

La reconnaissance automatique d'écriture

Théodore Bluche



DE QUOI VAIS-JE PARLER?

- Mon travail : un doctorat en informatique
- Mon sujet : la reconnaissance automatique d'écriture
- Ma collaboration avec des étudiants allemands

MON PARCOURS

Supélec - Diplôme d'Ingénieur
Electronique, Informatique, Automatique,
Electrotechnique, ...

Ingénieur en
Informatique

Baccalauréat
"Scientifique"



2005

2007

2009

2010

2011

Classes préparatoires
Maths / Physique
Préparation au concours d'entrée
aux Grandes Ecoles d'Ingénieurs



Oxford University
Master en
Informatique



Début de mon
doctorat

LE DOCTORAT



Moi



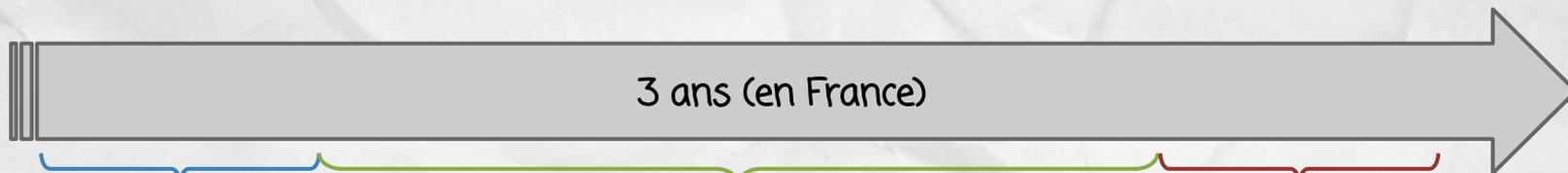
Directeur de thèse

Prof. Hermann Ney



Co - encadrant

Christopher Kermorvant



3 ans (en France)



Lire des articles
Connaitre le sujet

Exprériences



Ecrire la thèse



JE TRAVAILLE ...

Lieux de Travail

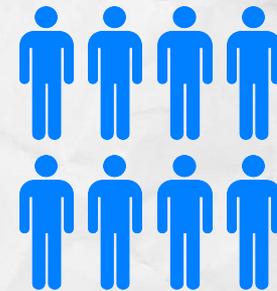


mardi et mercredi

Qu'est ce que c'est ?

un laboratoire de
recherche qui est
associé à l'université
de Paris

Mes collègues



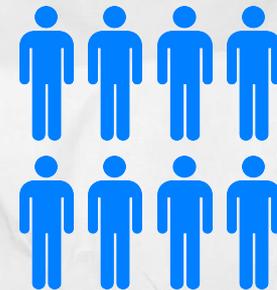
étudiants,
professeurs,
chercheurs



lundi, jeudi et
vendredi

une entreprise

je travaille dans le
département
Recherche et
Développement



ingénieurs

MON TRAVAIL...

Je fais de la *recherche* ...

... *en informatique*



Lire des articles scientifiques sur le sujet

- > pour savoir ce qui existe déjà
- > pour trouver des idées



Ecrire des programmes...

... et utiliser des programmes qui existent déjà



Mais aussi

- trouver des idées
- écrire des articles...



Assister à des conférences

- > pour connaître les dernières technologies
- > pour rencontrer des gens qui travaillent sur le même sujet



Exécuter les programmes --
c'est l'ordinateur qui travaille! même
la nuit et le week-end

MON SUJET

Termes techniques... -> méthodes utilisées

Nouveaux modèles à base de réseaux de neurones et de modèles de Markov cachés pour la reconnaissance d'écriture manuscrite à large vocabulaire

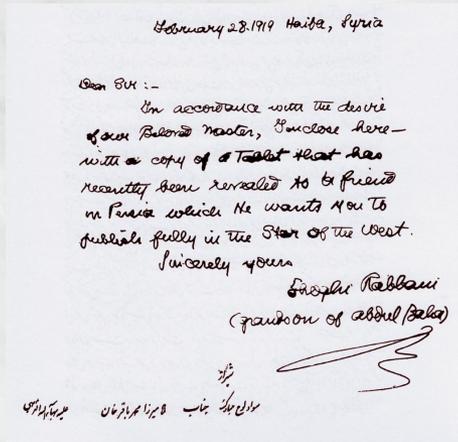
Sujet de recherche



POURQUOI UTILISER DES ORDINATEURS ?



A QUOI ÇA SERAIT?



-> Traitement automatisé du courrier entrant

... dans les grandes entreprises par exemple,
pour envoyer les courriers aux bons services

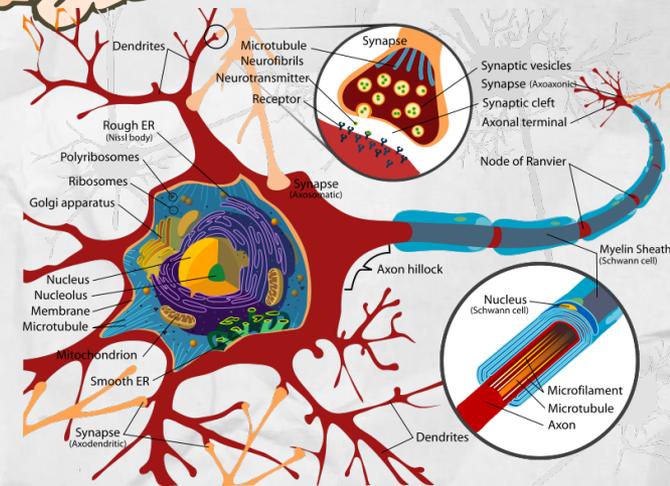
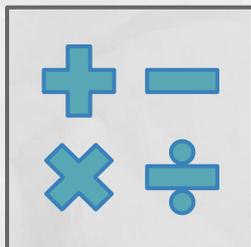
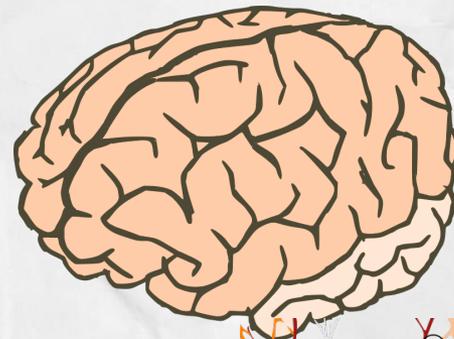
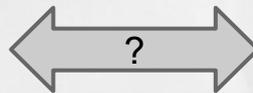
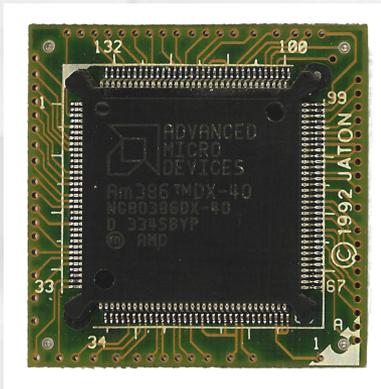
Numérisation d'archives <-

... ça prend moins de place dans un ordinateur

... ça permet de faire des recherches plus facilement avec un ordinateur



L'INTELLIGENCE ARTIFICIELLE



ORDINATEUR = MACHINE



vous

Image

L'Homme voit l'image, reconnaît que c'est du texte, sait interpréter que le noir correspond à l'écriture et le blanc au fond

VOUS

Texte

L'Homme sait lire le texte, et sait ce que ça veut dire

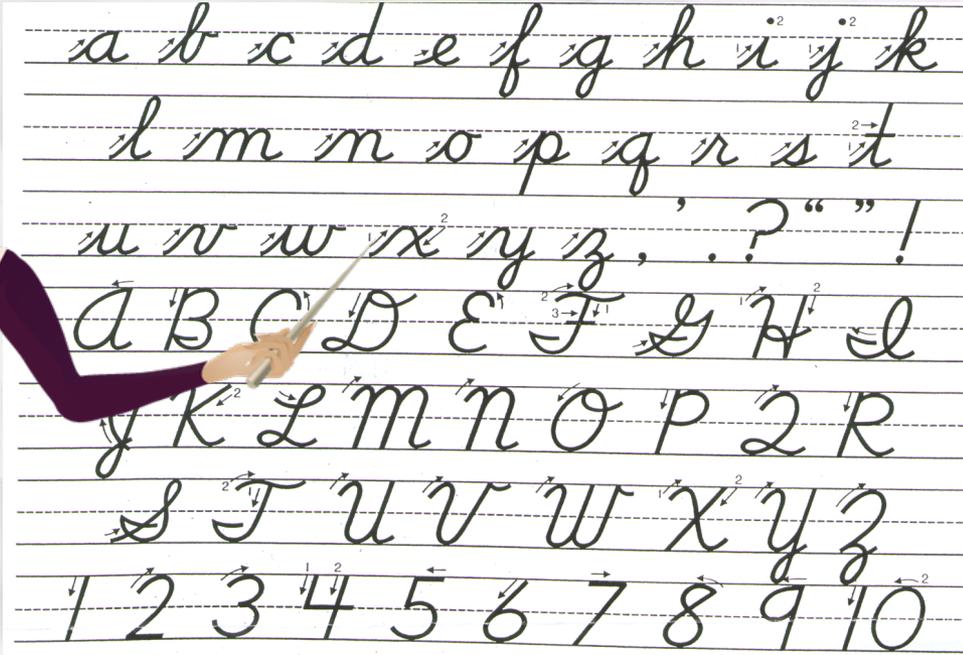
```
11100110010011010101011100000010001110011101101001011111  
10001110010111101111011001011001011000011110001100010001  
0110011001100010001000100110110101110011110001110011001  
00000010101011110100000000111001011100101010101010101  
010000100001110110000101110011110000...
```

L'ordinateur sait que c'est une image et sait comment l'afficher à l'écran

```
00010001101011000110111011100001110100110001000100111110  
1010001101000101011111000000001010100001001010010001110  
0010011101010001111010111101011110011010001011110101101000  
100000000110100100100001110010011110100100001101 ...
```

L'ordinateur sait que c'est du texte et sait comment l'afficher

APPRENDRE À LIRE ...



APPRENDRE À LIRE À UN ORDINATEUR

écrit

écrit

reconnait

reconnaisant

vous

vous

documentation

documentation

-> On utilise des bases de données de plus de 100,000 images + texte de mots

-> L'ordinateur sait traiter ces données très vite



EST-CE UN PROBLÈME FACILE?

Sentence Database

A01-122

"Mr. Powell finds it easier to take it out of mothers, children and sick people than to take on this war industry," Mr. Brown commented icily. "Let us have a full inquiry into the cost of drugs and the pharmaceutical industry." The health of children today owed much to the welfare food scheme. It was maintained during the war. Now in conditions of Tory affluence it seemed it could not be carried on.

„Mr. Powell finds it easier to take it out of mothers, children and sick people than to take on this war industry," Mr. Brown commented icily. „Let us have a full inquiry into the cost of drugs and the pharmaceutical industry." The health of children today owed much to the welfare food scheme. It was maintained during the war. Now in conditions of Tory affluence it seemed it could not be carried on.

Name: Alexander Debus

Il faut :

-> trouver les lignes de texte

-> les couper en mots

-> découper les mots en caractères

DEVINETTE :

Quelles sont ces lettres ... ?

ee e e

4

SOLUTION :

u m en n
ee u ee 4
↓ ↓ ↓ ↓
documentation

→ documentation

vous

111001100100110101010111000000100011100111011010010111
1110001110010111011110110010111000011110001100010
0010110011001100010001000100110101011100111100011100
1100100000001010101110100000000111001011100101010010
101011010000100001110110000101110011110000...

L'ordinateur ne "voit" pas : il ne sait pas interpréter ce qui est dans l'image, seulement l'afficher

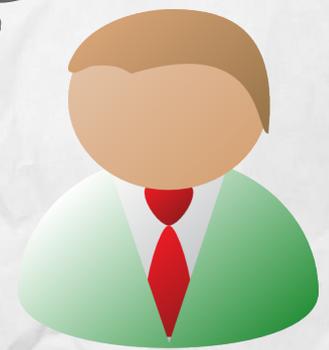
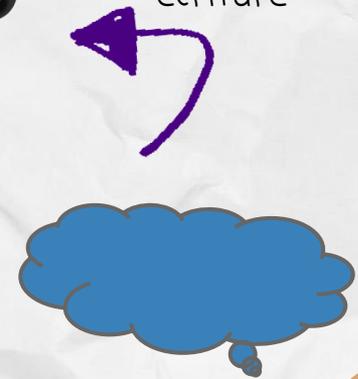
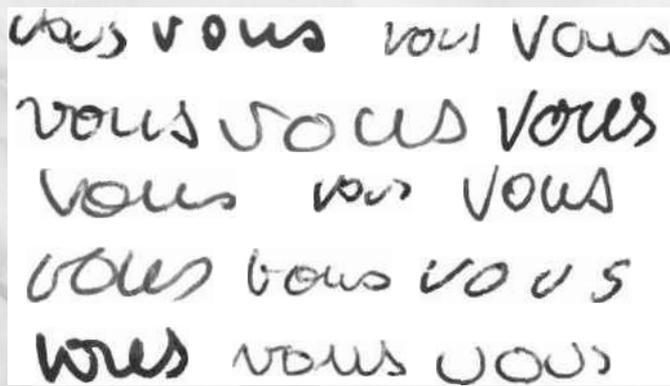
u m en n
ee ee ee
↓ ↓ ↓ ↓
documentation

C'est difficile de découper le mot en lettres



Nous allons utiliser ce que nous savons de ce problème pour construire un système facile à exécuter pour un ordinateur afin d'essayer de lui faire reconnaître l'écriture

Chaque personne a son propre style d'écriture



COMMENT FAIRE POUR QUE LES IMAGES SE "RESSEMBLENT" PLUS ?

responsable

Rendre l'image plus nette en **augmentant le contraste**

responsable

Forcer toutes les images à **avoir la même taille**

responsable

Faire en sorte que l'**écriture soit droite (et non penchée)**

Chaque personne a son propre style d'écriture

vous vous vous vous
vous vous vous
vous vous vous
vous vous vous
vous vous vous

responsable

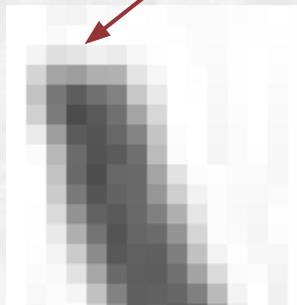
-> Il existe aussi **d'autres techniques ...**

vous

```
1100110010011010101011000000100011100111011010010111  
1100011100101110111011100101100101100001110001100010  
001011001100110001000100010011010101100111100011100  
11001000000101010111010000000011001011100101010010  
101011010000100001110110000101110011110000...
```

L'ordinateur ne "voit" pas : il ne sait pas interpréter ce qui est dans l'image, seulement l'afficher

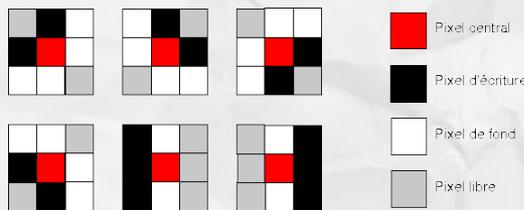
vous



Une image est composée de **pixels** :
-> pour l'ordinateur, cela correspond à un **nombre**, qui représente la couleur, l'intensité, ...



On peut utiliser notre connaissance du problème pour **extraire des caractéristiques (des nombres)** :
-> le nombre de pixels noirs par rapport au nombre de pixels blancs
-> la façon dont les pixels noirs sont organisés : en **comptant** combien de fois on rencontre des configurations simples



Les ordinateurs aiment les nombres, et savent faire beaucoup de choses avec...

-> on peut utiliser la valeur des pixels tout simplement

COMMENT AIDER L'ORDINATEUR À INTERPRÉTER LE CONTENU DE L'IMAGE ?

COMMENT RECONNAÎTRE LES LETTRES DONT LE MOT EST COMPOSÉ ?



C'est difficile de découper le mot en lettres

On pourrait essayer de reconnaître le mot complet ...

-> c'est possible pour les chèques car il n'y a pas beaucoup de mots (moins de 100)

-> en général, **il y a beaucoup de mots** possibles (plus de 50,000), alors qu'il y a **moins de 100 caractères différents** (A,B,C,...,a,b,c,0,1,...)

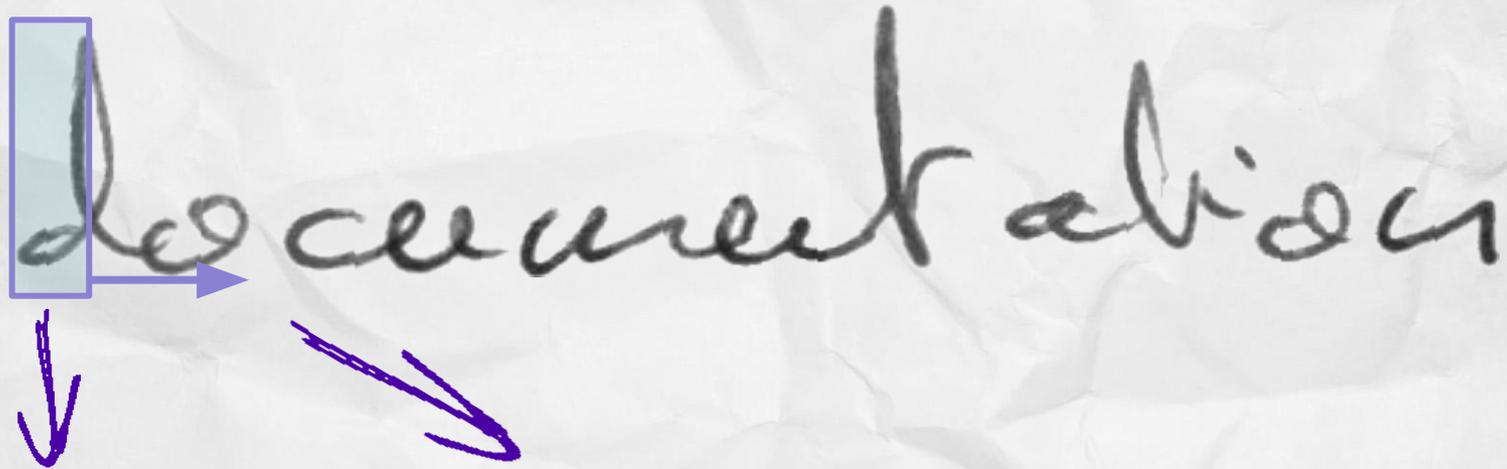
-> le "m" de "maman" et celui de "maison" se ressemblent : on veut en profiter!

COMMENT RECONNAÎTRE LES LETTRES DONT LE MOT EST COMPOSÉ ?

u m en n
ee ee e
↓ ↓ ↓ ↓
documentation

C'est difficile de découper le mot en lettres

Nous allons lire le mot et essayer de reconnaître les lettres au fur et à mesure :



On commence à gauche, et on regarde une petite partie de l'image ...

... ensuite, on se déplace un petit peu vers la droite ... et on continue comme ça jusqu'à la fin...

COMMENT RECONNAITRE LES LETTRES DONT LE MOT EST COMPOSÉ ?

MAIS ...

documentation

Des fois, on ne regarde qu'un bout de la lettre...

documentation

Des fois, on regarde plus qu'une seule lettre ...

documentation

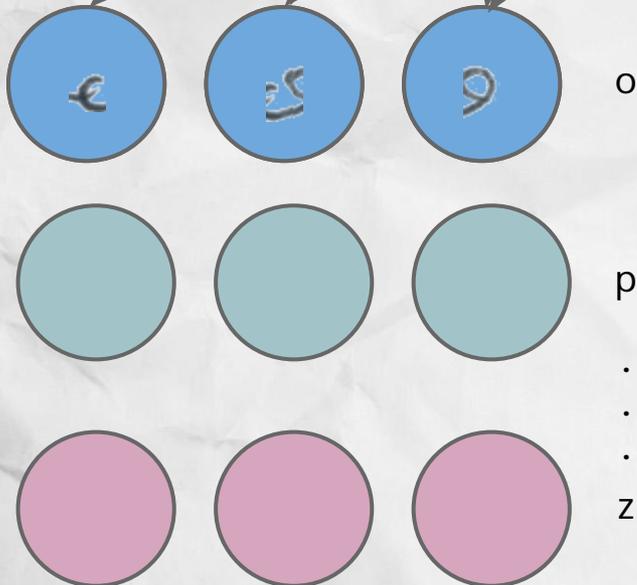
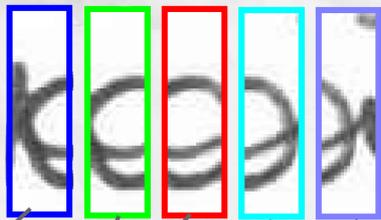
Des fois, on tombe entre deux lettres ...

u m en n
ee ee e
↓ ↓ ↓ ↓
documentation

C'est difficile de découper le mot en lettres

COMMENT RECONNAITRE LES LETTRES DONT LE MOT EST COMPOSÉ ?

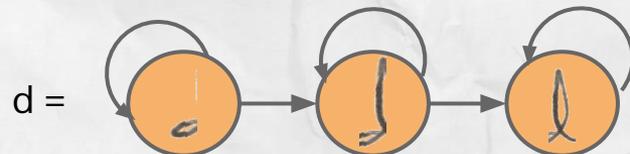
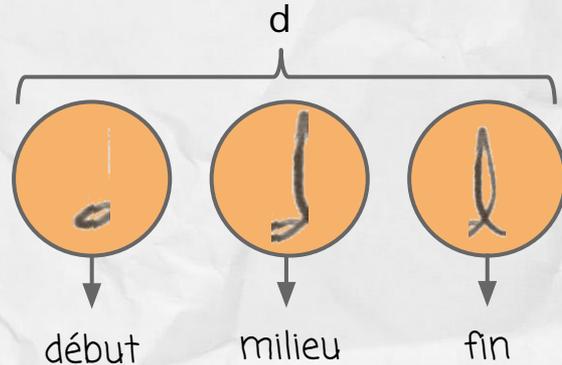
voisie



On ne va pas reconnaître directement des lettres mais plutôt des morceaux de lettres

u m en n
ee ee ee
↓ ↓ ↓ ↓
doceurment abien

C'est difficile de découper le mot en lettres



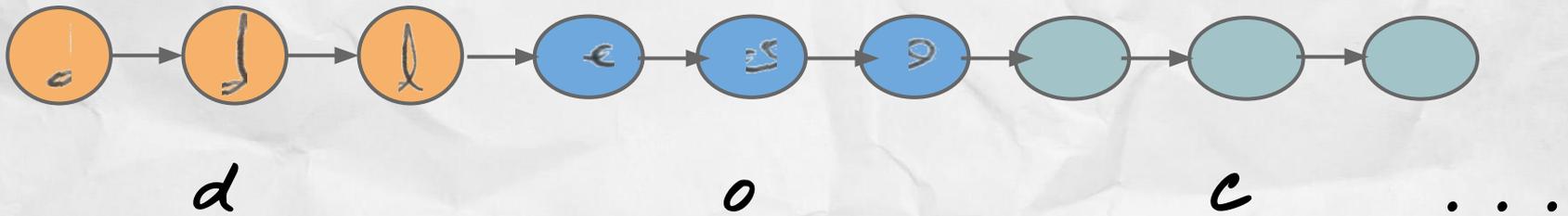
COMMENT RECONNAITRE LES LETTRES DONT LE MOT EST COMPOSÉ ?

u m en n
ee ee ee n
↓ ↓ ↓ ↓
do ce u ment ab i on

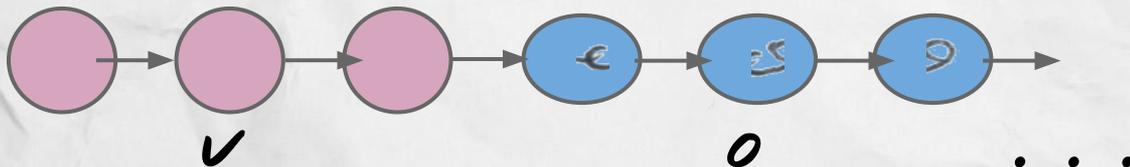
C'est difficile de découper le mot en lettres

On reconnaît des morceaux de lettres ... cela permet de reconnaître des lettres, puis des mots

do ce u ment ab i on



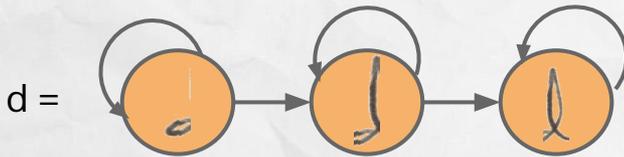
vous



COMMENT RECONNAITRE LES LETTRES DONT LE MOT EST COMPOSÉ ?

u m en n
ee ee ee n
↓ ↓ ↓ ↓
do ce ment a bio n

C'est difficile de découper le mot en lettres

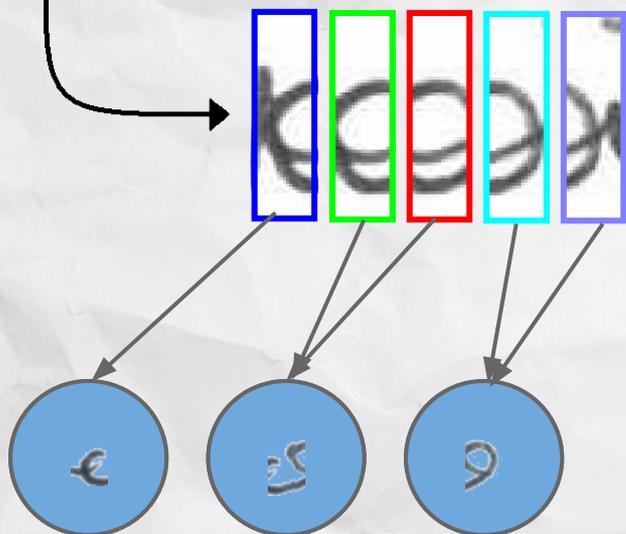


MON SUJET

Nouveaux modèles à base de réseaux de neurones et de modèles de Markov cachés pour la reconnaissance d'écriture manuscrite à large vocabulaire

COMMENT RECONNAITRE LES LETTRES DONT LE MOT EST COMPOSÉ ?

voisie



Comment fait-on pour associer un morceau d'image au bon modèle ?

Souvenez-vous ...



Extraction de caractéristiques

Nombres

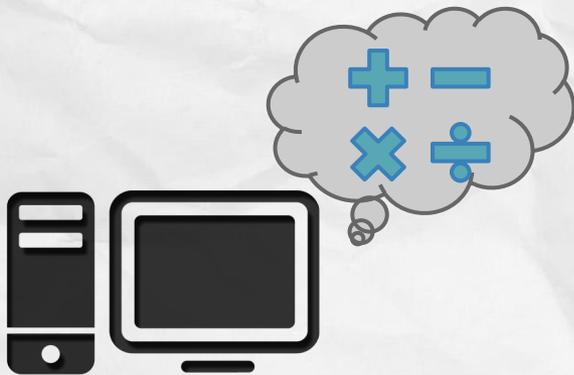
- 1.28
- 5.42
- 3.14
- 0.69
- 7.98

on utilise ces nombres plutôt que l'image

u m en n
ee ee ee
↓ ↓ ↓ ↓
doceumentation

C'est difficile de découper le mot en lettres

LES RÉSEAUX DE NEURONES



L'ordinateur est très fort pour faire des opérations mathématiques simples ...

Nombres

on multiplie chaque nombre

$$\times 5 = 6.40$$

$$\times 2 = 10.84$$

$$\times 1 = 3.14$$

$$\times 3 = 2.97$$

$$\times 4 = 31.92$$

1.28

5.42

3.14

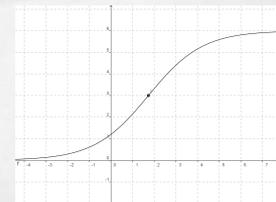
0.69

7.98

on obtient un nouveau nombre



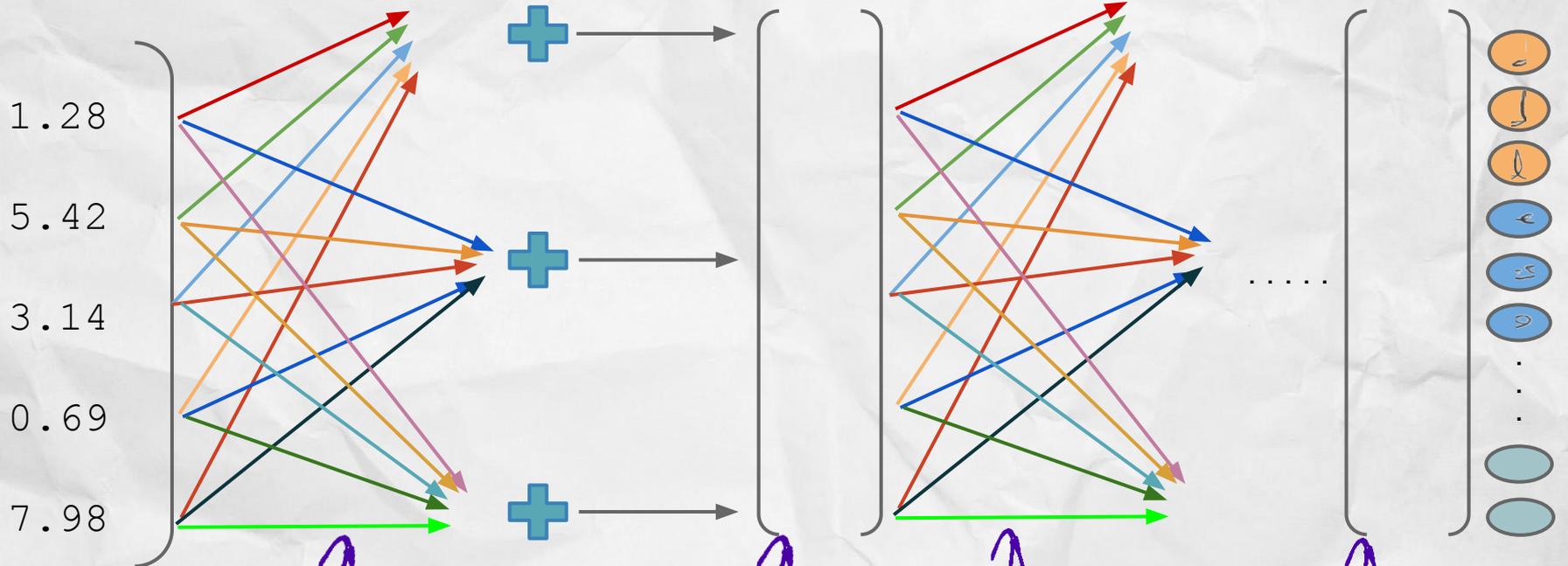
on additionne tous les nombres



-> pourquoi fait-on cela ?

LES RÉSEAUX DE NEURONES

-> pourquoi fait-on cela ?



On peut faire ça plusieurs fois, en multipliant par des nombres différents...

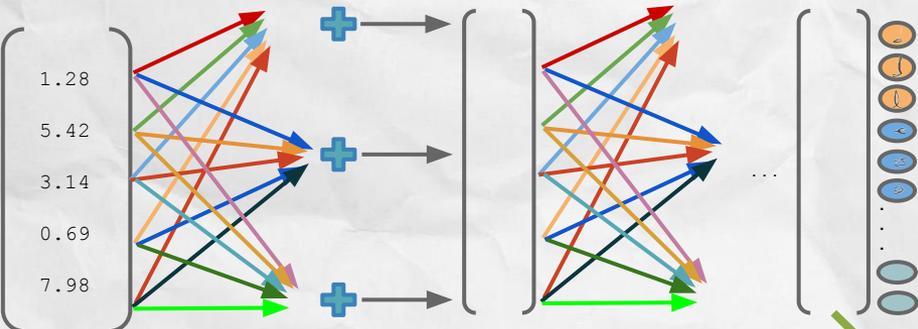
... ainsi, on obtient plein de nouveaux nombres à utiliser ...

... on peut faire ça plusieurs fois, en multipliant par des nombres différents...

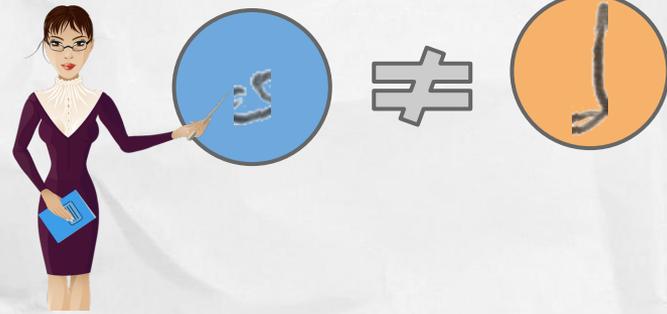
A la fin, on veut que chaque nombre représente la probabilité d'un morceau de lettre

comment faire ?

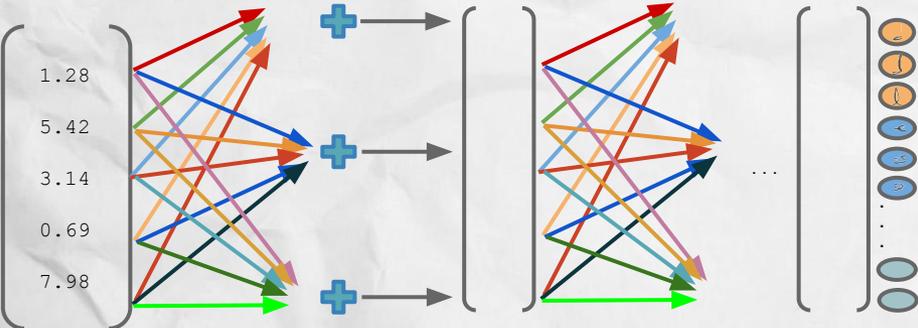
LES RÉSEAUX DE NEURONES



- > On commence avec des nombres au hasard...
- > On fait toutes les opérations...
- > On regarde les probabilités obtenues...



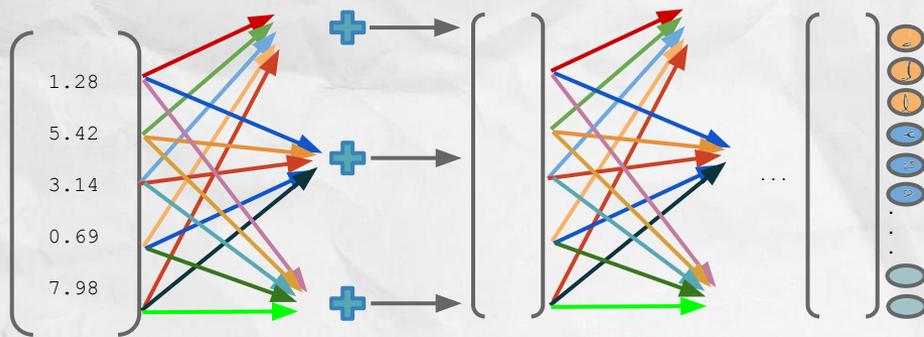
-> On compare avec le résultat attendu...



On modifie les nombres par lesquels on multiplie pour augmenter les chances de trouver la bonne réponse la prochaine fois...

 c'est mathématique !

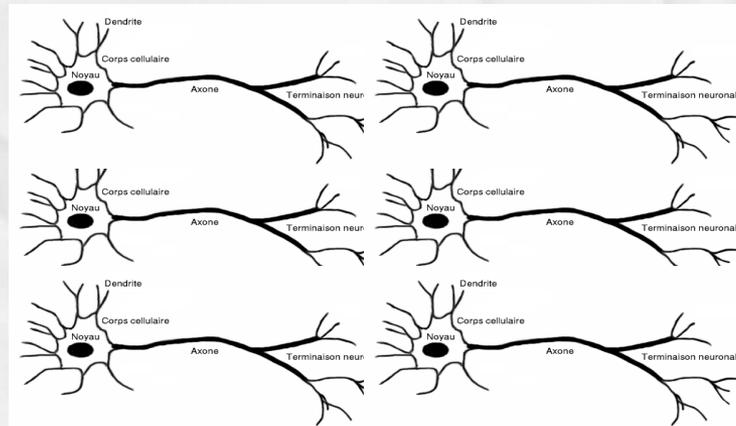
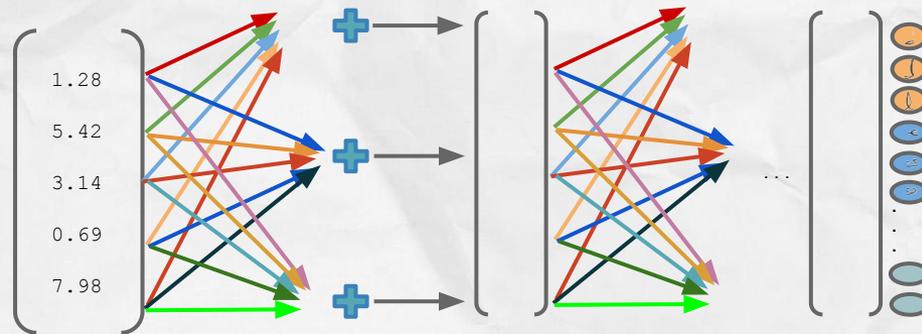
LES RÉSEAUX DE NEURONES



MON SUJET

Nouveaux modèles à base de **réseaux de neurones**
et de **modèles de Markov cachés** pour la
reconnaissance d'écriture manuscrite à large
vocabulaire

LES RÉSEAUX DE NEURONES



COMMENT PEUT-ON ENCORE AIDER L'ORDINATEUR ?

Souvenez-vous ...



L'ordinateur ne "sait" pas que **Bonjour** est un mot correct et que **jBoonru** ne veut rien dire ...



Mais c'est un problème plus simple ...

-> On n'autorise le système qu'à reconnaître des mots parmi une liste de mots autorisés

= un **dictionnaire**

COMMENT PEUT-ON ENCORE AIDER L'ORDINATEUR ?

L'ordinateur ne comprend pas les textes, il ne "sait" pas que

Bonjour, je suis heureux d'être ici.

est correct alors que

Bon jour. Je Suisse heure eux de tri si

ne veut rien dire ...



On va construire un **modèle de langue**

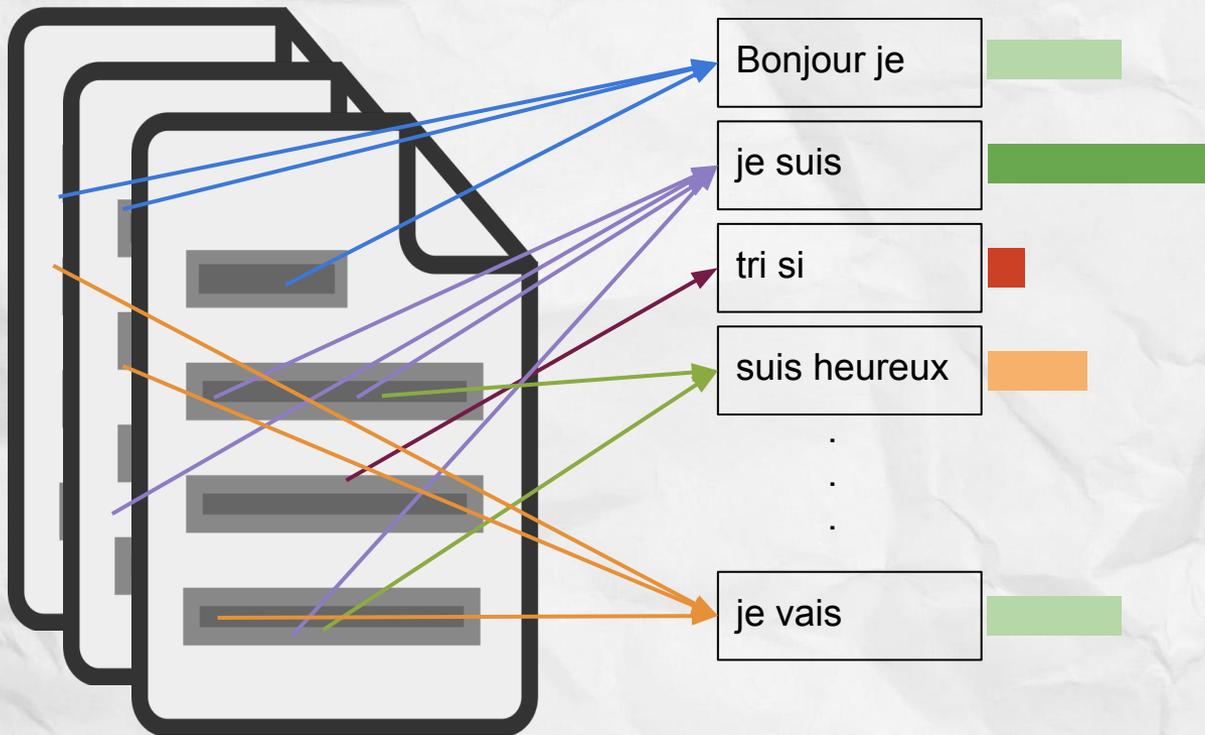
LES MODÈLES DE LANGUE

Les modèles statistiques

- > Pour dire qu'une phrase est correcte ou non, **il faut des connaissances** (vocabulaire, grammaire, compréhension de la signification de la phrase)
- > Au lieu d'essayer de donner ces connaissances à un ordinateur, on lui fait lire beaucoup de texte
- > On utilise des outils **statistiques** pour l'aider à se souvenir des suites de mots qu'il a lu
- > Ensuite, il sera capable de dire **s'il a vu** ces mots-là, dans cet ordre-là, **souvent, rarement ou jamais**

LES MODÈLES DE LANGUE

Les modèles statistiques

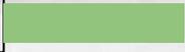


Beaucoup de textes (faciles à trouver grâce à internet)

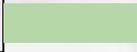
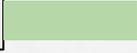
On compte combien de fois on voit chaque séquence
-> ensuite, on peut dire quelles séquences on voit souvent

LES MODÈLES DE LANGUE

1 mot

et	
de	
le	
vous	

2 mots

et le	
de vous	
je suis	
il est	

trop de mots

et le moment de vous demander	
je vous prie de bien vouloir	
je suis un peu déçu de	
il est toujours très en retard	

commande 



25 juin 



le 25 juin 1987, la commande 



Permet de savoir quels mots sont fréquents, mais ça ne nous aide pas beaucoup pour des suites de mots (des phrases par exemple)

Permet de savoir quels mots suivent souvent d'autres mots, mais ça nous aide un peu plus

On verra trop rarement plusieurs fois la même suite de mots, et cela ne nous aidera pas beaucoup.

LES MODÈLES DE LANGUE

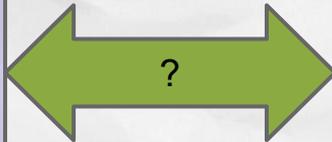
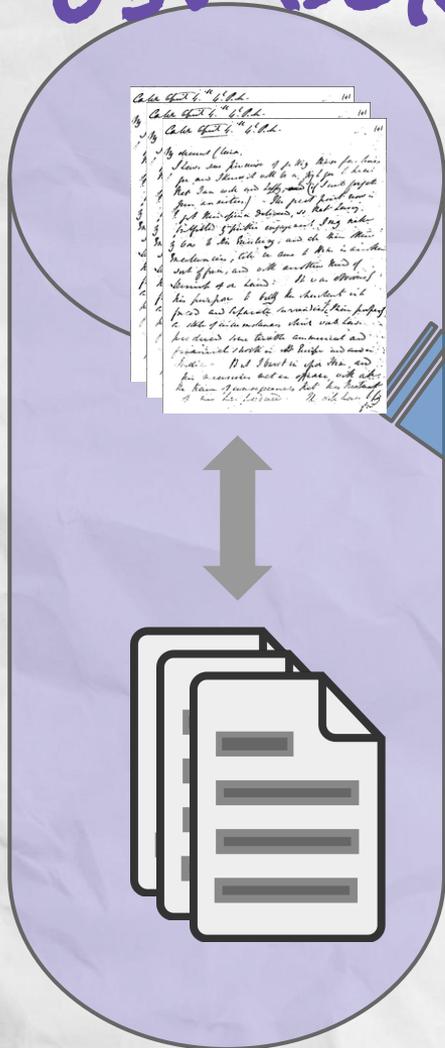
Bonjour, je suis heureux d'être ici.



Bon jour. Je Suisse heure eux de tri si



COMMENT SAVOIR SI LE SYSTÈME EST BON ?



15 %



-> On utilise des **bases de données publiques** avec des images et le texte correspondant

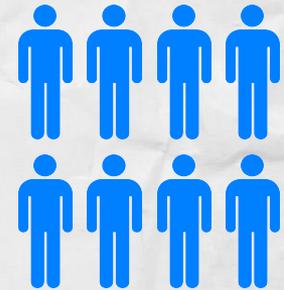
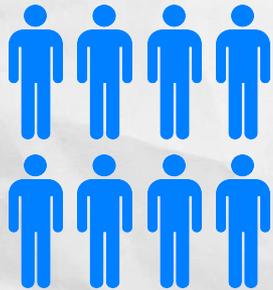
-> Il y a aussi des compétitions !



On **compte le nombre d'erreur** faites par l'ordinateur

LA COLLABORATION AVEC LES ETUDIANTS ALLEMANDS

On échange ...



On se **rencontre** régulièrement pour **présenter** nos résultats, **échanger** nos idées, **parler** des expériences qu'on pourrait faire...

Merci de votre attention !

Théodore Bluche
tb@azia.com