

LU3I026 - Science des données

Des données à l'apprentissage

Olivier Schwander <olivier.schwander@sorbonne-universite.fr>

Sorbonne Université

2025-2026

Formaliser le problème d'IA

Discussion informelle

- ▶ «On veut de l'IA»
- ▶ «sur nos données»
- ▶ «pour faire ci ou ça»
- ▶ «avec telle et telle contraintes»

Objectif

- ▶ Traduction dans un vocabulaire IA
- ▶ Correspondance avec les concepts scientifiques

Première question à poser

- ▶ «Cher.e chef.fe , avez-vous des données ?»
- ▶ Parfois, la réponse est non...
- ▶ Or c'est un besoin capital !

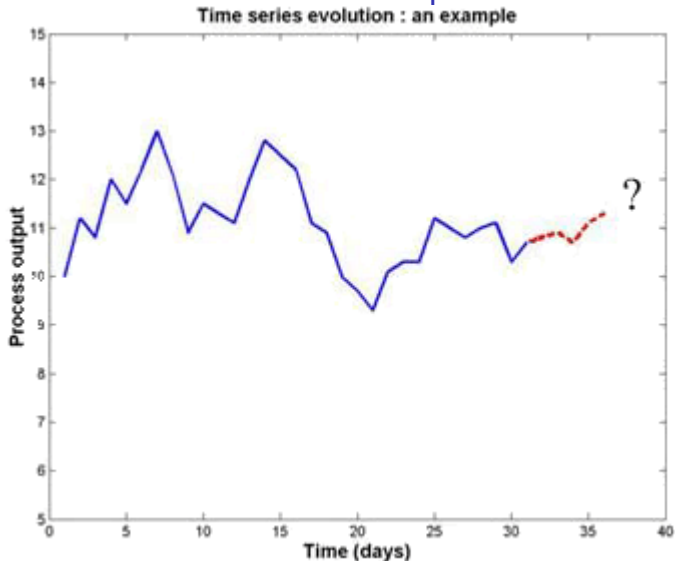
Différents données: texte

Extrait de wikipedia

ChatGPT est un [prototype d'agent conversationnel](#) utilisant l'[intelligence artificielle](#), développé par [OpenAI](#), et spécialisé dans le dialogue. L'agent conversationnel est un [modèle de langage](#) affiné par [apprentissage supervisé](#) et par [apprentissage par renforcement](#).

Le sigle *ChatGPT* signifie « *Generative pre-trained transformer* » (« *Transformateur générique pré-entraîné* »).

Différentes données: séries temporelles



Différentes données: autres

Graphes

- ▶ Réseaux sociaux
- ▶ Échanges téléphoniques, sms
- ▶ Co-auteurs de publications scientifiques

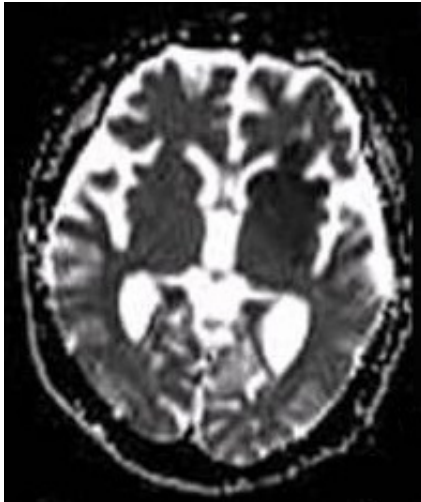
Vidéos

- ▶ Image + aspect temporel

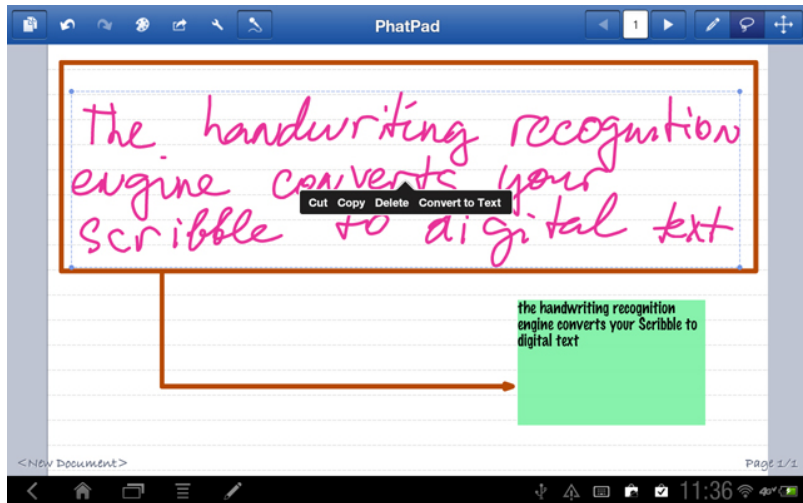
Conversations

- ▶ Échanges de type question-réponse (ChatGPT)

Différents domaines: image médicale



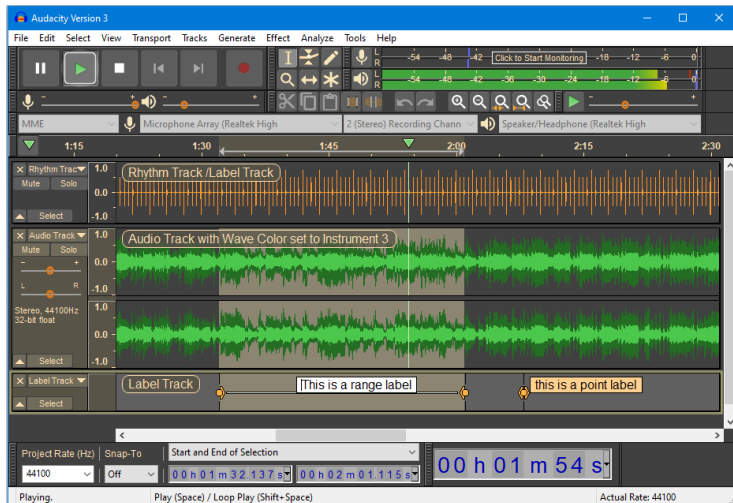
Différents domaines: écriture manuscrite



Différentes domaines: images satellites



Différentes domaines: audio



Différents domaines

Voiture autonome: vision

- ▶ Été – hiver
- ▶ Ville – campagne
- ▶ Amérique du nord – Europe

Contrôle

- ▶ Voiture autonome – drone – autre robot

Tâches de machine learning

Choix à faire

- ▶ Formalisation
- ▶ Définir les objectifs
- ▶ Définir le protocole expérimental

Exemple

- ▶ *Classification*
- ▶ *binaire*
- ▶ évaluation avec *accuracy*
- ▶ *découpage train/test* fixé de telle façon

Données

Matériau de base

- ▶ **Incourtournable**
- ▶ (sauf cas particulier)
- ▶ Vérifier: quantité, qualité

Avant de se lancer avec une méthode IA

- ▶ Charger les données
- ▶ Compréhension minimale
- ▶ Réfléchir aux biais
- ▶ Statistiques simples, visualisations
- ▶ Réfléchir aux prétraitements

Sources de données

Monde académique

- ▶ Données jouets
- ▶ Données réelles
- ▶ Souvent bien connu et maîtrisés: stockage propre, exemples d'utilisation, bibliothèques

Monde réel

- ▶ Base de données
- ▶ Fichiers structurés: CSV, JSON, YAML, tableurs, autres formats
- ▶ Web: services web, API, scrapping

Exemple du CSV

1	region	total_pop	percent_white	percent_black	percent_hispanic	per_capita_income	med_hh_income	med_household_size
2	Alabama	4709277	61	12	8	23680	361	2.61
3	Alaska	726226	64	3	5	10052	570	2.16
4	Arizona	6479783	57	4	10	25558	747	2.63
5	Arkansas	2933369	74	15	1	22170	480	2.13
6	California	3765181	57	4	13	36	26537	3.19
7	Colorado	5119329	78	4	3	21	31189	2.52
8	Connecticut	3583561	76	8	14	19862	480	2.62
9	Delaware	989846	65	21	8	29818	638	2.69
10	District of Columbia	653371	55	2	12	45296	554	3.16
11	Florida	19891356	57	15	2	25	36296	2.68
12	Georgia	10018154	62	12	8	21842	470	2.69
13	Hawaii	1274298	23	1	17	9	28989	3.22
14	Idaho	1581862	84	1	1	11	31568	2.67
15	Illinois	12848554	63	14	5	14	29848	2.79
16	Indiana	6549862	81	9	2	24809	597	2.12
17	Iowa	3046253	88	3	2	5	17027	2.14
18	Kansas	3699187	76	6	2	11	28929	2.51
19	Kentucky	4361333	86	8	1	8	19462	2.68
20	Louisiana	4602553	60	22	3	27827	534	2.61
21	Maine	1328328	94	1	1	1	28824	2.64
22	Maryland	5843795	54	29	6	6	36354	3.04
23	Massachusetts	6860269	76	6	6	18	25763	3.26
24	Michigan	9898905	76	14	3	5	25883	3.23
25	Minnesota	5547746	83	5	5	38913	274	2.78
26	Mississippi	2976672	58	27	3	28048	518	2.62
27	Missouri	5887182	81	11	2	4	25849	2.49
28	Montana	1005564	87	1	1	3	23273	2.77
29	Nebraska	1841825	82	4	2	9	26899	2.63
30	Nevada	2798865	53	8	7	27	26589	2.64
31	New Hampshire	1321211	93	1	2	3	33334	2.76
32	New Jersey	8852486	59	13	9	18	36827	3.02
33	New Mexico	2069786	68	2	1	47	22763	2.65
34	New York	19479538	58	14	8	18	32382	3.13

```

1 Demographics_State.csv
2 "region","total_population","percent_white","percent_black","percent_hispanic","per_capita_income","me
3 "Alabama",4709277,61,12,8,23680,361
4 "Alaska",726226,64,3,5,10052,570
5 "Arizona",6479783,57,4,10,25558,747
6 "Arkansas",2933369,74,15,1,22170,480
7 "California",3765181,57,4,13,36,26537,3.19
8 "Colorado",5119329,78,4,3,21,31189,2.52
9 "Connecticut",3583561,76,8,14,19862,480
10 "District of Columbia",653371,55,2,12,45296,554
11 "Florida",19891356,57,15,2,25,36296,2.68
12 "Georgia",10018154,62,12,8,21842,470
13 "Hawaii",1274298,23,1,17,9,28989,3.22
14 "Idaho",1581862,84,1,1,11,31568,2.67
15 "Illinois",12848554,63,14,5,14,29848,2.79
16 "Indiana",6549862,81,9,2,24809,597
17 "Iowa",3046253,88,3,2,5,17027,2.14
18 "Kansas",3699187,76,6,2,11,28929,2.51
19 "Kentucky",4361333,86,8,1,8,19462,2.68
20 "Louisiana",4602553,60,22,3,27827,534
21 "Maine",1328328,94,1,1,1,28824,2.64
22 "Maryland",5843795,54,29,6,6,36354,3.04
23 "Massachusetts",6860269,76,6,6,18,25763,3.26
24 "Michigan",9898905,76,14,3,5,25883,3.23
25 "Minnesota",5547746,83,5,5,38913,274
26 "Mississippi",2976672,58,27,3,28048,518
27 "Missouri",5887182,81,11,2,4,25849,2.49
28 "Montana",1005564,87,1,1,3,23273,2.77
29 "Nebraska",1841825,82,4,2,9,26899,2.63
30 "Nevada",2798865,53,8,7,27,26589,2.64
31 "New Hampshire",1321211,93,1,2,3,33334,2.76
32 "New Jersey",8852486,59,13,9,18,36827,3.02
33 "New Mexico",2069786,68,2,1,47,22763,2.65
34 "New York",19479538,58,14,8,18,32382,3.13
    
```

Un fichier texte: des lignes

- ▶ /valeur1 VIRGULE Valeur2 VIRGULE etc
- ▶ /valeur1 VIRGULE Valeur2 VIRGULE etc

Quelques outils

Éviter Python pur

Pandas – Polars

- ▶ Réflexe à avoir: chargement, filtrage, statistiques, visualisation
- ▶ Point de vue base de données
- ▶ `read_csv` (et plein d'autres)

Numpy – Scipy

- ▶ Implémentations d'algorithmes
- ▶ Calcul scientifique
- ▶ `loadtxt`

Requests – HTTPX – BeautifulSoup – Playwright

- ▶ HTTP, HTML, Javascript

Regarder: fast-checking

- ▶ Dimensions des structures de données
- ▶ En-têtes des colonnes
- ▶ print tous simples

	region	total_population	percent_white	percent_black	percent_
0	alabama	4799277	67	26	
1	alaska	720316	63	3	
2	arizona	6479703	57	4	
3	arkansas	2933369	74	15	
4	california	37659181	40	6	

	percent_hispanic	per_capita_income	median_rent	median_age
0	4	23680	501	38.1
1	6	32651	978	33.6
2	30	25358	747	36.3
3	7	22170	480	37.5
4	38	29527	1119	35.4

Indicateurs statistiques

Outils mathématiques simples

- Moyennes, médianes, écart-types, percentiles

Outils informatique

Pandas describe (attention: pour vous ! pas pour un rapport propre !)

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
count	364	364	364.000000	364	364	364	364.000000	364	3.640000e+02
unique	364	30	NaN	5	22	17	NaN	115	NaN
top	Cleanthony Early	New Orleans Pelicans	NaN	SG	24.0	6-9	NaN	Kentucky	NaN
freq	1	16	NaN	87	41	49	NaN	22	NaN
mean	NaN	NaN	16.829670	NaN	NaN	NaN	219.785714	NaN	4.620311e+06
std	NaN	NaN	14.994162	NaN	NaN	NaN	24.793099	NaN	5.119716e+06
min	NaN	NaN	0.000000	NaN	NaN	NaN	161.000000	NaN	5.572200e+04
20%	NaN	NaN	4.000000	NaN	NaN	NaN	195.000000	NaN	9.472760e+05
40%	NaN	NaN	9.000000	NaN	NaN	NaN	212.000000	NaN	1.638754e+06
50%	NaN	NaN	12.000000	NaN	NaN	NaN	220.000000	NaN	2.515440e+06
60%	NaN	NaN	17.000000	NaN	NaN	NaN	228.000000	NaN	2.429924e+06

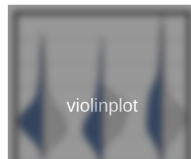
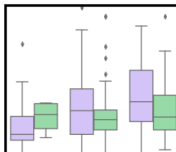
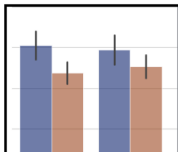
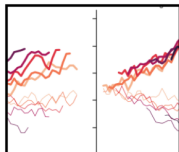
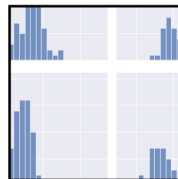
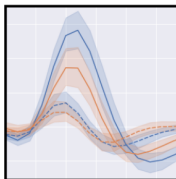
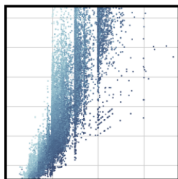
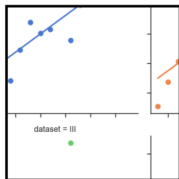
Visualisation

Graphiques

- ▶ Histogrammes, courbes, nuages de points, etc

Outils informatique

- ▶ Matplotlib, Bokeh, Seaborn, Altair



Autres types de données

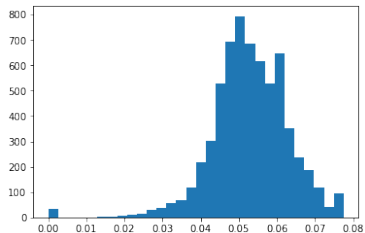
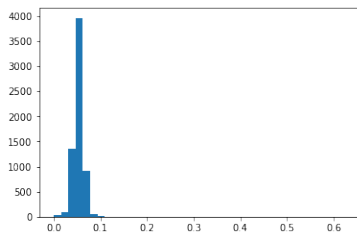
Visualisation adaptée !

- ▶ Images: afficher des images
- ▶ Texte: le texte...
- ▶ Son: écouter
- ▶ etc

Oui mais: retour sur la qualité des données

Valeurs aberrantes (outliers)

Histogramme des prix au km des courses Blablacar



(bonus: dire du mal de cette figure)

Comment faire ?

- ▶ À la main: pas évident
- ▶ Automatiquement: difficile

Biais éventuels

Problèmes de société

- ▶ Des images de femmes dans le dataset ?
- ▶ Toutes les couleurs de peau ?
- ▶ Uniquement de l'anglais ?

Rien d'objectif

- ▶ Un dataset est *situé* (au sens de la sociologie)
- ▶ Discriminations systémiques: racisme, sexisme, validisme, etc

Conséquences

- ▶ Technique: un modèle qui marche moins bien (parce que pas réaliste)
- ▶ Sociale: un modèle **complice** des discriminations

Il faut regarder dans le dataset !!!

Contextes d'apprentissage

Apprentissage supervisé

- ▶ Données
- ▶ Étiquettes
- ▶ Exemple: des images, et le nom de l'objet dans l'image
- ▶ Difficulté: obtenir les étiquettes...

Apprentissage non-supervisé

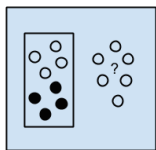
- ▶ Juste les données
- ▶ Pas d'étiquettes

Autres

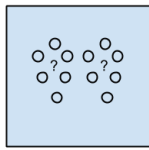
- ▶ Semi-supervisé
- ▶ Renforcement
- ▶ Auto-supervisé

Contextes d'apprentissage

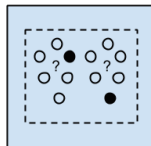
Supervisé Non-supervisé Semi-supervisé Renforcement



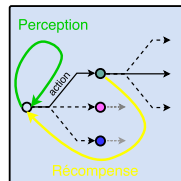
Supervised Learning Algorithms



Unsupervised Learning Algorithms



Semi-supervised Learning Algorithms



Première chose à déterminer

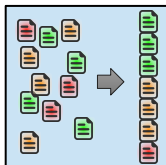
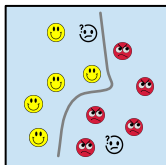
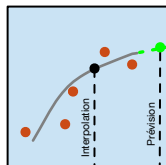
- ▶ conditionne les choix d'algorithmes et de méthodes
- ▶ et les mesures de performances

Tâches d'apprentissage

Régression

Classification

Ordonnancement



Autres tâches

- ▶ Souvent des cas particuliers des tâches précédentes
- ▶ Image: segmentation (classif), localisation (régression)
- ▶ Texte: génération (classif), dialogue (ordonnancement)

Deuxième chose à déterminer

- ▶ Conditionne les choix précis de modèles
- ▶ Permet de fixer le protocole expérimental
- ▶ Ensuite, on peut réfléchir aux modèles

Premier modèle: linéaire

Dépendance linéaire entre l'entrée et la sortie

- ▶ Dans un espace vectoriel de dimension d , soit un point $x \in \mathbb{R}^d$
- ▶ $f_w(x) = \sum_{i=1}^d w_i x_i$
- ▶ f est le *modèle*
- ▶ w est le vecteur de *paramètres*

Apprentissage

- ▶ À partir des données
- ▶ Trouver le meilleur w

Abus de langage

- ▶ Linéaire: $y = ax$
- ▶ Affine: $y = ax + b$
- ▶ **Attention** en TP !

Régression linéaire

- ▶ Avec **UN** point: $\mathbf{x} = (x_1, \dots, x_d)$
- ▶ Une étiquette $y \in \mathbb{R}$
- ▶ Approcher l'étiquette y avec un modèle $f(x)$

Apprentissage supervisé

- ▶ À partir de N points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
- ▶ Et leur label y_1, \dots, y_N
- ▶ Chercher le meilleur compromis

Classification linéaire

- ▶ Avec **UN** point: $\mathbf{x} = (x_1, \dots, x_d)$
- ▶ Une étiquette $y \in \{-1, +1\}$
- ▶ Retrouver l'étiquette y avec un modèle $\text{signe}(f(\mathbf{x}))$

avec $\text{signe}(x) = +1$ si $x \leq 0$ et -1 sinon

Apprentissage supervisé

- ▶ À partir de N points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
- ▶ Et leur label y_1, \dots, y_N
- ▶ Chercher le meilleur compromis

Géométriquement

- ▶ Tracer les points
- ▶ Tracer la frontière de décision
- ▶ À quoi correspond la somme ?

Plus proche voisin

Nearest neighbor

Idée

Associer un nouveau point à l'exemple connu le plus proche
[Au tableau]

k -plus proches voisins

Algorithme

- ▶ Pour un nouveau point
- ▶ Calculer toutes les distances entre ce point et les points connus
- ▶ Rechercher les k plus proches voisins
- ▶ Aggréger: vote (classification), moyenne (régression)

[Au tableau]

Analyse

Complexité algorithmique

- ▶ Temps, espace

Performances

- ▶ Influence du k
- ▶ Score et généralisation

[Au tableau]

Prétraitements

[Au tableau]

Conclusion méthodologique

Contexte

- ▶ Apprentissage supervisé ou non
- ▶ Tâche de classification ou de régression

Modèles

- ▶ Linéaire
- ▶ k -Plus proches voisins

Prétraitements

- ▶ Données catégorielles
- ▶ Normalisation
- ▶ Transformation des données

Conclusion mathématique

Outils statistiques

- ▶ Moyennes, écart-type
- ▶ Distribution gaussienne
- ▶ Corrélation

Outils géométriques

- ▶ Droites, hyper-plans, surfaces
- ▶ Distances, produits scalaires

Conclusion informatique

Pandas

- ▶ Chargement des données
- ▶ Analyse préliminaire
- ▶ Statistique descriptive

Numpy

- ▶ Implémentation des algorithmes
- ▶ Bas-niveau
- ▶ Calcul matriciel **efficace**
- ▶ **JAMAIS** de boucle for de Python