “approxime” f à partir d'exemple. Quel type de fonction ?

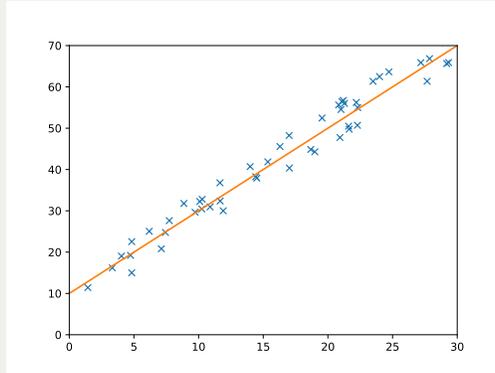


Qu'est-ce que c'est que cette fonction ?

23.1

Régression

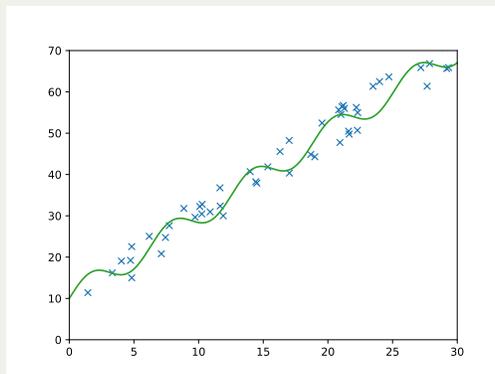
On veut trouver h qui “approxime” f à partir d'exemple. Quel type de fonction ?



Une droite $y = ax + b$

Régression

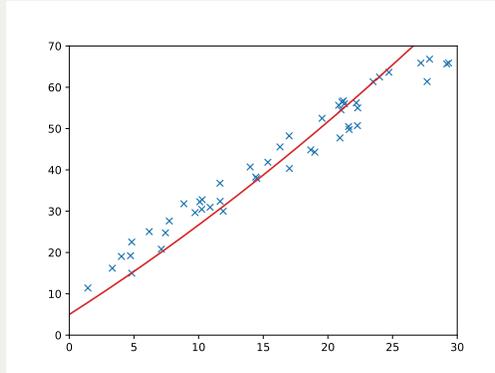
On veut trouver h qui “approxime” f à partir d'exemple. Quel type de fonction ?



Une droite sinusoïdale $y = ax + b\sin(x) + c$

Régression

On veut trouver h qui “approxime” f à partir d'exemple. Quel type de fonction ?

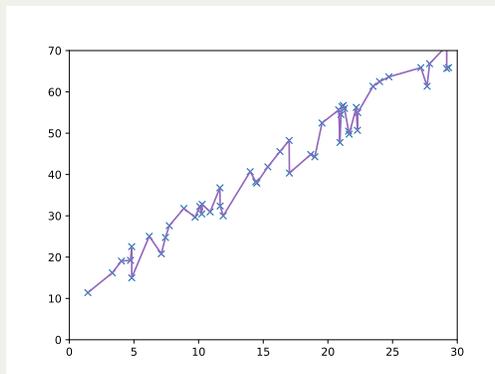


Une parabole $y = ax^2 + bx + c$

23.4

Régression

On veut trouver h qui “approxime” f à partir d'exemple. Quel type de fonction ?



Euh ? Tout plein de petits segments!

23.5

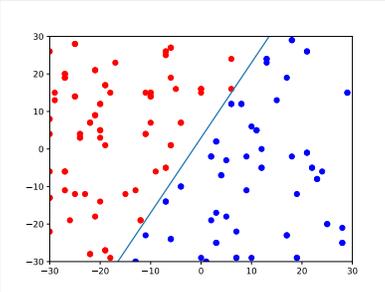
Famille d'hypothèses

Pour trouver une bonne hypothèse, on fixe une famille \mathcal{H} d'hypothèses et on cherche le “meilleur” $h \in \mathcal{H}$ pour nos données.

Exemples pour la classification (deux classes VRAI / FAUX):

Exemples pour la classification (deux classes VRAI / FAUX):

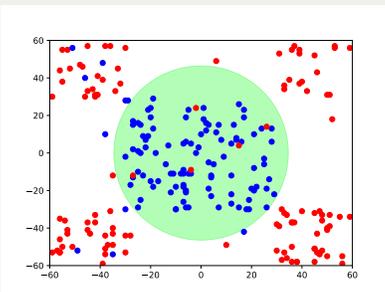
- Une droite $D(a,b,c) : ax + by + c \geq 0$?
 - **On cherche les meilleurs (a,b,c) pour séparer nos données (VRAI au-dessus de la droite, FAUX en dessous) !**



25.1

Exemples pour la classification (deux classes VRAI / FAUX):

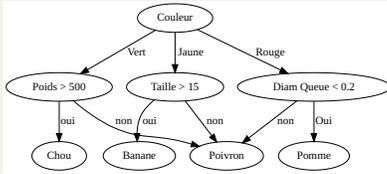
- Une droite $D(a,b,c) : ax + by + c \geq 0$?
 - **On cherche les meilleurs (a,b,c) pour séparer nos données (VRAI au-dessus de la droite, FAUX en dessous) !**
- Un cercle de rayon r centrée en $(a,b) : C(r,a,b) : (x-a)^2 + (y-b)^2 \leq r^2$
 - **On cherche le meilleur rayon r et le meilleur centre (a,b) pour séparer nos données !**



Nos hypothèses sont définies à partir de paramètres. On essaie d'optimiser ces paramètres.

25.2

Exemples pour la classification (deux classes VRAI / FAUX):



- Une droite $D(a,b,c) : ax + by + c \geq 0$?
 - **On cherche les meilleurs (a,b,c) pour séparer nos données (VRAI au-dessus de la droite, FAUX en dessous) !**
- Un cercle de rayon r centrée en $(a,b) : C(r,a,b) : (x-a)^2 + (y-b)^2 \leq r^2$
 - **On cherche le meilleur rayon r et le meilleur centre (a,b) pour séparer nos données !**
- Arbres de décisions (voir TP 3)
 - **on classe en répondant à une suite de questions ; permet de prendre des attributs discrets en compte.**

Nos hypothèses sont définies à partir de paramètres. On essaie d'optimiser ces paramètres.

Évaluation

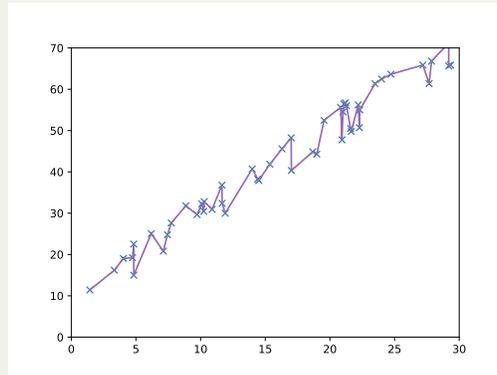
A-t-on bien appris ?

- on a une **petite erreur** sur les données (erreur **apparente**)
- on **généralise** à des exemples qu'on a jamais vu (erreur **réelle**)

Évaluation

A-t-on bien appris ?

- on a une **petite erreur** sur les données (erreur **apparente**)
- on **généralise** à des exemples qu'on a jamais vu (erreur **réelle**)



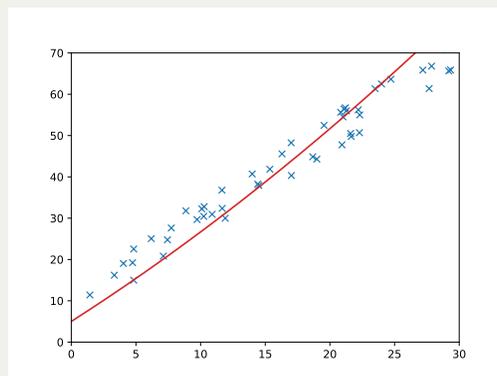
You overfit it: petite erreur apparente, grande erreur réelle

26.1

Évaluation

A-t-on bien appris ?

- on a une **petite erreur** sur les données (erreur **apparente**)
- on **généralise** à des exemples qu'on a jamais vu (erreur **réelle**)



You underfit it: grande erreur tout court

26.2

Évaluation d'un classifieur

Erreur **observée** d'une hypothèse h :

- $\hat{E}(h) = \#\{x_i \mid y_i \neq h(x_i)\}/n$,
- ie, la proportion de points où h se trompe.
- **facile à calculer**

Erreur **réelle** d'une hypothèse h :

- $E(h) = \Pr(y \neq h(x) \mid f(x)=y)$
- ie, la probabilité que h se trompe.
- **Comment l'estimer?**

Erreur observée n'a aucun sens si on inclut des données apprises.

27

Séparer apprentissage du test

- On apprend avec un sous-ensemble $A \subseteq D$ des données **tiré aléatoirement**
- On estime l'erreur avec $D \setminus A$

Airline	Flight	AirportFrom	AirportTo	DayOfWeek	Time	Length	Delay
CO	269	SFO	IAH	3	15	205	1
US	1558	PHX	CLT	3	15	222	0
...
AA	2400	LAX	DFW	3	20	165	0
AA	2466	SFO	DFW	3	20	195	1
...

Problème : une grosse partie des données n'est pas utilisée pour l'entraînement

28

Validation croisée

- On sépare les données D en deux ensembles D_1, D_2 de taille égale
- On apprend :
 - h sur D ,
 - h_1 sur D_1 , on calcule l'erreur e_1 de h_1 sur D_2
 - h_2 sur D_2 , on calcule l'erreur e_2 de h_2 sur D_1
- L'**erreur** de h est $e_1 + e_2$.

Airline	Flight	AirportFrom	AirportTo	DayOfWeek	Time	Length	Delay
CO	269	SFO	IAH	3	15	205	1
US	1558	PHX	CLT	3	15	222	0
...
AA	2400	LAX	DFW	3	20	165	0
AA	2466	SFO	DFW	3	20	195	1
...

Coûteux, non prouvé en théorie mais fonctionne en pratique (en coupant en plus que 2)

Mettre en place une solution d'apprentissage supervisé

Données

- Choisir un modèle de **représentation** : **quels attributs ?**
- Collecter des **données** : **échantillon représentatif ? assez grand ?**

Apprentissage

- Choisir une famille d'hypothèses : **exploration préliminaire des données, évaluation sur des petits sous-ensembles**
- Trouver la meilleure hypothèse : **algorithme d'apprentissage adapté à la famille**

Évaluation

- Choisir une méthode d'évaluation
- Évaluer les performances sur une partie des données **non utilisées pour l'entraînement**.

Exemple complet : la séparation linéaire

31

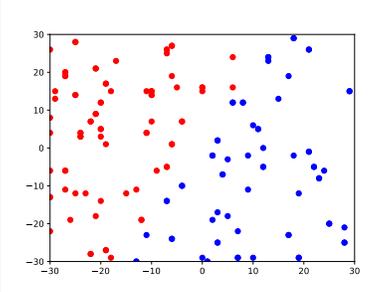
Classification binaire

- On a deux paramètres (abscisse / ordonnée)
- Deux classes : **+1** / **-1**

32

Classification binaire

- On a deux paramètres (abscisse / ordonnée)
- Deux classes : **+1** / **-1**

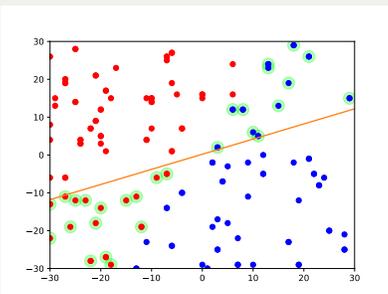


- Hypothèse : on cherche une hypothèse de la forme $SIGNE(y-(ax+b)) \in \{-1, +1\}$

32.1

Classification binaire

- On a deux paramètres (abscisse / ordonnée)
- Deux classes : **+1** / **-1**

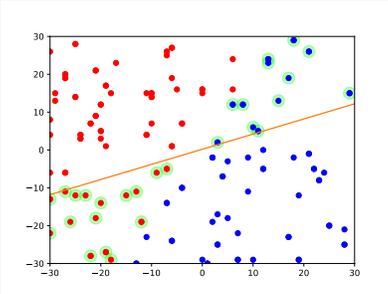


- Hypothèse : on cherche une hypothèse de la forme $SIGNE(y-(ax+b)) \in \{-1, +1\}$
- Erreur $E(a,b)$: nombre de point mal classés

32.2

Classification binaire

- On a deux paramètres (abscisse / ordonnée)
- Deux classes : **+1** / **-1**

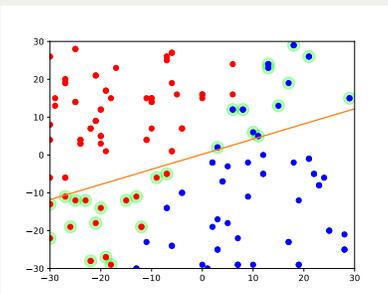


- Hypothèse : on cherche une hypothèse de la forme $SIGNE(y-(ax+b)) \in \{-1, +1\}$
- Erreur $E(a,b)$: nombre de point mal classés
- Apprendre = problème d'optimisation. Trouver les meilleurs a, b !

32.3

Classification binaire

- On a deux paramètres (abscisse / ordonnée)
- Deux classes : **+1** / **-1**



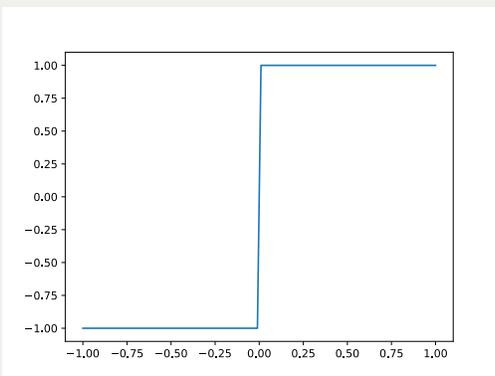
- Hypothèse : on cherche une hypothèse de la forme $SIGNE(y-(ax+b)) \in \{-1, +1\}$
- Erreur $E(a,b)$: nombre de point mal classés
- Apprendre = problème d'optimisation. Trouver les meilleurs a, b !
- Problème: **E n'est pas dérivable, on ne sait rien faire...**

32.4

Erreur dérivable et descente de gradient

33

Erreur dérivable et descente de gradient

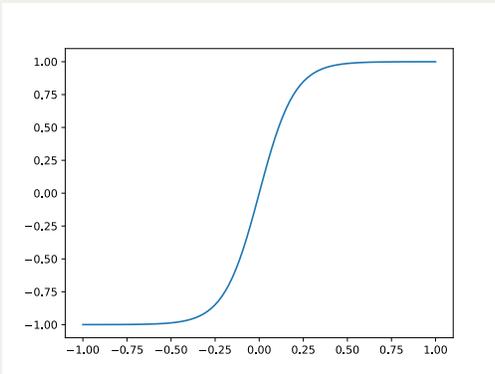


SIGNE(x)

$$E(a,b) = (1/4)\sum_i(c_i - \text{SIGNE}(y_i - (ax_i + b)))^2$$

33.1

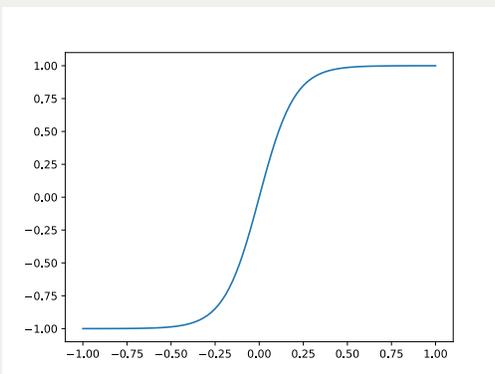
Erreur dérivable et descente de gradient



$$\sigma(x) = 2e^{10x}/(1+e^{10x}) - 1$$

$$\hat{E}(a,b) = (1/4)\sum_i(c_i - \sigma(y_i - (ax_i + b)))^2$$

Erreur dérivable et descente de gradient

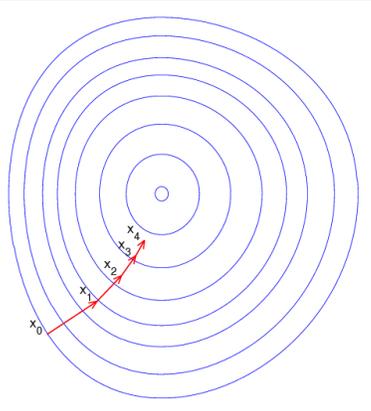


$$\sigma(x) = 2e^{10x}/(1+e^{10x}) - 1$$

$$\hat{E}(a,b) = (1/4)\sum_i(c_i - \sigma(y_i - (ax_i + b)))^2$$

On va minimiser \hat{E} par descente de gradient !

Descente de gradient (Méthode de Newton)



- On part avec (a_0, b_0) aléatoires
- On définit : $(a_{n+1}, b_{n+1}) = \hat{E}(a_n, b_n) - \alpha \nabla \hat{E}(a_n, b_n)$
où $\nabla \hat{E} = (\partial \hat{E} / \partial a, \partial \hat{E} / \partial b)$ et $\alpha > 0$ est un paramètre de vitesse

Converge vers un minimum local.

Voir animation *Why Momentum Really Work*, Gabriel GOH

Animation

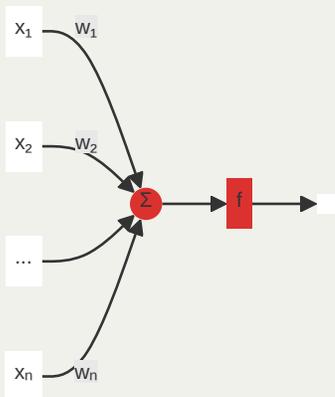
TensorFlow

Réseaux de neurones

36

Perceptron (1957)

L'algorithme qu'on vient de voir est exactement comment un neurone formel va apprendre :

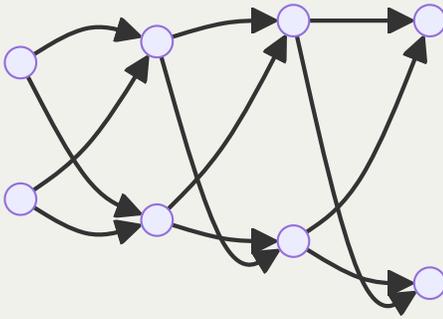


- $F = \sigma(\sum_i w_i x_i)$
- But : Minimiser $\sum_{(x,y) \in D^y} -F(x)$

Limites : on n'apprend que des données linéairement séparables (voir exemple précédent).

37

Réseau de neurone (artificiels)



Chaîner des neurones pour faire des fonctions plus complexes.

- Fonction calculée est toujours dérivable
- On sait la dériver efficacement (dérivation chaînée des neurones, avec mémorisation)
- On sait mettre les poids à jour via la méthode du gradient

On a une méthode qui fait converger le réseau de neurone vers un optimum local!

38

Réseaux de neurones

Intêret pour les réseaux de neurones renouvelé :



Un GPU (ou carte graphique, crédit : Henry Mühlpfordt)

- Avancée technologique ont permis de passer à une échelle supérieur (GPU)
- Applications très prometteuses (reconnaissance d'images, transformers, etc.)
- Données d'apprentissage bien plus facile d'accès (internet)

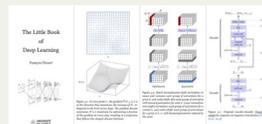
39

Quelques références

Livres



Apprentissage machine, Clé de l'Intelligence Artificielle, Rémi Gilleron



The Little Book of Deep Learning, François Fleuret

Neural Networks and Deep Learning

Neural Networks and Deep Learning is a free online book. The book will teach you about:

- Neural networks, a beautiful biologically-inspired programming paradigm which enables a computer to learn from observational data
- Deep learning, a powerful set of techniques for learning in neural networks

Neural networks and deep learning currently provide the best solutions to many problems in image recognition, speech recognition, and natural language processing. This book will teach you many of the core concepts behind neural networks and deep learning.

For more details about the approach taken in the book, [see here](#). Or you can jump directly to [Chapter 1](#) and get started.

Neural Network and Deep Learning, Michael Nielsen

Autres

- [Série de 3Blue1Brown sur les réseaux de neurones](#)
- [Machine Learning Crash Course \(Google\)](#)