

Machine Learning – Outils pour l'enseignement

Ricco Rakotomalala

Université Lumière Lyon 2 – Data, informatique et statistique

<http://dis.univ-lyon2.fr/>

- Formation en économétrie (statistique, économie mathématique)
- Thèse de doctorat en Machine Learning ([Apprentissage statistique](#))
- Enseignant chercheur, en poste à l'Université Lumière Lyon 2
- Spécialité : statistique et informatique, data mining et ses applications - [Data Science](#)
- « Père » des logiciels gratuits [SIPINA v.3](#) et [TANAGRA](#) (open source)
- Auteur d'une dizaine d'[ouvrages libres](#)
- Auteur de **plus de 500** [supports de cours](#), de tutoriels en [français](#) et en [anglais](#)
- [630 sessions par jour](#) sur plus de 10 années (depuis 01.02.2008 - Compteur Google Analytics)
- Responsable du Master [SISE](#) (Statistique et Informatique pour la Science des donnÉEs)

Interrogations...

Peut-on « se contenter » de logiciels libres pour l'enseignement du machine learning dans nos formations ?

Si oui, lesquels de logiciels libres ? Et pourquoi ?

Pédagogie

Praticabilité dans les enseignements

Richesse fonctionnelle

Insertion professionnelle des étudiants

R et Python

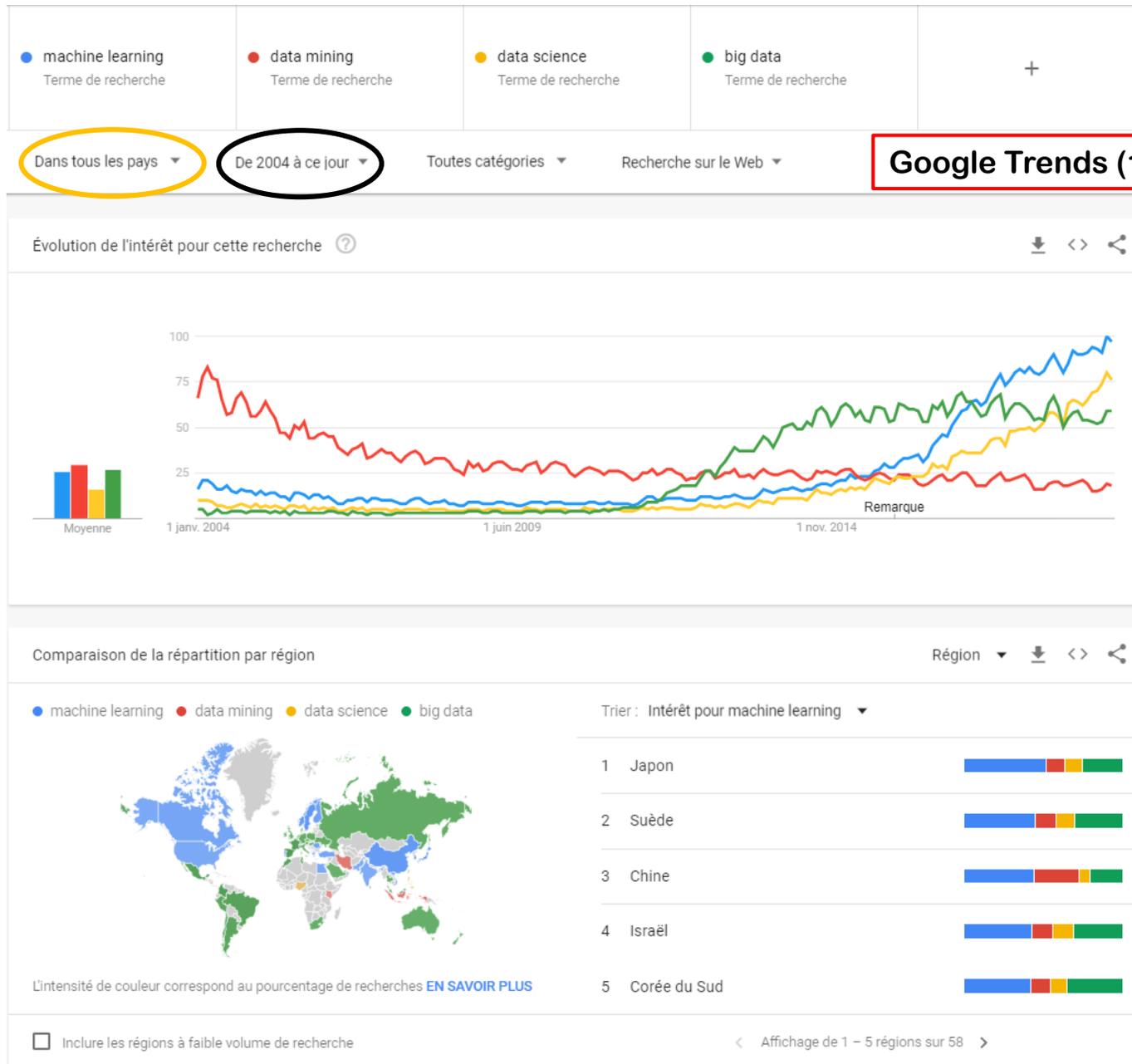
Plan

1. Machine Learning, Data Mining, Data Science...
2. Logiciels – Attentes pédagogiques
3. Réalisations étudiantes sous R et Python
4. Conclusion

Machine Learning

Faire du neuf avec du vieux ? Big Data Analytics

Machine Learning – Regain de popularité récent



« **Machine Learning** »
(apprentissage automatique, apprentissage machine, apprentissage statistique),
branche de l'intelligence artificielle, qui trouve ses
origines **dans les années 40 - 50**
([Wikipédia](#)).

Machine Learning vs. Statistique

The screenshot shows the Coursera website for the 'Apprentissage automatique' course. The page is divided into two weeks. Week 1 includes 'Introduction', 'Linear Regression with One Variable', and 'Linear Algebra Review'. Week 2 includes 'Linear Regression with Multiple Variables'. The course has a 4.9 rating and 28,536 reviews. A red dashed box highlights the Coursera logo and the search bar.

Apprentissage automatique
★★★★★ 4.9 116 277 notes • 28 536 avis

S'inscrire gratuitement
Commence le oct. 14

Alde financière disponible

2 574 065 déjà inscrits !

À propos Programme de cours Avis Enseignants

SEMAINE 1

2 heures pour terminer

Introduction

Welcome to Machine Learning! In this module, we introduce the core idea of teaching a computer to learn concepts using data—without being explicitly programmed. The Course Wiki is under construction. Please visit the resources tab for the most complete and up-to-date information.

5 vidéos (Total 42 min), 9 lectures, 1 quiz **VOIR TOUT**

2 heures pour terminer

Linear Regression with One Variable

Linear regression predicts a real-valued output based on an input value. We discuss the application of linear regression to housing price prediction, present the notion of a cost function, and introduce the gradient descent method for learning.

7 vidéos (Total 70 min), 8 lectures, 1 quiz **VOIR TOUT**

2 heures pour terminer

Linear Algebra Review

This optional module provides a refresher on linear algebra concepts. Basic understanding of linear algebra is necessary for the rest of the course, especially as we begin to cover models with multiple variables.

6 vidéos (Total 61 min), 7 lectures, 1 quiz **VOIR TOUT**

SEMAINE 2

3 heures pour terminer

Linear Regression with Multiple Variables

What if your input has more than one value? In this module, we show how linear regression can be extended to accommodate multiple input features. We also discuss best practices for implementing linear regression.

8 vidéos (Total 65 min), 16 lectures, 1 quiz **VOIR TOUT**

« On aimerait faire de l'IA et du Machine Learning... où il est question de convolutions et de gradient... avec du traitement d'images... »

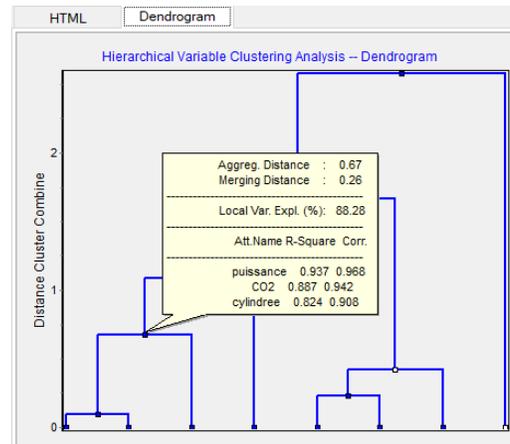
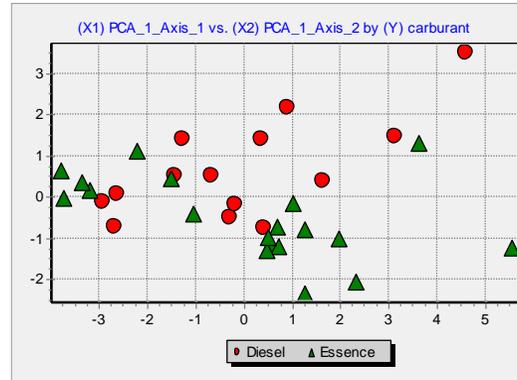
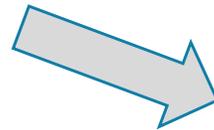
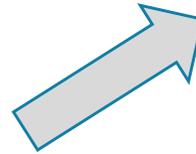
Dépasser les chapelles de naguère : statistique, économétrie, analyse de données, apprentissage automatique....

Statistique – Statistique exploratoire

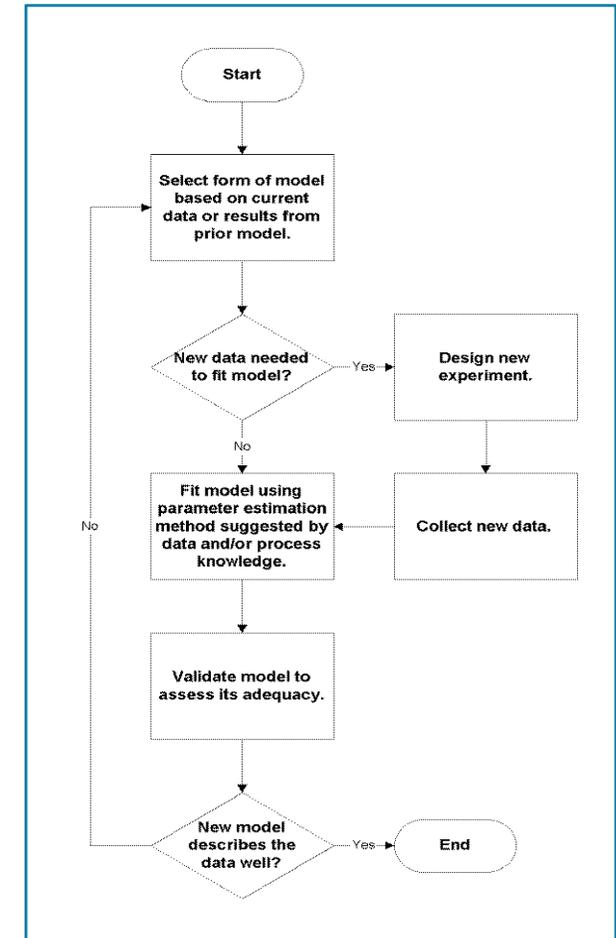


Les données sont spécifiquement recueillies à des fins d'étude (ex. enquête, expérimentations, etc.)

- Bonne qualité souvent
- Faible volumétrie (rareté)



Application des techniques de modélisation et de statistique

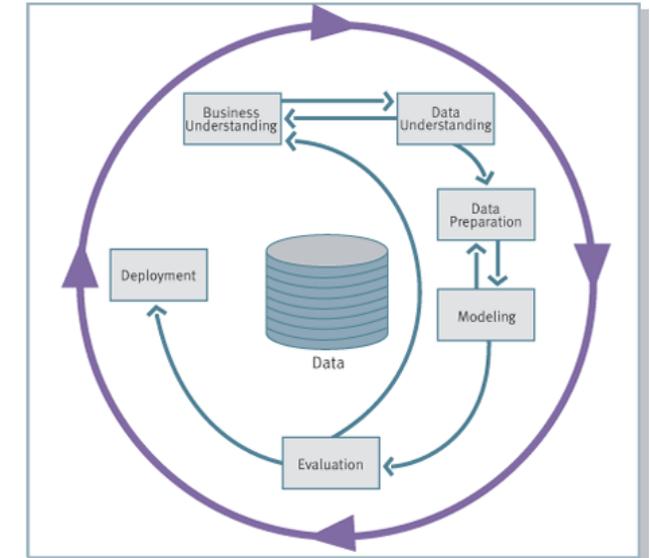
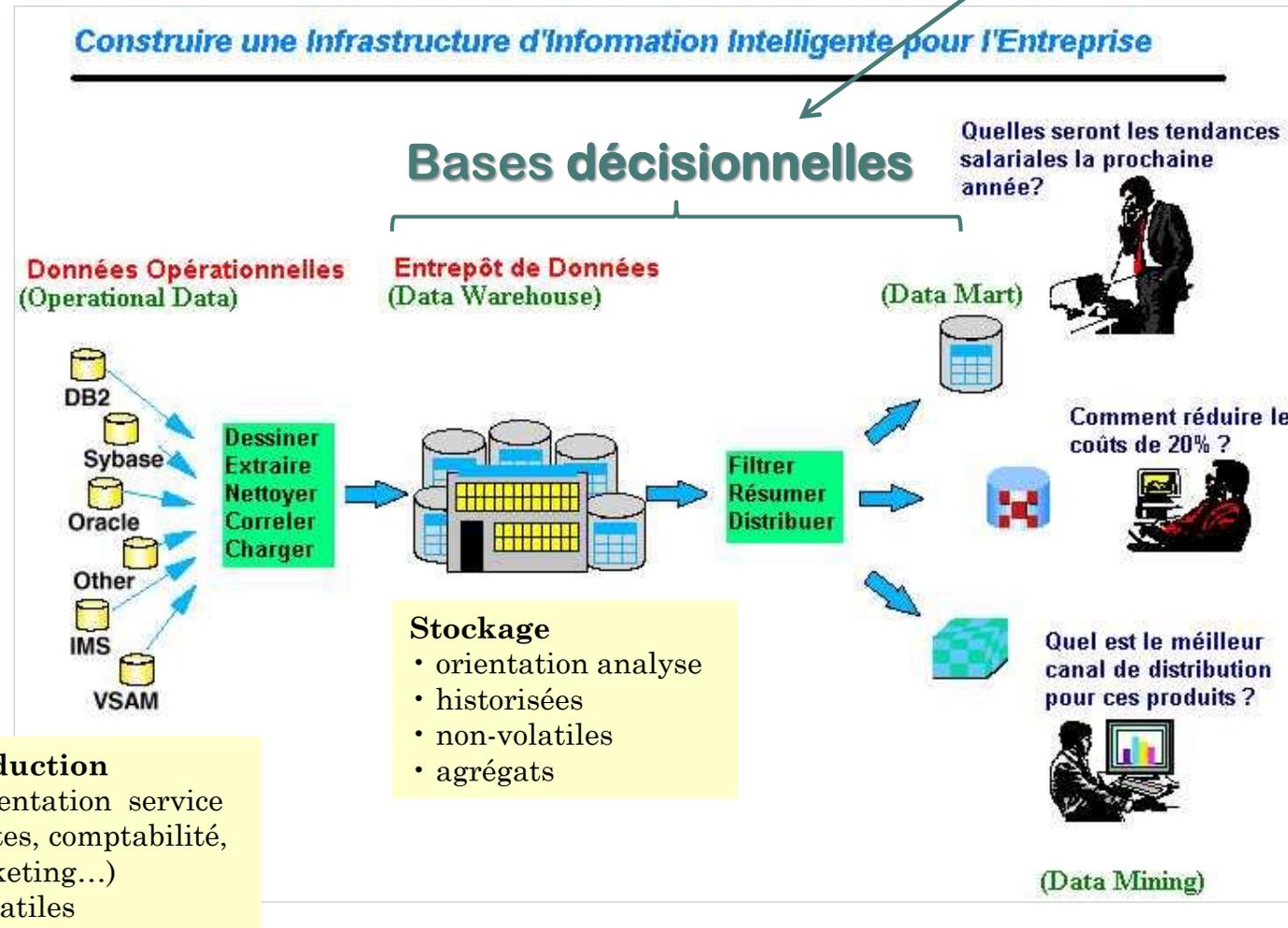


NIST – e-Handbook of Statistical Methods

Volume de traitements – de toute manière – limité par les capacités des outils informatiques disponibles (à l'époque). !

Vague « Data Mining » à partir de la fin des années 90

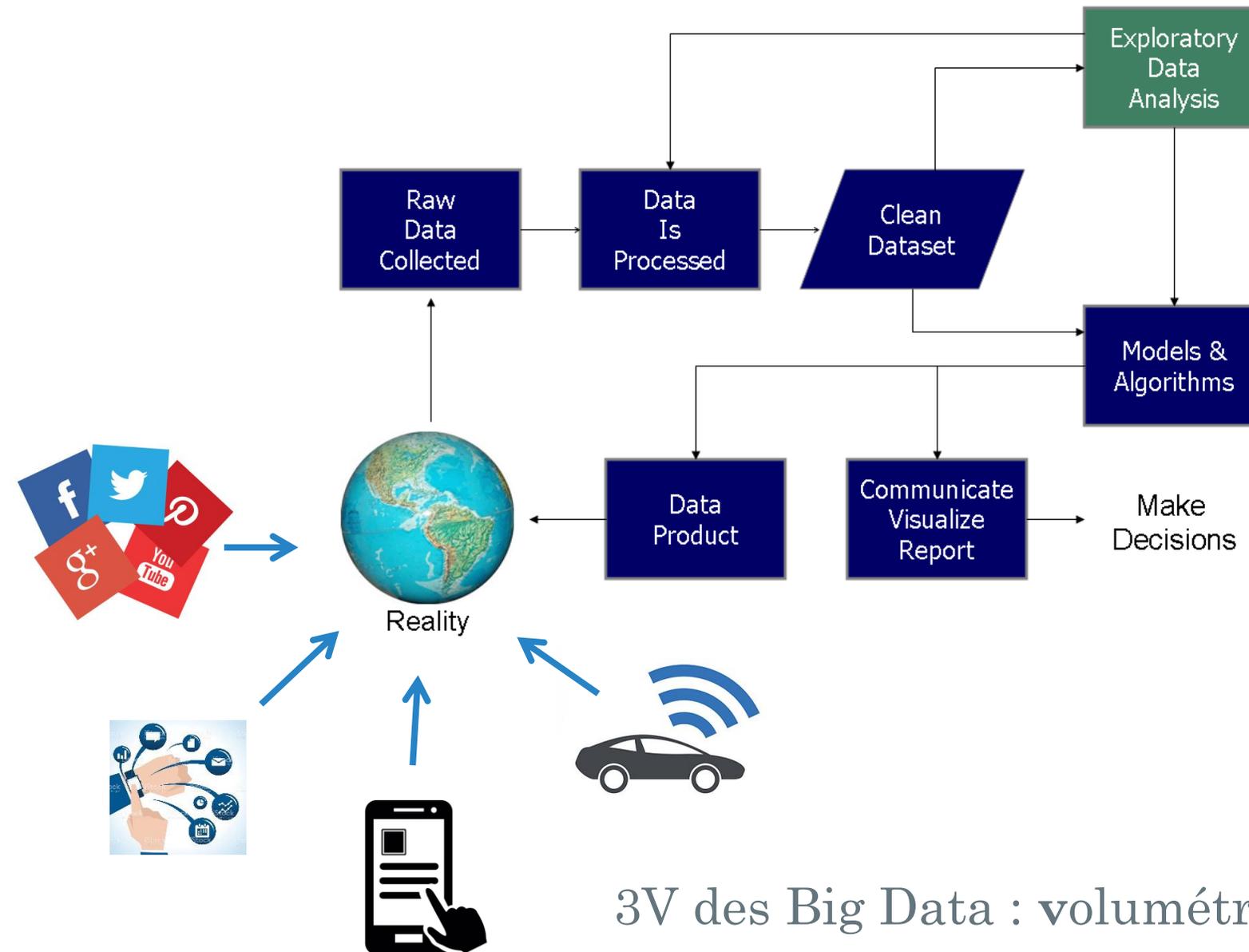
Les données sont organisées et stockées de manière à ce que nous puissions mener des analyses.



Les sources d'information et les technologies évoluent. Elément clé : l'entrepôt de données.

Vague « Big Data Analytics » (Data Science) actuelle

Data Science Process



Sources additionnelles d'information **externes** à l'entreprise, multiplicité des formats, **nouveaux enjeux technologiques** pour le stockage et le traitement (stockage NoSQL, data lake, tech. cloud, info. distribuée [Hadoop, Spark]...) et **nouvelles opportunités d'analyse** ! (cf. Rapports Lauvergeon, Villani)

3V des Big Data : volumétrie, variété, vélocité.

Exemple. Filtrage collaboratif, Systèmes de recommandation.

amazon.fr Premium
Toutes nos boutiques | gil jourdan | Les séries Amazon Original avec Amazon Premium
Commencez votre essai gratuit de 30 jours

Parcourir les boutiques | Chez vous | Ventes Flash | Chèques-cadeaux | Bonjour, Identifiez-vous | Testez Premium | Vos Listes | Panier

Livres | Recherche détaillée | Nos rubriques | Livres de l'hiver | Meilleures ventes | Nouveautés | Précommandes | Livres Poche | Livres anglais et étrangers

Retour aux résultats de la recherche pour « gil jourdan »

Gil Jourdan : L'Intégrale 1 Album – 5 juin 2009
de Maurice Tillieux (Auteur)
★★★★★ 9 commentaires client

Album
EUR 24,00

9 d'occasion à partir de EUR 12,00
9 neufs à partir de EUR 24,00

Voulez-vous le faire livrer plus vite ?
15 h et 5 mins et c'est gratuit de votre commande

Note: Cet article est éligible à la livraison gratuite Amazon Prime

Partager

EUR 24,00
Tous les prix incluent la TVA.

Livraison à EUR 0,01 en France métropolitaine.

Produits fréquemment achetés ensemble
Prix total: EUR 72,00
Ajouter ces trois articles au panier

- Cet article : Gil Jourdan : L'Intégrale 1 par Maurice Tillieux Album EUR 24,00
- Gil Jourdan - L'Intégrale - tome 2 - Gil Jourdan 2 (intégrale) 1960 - 1963 par Maurice Tillieux Album EUR 24,00
- Gil Jourdan : L'Intégrale 3 par Maurice Tillieux Relié EUR 24,00

Recommandation basée sur les transactions.

Recommandation basée sur les utilisateurs (clients).

Les clients ayant acheté cet article ont également acheté

Page 1 sur 19

- Gil Jourdan - L'Intégrale - tome 2 - Gil Jourdan 2 (intégrale) 1960 - 1963 par Maurice Tillieux Album EUR 24,00 Premium
- Gil Jourdan : L'Intégrale 3 Maurice Tillieux Relié EUR 24,00 Premium
- Gil Jourdan - L'Intégrale - tome 4 - Gil Jourdan 4 (intégrale) 1970 - 1979 par Maurice Tillieux Album EUR 24,00 Premium
- Johan et Pirlouit - L'Intégrale - tome 1 - Johan et Pirlouit intégrale 1... Peyo Album EUR 20,50 Premium
- Johan et Pirlouit - L'Intégrale - tome 2 - Johan et Pirlouit intégrale 2 réédition Peyo Album EUR 24,00 Premium

Evaluations des produits
Commentaires des clients

Attentes vis-à-vis des logiciels de data science / machine learning

Qu'attendre aujourd'hui des logiciels pour l'enseignement ?

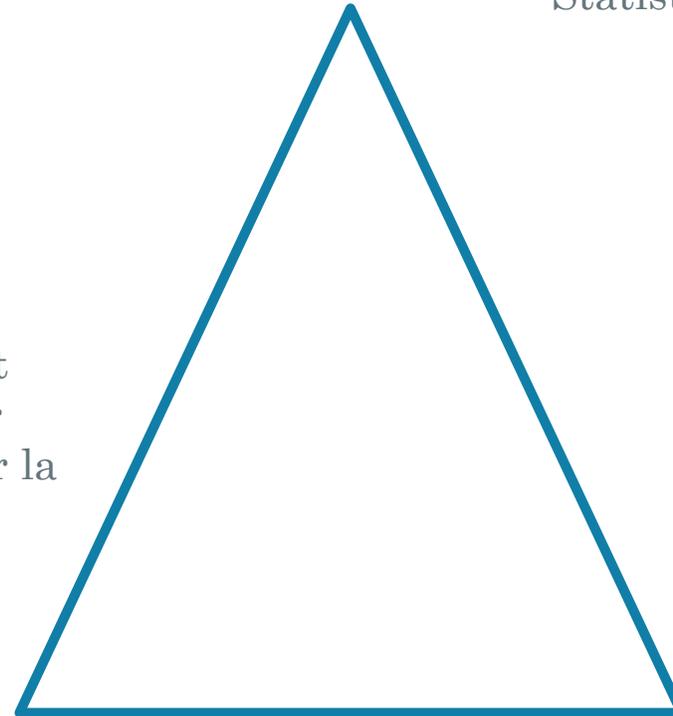
Enseignement de la data science

**STATISTIQUE
DATA MINING**

Connaître et comprendre les techniques de modélisation, d'analyse de données, d'inférence... savoir exploiter les régularités « cachées » dans les données, pourvoyeuses de connaissances.
Statistique, Data mining, Machine Learning.

Maîtriser les outils pour accéder et manipuler les données, développer des stratégies nouvelles pour gérer la profusion de l'information,...
Technologies big data

INFORMATIQUE



Toute analyse s'inscrit dans un domaine d'application : données de sécurité, données du web, analyse des réseaux sociaux, etc.

APPLICATIONS

Le logiciel joue un rôle très important



Critères pour les logiciels de data science

1. Architecture : stand-alone, client-serveur, via un navigateur, ...
2. Mode opératoire : diagramme de traitements, langage de script, pilotage par menu,...
3. Performances, capacités de traitement, temps de calcul
4. Accès aux données : fichiers textes, Excel, accès aux bases de données,...
5. Solutions pour la volumétrie, **technologies big data** (hadoop, spark,...)
6. Accès aux **données non structurées** et primitives de traitements (texte, image, ...)
7. Interfaçage avec les **API du web** (ex. Twitter, Google Analytics, OpenStreetMap, ...)
8. Manipulation des données : transformations, recodage,...
9. Exploration graphique : représentations, visualisations, interactions,...
10. **Bibliothèques de techniques de machine learning** : supervisées, ..., **deep learning**,...
11. Evaluation et comparaisons : comparaison des approches, benchmarking...
12. Reporting et solutions pour le **déploiement** (PMML,...)

L'utilisation du logiciel doit s'inscrire dans une démarche pédagogique. Deux extrêmes à éviter :
Qu'est-ce qu'il faut faire là ? vs. On recopie sans rien comprendre... (cf. Fiches de TD en Master)

Objectifs d'une séance de TD sur machine :

- Mettre en œuvre une technique de machine learning vue en cours c.-à-d.
 - Charger les données
 - Lancer l'algorithme en respectant un schéma prédéfini (ex. app-test en prédictif)
 - Savoir lire et inspecter / expertiser les résultats
 - Déployer les modèles
- Plus loin : apprendre à manipuler les paramètres des algorithmes
- Plus loin encore : s'intéresser aux problèmes pratiques (recodage, données manquantes, comparaison des approches, solutions pour volumétrie, etc.)

Nous formons des étudiants qui vont en entreprise. Il ne faut pas que nos choix les impactent négativement en les emmenant sur des voies de garage.



Une utilisation **conforme aux standards du domaine** et répondant aux objectifs pédagogiques :

- S'attacher **au fond et non la forme** (la base : charger les données [aux formats usuels], appliquer les méthodes, restituer les résultats. Cf. **les objectifs d'une séance de TD**).
- **Pas de mode opératoire spécifique**, nécessitant des compétences supplémentaires pouvant parasiter le discours (Programmation vs. Diagramme de traitement).
- Démarche harmonisée pour un **large spectre d'applications** (même environnement pour le machine learning « classique », le text mining, image mining, etc.). Ne pas multiplier les outils en fonction de la matière étudiée. Intérêt des dispositifs à packages.

Question ancienne : opportunité d'utiliser des logiciels libres pour l'enseignement du data mining ([Déc. 2005](#)).

OUI pour les aspects méthodologiques, MAIS attention aux aspects opérationnels (ex. reporting, déploiement)

[WEKA, ORANGE ML, TANAGRA], et *pas vocation à être utilisés en entreprise*.

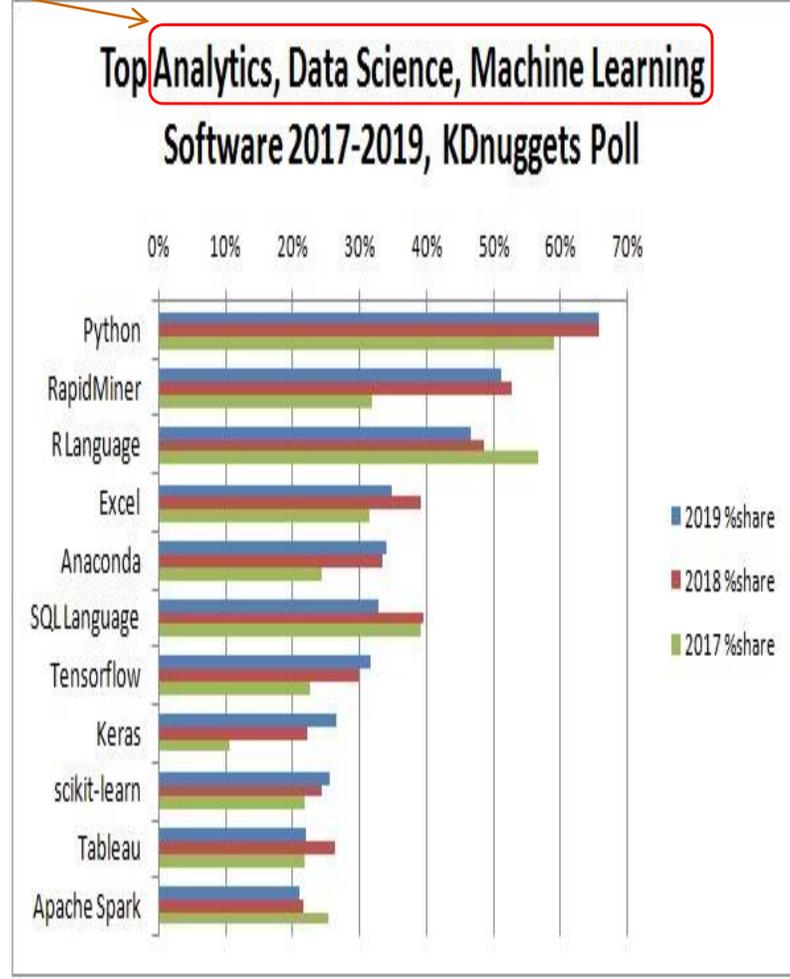
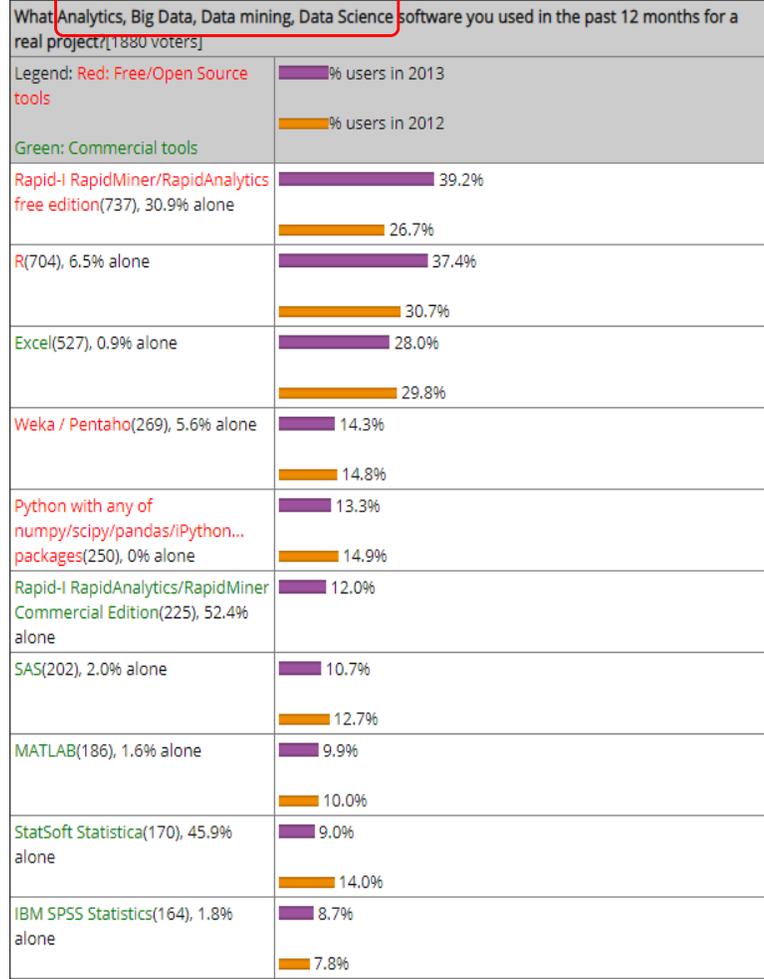
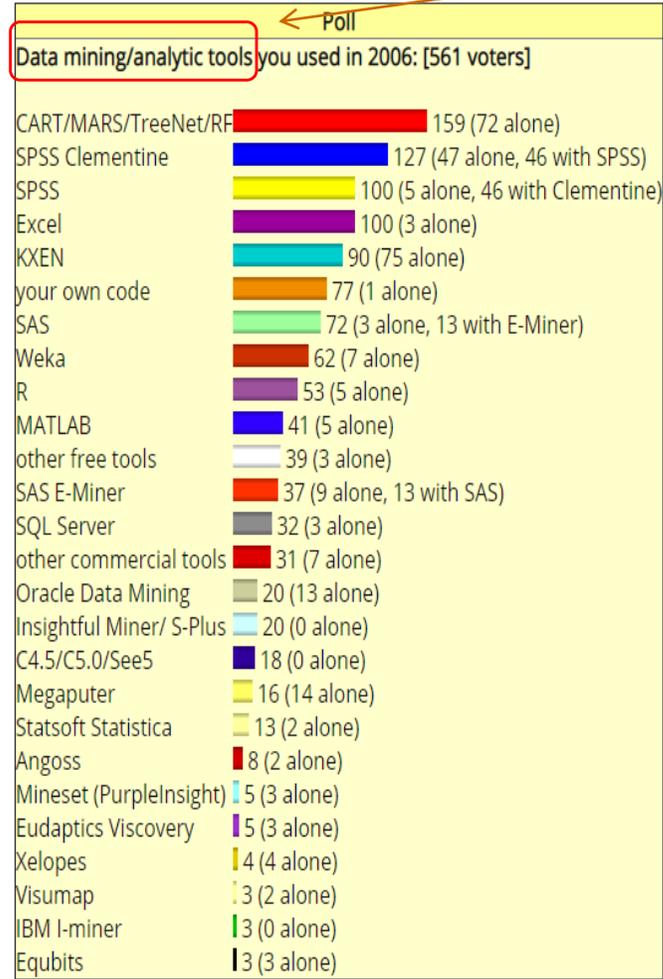
Aujourd'hui, **R** et **Python** s'imposent dans nos formations parce que :

- **Gratuits**, totalement, ce n'est pas moindre de leurs qualités.
- Forte **pénétration de ces outils dans les entreprises** (cf. les offres d'emploi sur le site de l'APEC, avec les mots-clés « statistique », « data science », etc.).
- Le **niveau en programmation** des étudiants a considérablement évolué. Et on peut s'appuyer sur ces outils pour l'améliorer encore.
- **Outils polyvalents**. Très large spectre d'utilisation.



Compétence sur ces outils devient un marqueur fort dans les CV.

Les appellations changent au fil des années...



Il y a matière à réflexion quand on regarde cette évolution



Constat à ce jour :

- L'intérêt exacerbé pour Python (cf. mon site des tutoriels) ----->
- Dans l'esprit des entreprises en France : Machine Learning = Python (cf. démo de l'étude sous R-Shiny, les stages encadrés, enfin les offres d'emploi sur l'APEC)
- La « lenteur » de R est souvent décriée (mais est-ce vraiment justifié...)

Ce que je pense :

- Selon le thème, le contexte, les objectifs... et les packages disponibles, on peut préférer l'un ou l'autre
- On peut facilement développer une compétence conjointe forte pour R et Python

UE Informatique appliquée (9 ECTS)

- Programmation Statistique sous R (R. Rakotomalala)

Apprentissage de la programmation sous R. Structures avancées. Programmation des algorithmes de statistique et de data mining sous R. Modèle objet sous R.

Programmation big data (map reduce) sous hadoop. Programmation R sous spark.

Création de packages.

M2 SISE 2019-2020

- Machine Learning sous Python (N. Sawadogo)

Bases de la programmation python, structures vectorielles et matricielles. Algorithmes de machine learning d'apprentissage supervisé et non supervisé (svm - support vector machine, dbscan, birch,...). Image mining, traitement des données images. Projets de ces dernières années : reconnaissance faciale, reconnaissance et recommandation musicale, programmation d'un chatbot.

Python - Machine learning avec scikit-learn
Honnêtement, mon intérêt pour Python doit beaucoup à la découverte des packages de statistique et de data mining qui l'accompagnent. « sciki...

Descente de gradient stochastique sous Python

Ce tutoriel fait suite au support de cours consacré à l'application de la méthode du gradient en apprentissage supervisé. Nous travaillons ...

Python - Statistiques avec SciPy
SciPy est une bibliothèque de calcul scientifique pour Python. Elle couvre de nombreux domaines (intégration numérique, interpolation, opti...

Programmation Python sous Spark avec PySpark

Dans la série « Je découvre Spark », voici un tutoriel consacré à la librairie PySpark pour la programmation Python sous Spark. Il vient en ...

Python : Manipulations des données avec Pandas

La manipulation des données est la base de l'activité du data scientist. Si on ne sait pas charger un fichier, exécuter des restrictions et ...

Analyse en composantes principales - Diapos

Mon premier contact avec l'analyse en composantes principales, technique populaire s'il en est, a été l'excellent ouvrage (pour l'économ...

Deep Learning avec Tensorflow et Keras (Python)

Tensorflow est une bibliothèque open-source développée par l'équipe Google Brain qui l'utilisait initialement en interne. Elle implémente d...

Excel avancé - Cours et exercices corrigés

J'ai dû arrêter mon cours d'Excel avancé (outils d'analyse et programmation VBA - Visual Basic pour Applications) cette année. Avec un p...

L'add-in Tanagra pour Excel 2007 et 2010

La macro complémentaire (" add-in " en anglais) " tanagra.xls " participe grandement à la diffusion du logiciel Tanagra....

Deep Learning avec Keras sous Knime

"En dehors de R et Python, point de salut alors pour le deep learning ?" me demande, un brin inquiet, un internaute. J'ai comp...

Et les outils payants ?

Positionnement des outils analytiques qui proposent (au moins) une licence commerciale - [Gartner Magic Quadrant for Data Science and Machine Learning Platforms](#), KDnuggets, Février 2019.

Codes de lecture (Gartner.com) :

- Leaders** execute well against their current vision and are well positioned for tomorrow.
- Visionaries** understand where the market is going or have a vision for changing market rules, but do not yet execute well.
- Niche Players** focus successfully on a small segment, or are unfocused and do not out-innovate or outperform others.
- Challengers** execute well today or may dominate a large segment, but do not demonstrate an understanding of market direction.

*Parlons-en de **KNIME**...*



Etudes de cas

Sous R et Python – Projets POC effectués par les étudiants

Traitements amusants sur le web

Reconnaissance d'objets
avec Google Image



<https://www.google.fr/imghp>

Détection de l'âge avec
la technologie Microsoft

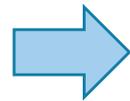


<https://www.how-old.net/#>

Analyse des offres d'emploi

Analyse de documents textuels (text mining) et classement / classification. Projet sous R (Shiny)

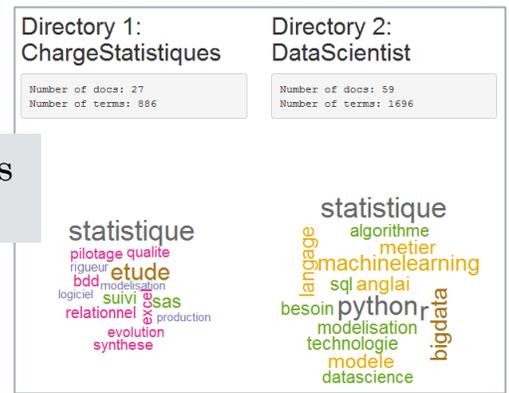
Offres d'emploi qui ont été étiquetées manuellement.



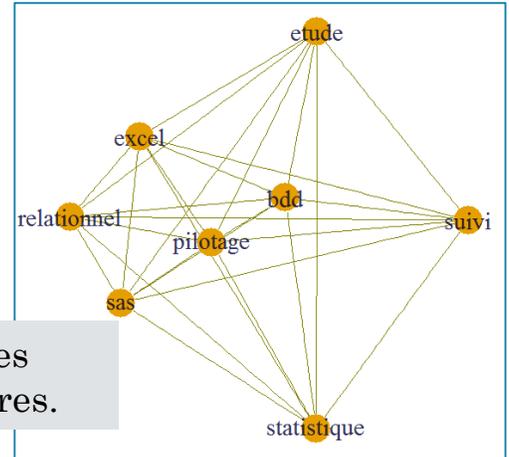
Analyse et développement d'un application Shiny

Métiers : Chargés d'études statistique, consultant BI, data analyst, data engineer, data manager, data miner, data scientist, data visualisation

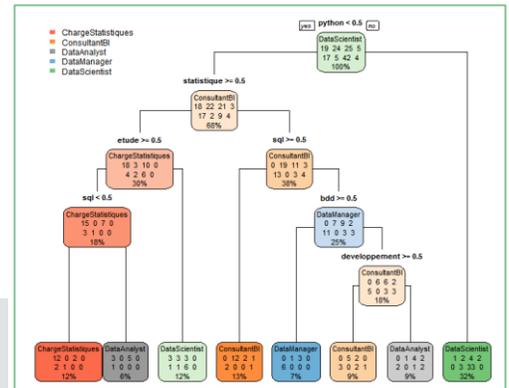
Mots clés fréquents selon les métiers



Association entre les termes dans les offres.



Identification des métiers selon les termes de l'offre.



Reconnaissance faciale

Démarche de recherche d'information par le contenu. Projet en Python (2016, 2019).

Disposer d'une banque d'images



Extraction de caractéristiques



Matrice de description, ligne : individus, colonnes : caractéristiques.

x1	x2	x3	x4	x5

Extraction de caractéristiques



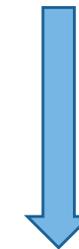
x1	x2	x3	x4	x5

Image « requête »



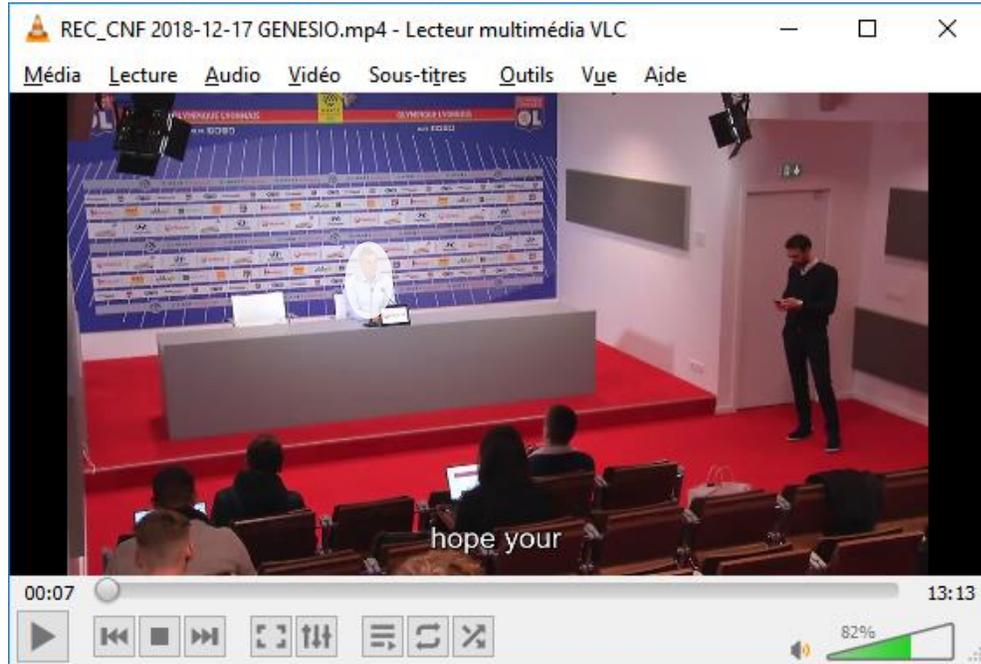
Vecteur de description de l'individu « requête »

Recherche de similarités.



Identification avec degré de fiabilité.

Traduction automatique de vidéos



Traduction automatique de vidéos
FR pour sous-titrage en anglais

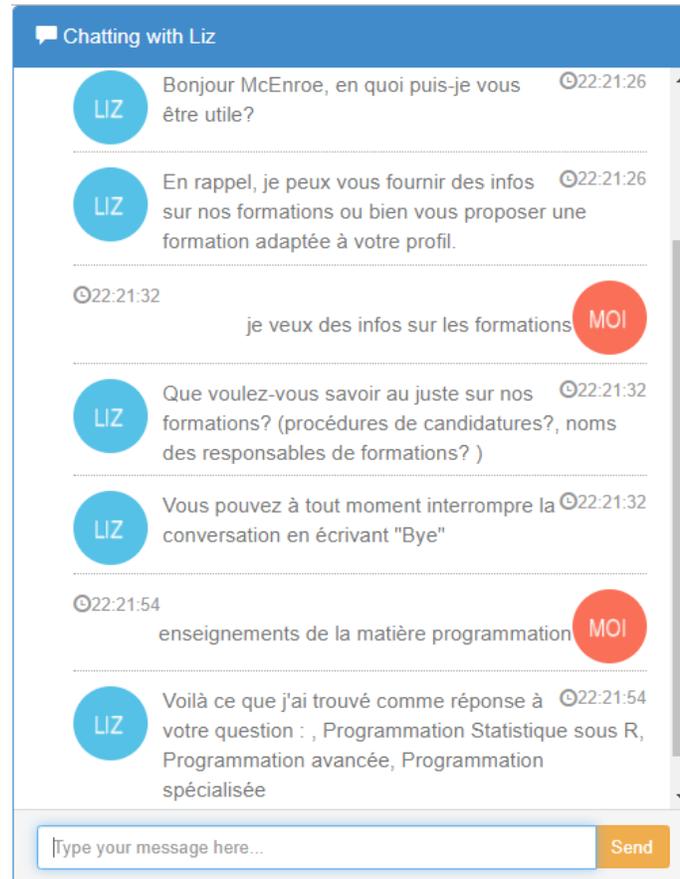
Speech-to-text

- Gestion du niveau de bruit
- Découpage de la bande son pour extraction des phrases

Utilisation d'outils de traduction

- Choix du plus performant
- Apprentissage du vocabulaire spécifique au domaine

Programmation d'un chatbot



Un chatbot pour aiguiller les étudiants dans leurs choix de formation

Traduire les questions formulées en langage naturel sous forme de requêtes SQL, avec identification des mots-clés, gestion des synonymes, des stopwords, etc.

Utilisation du framework Flask pour une architecture C/S.

Conclusion et bibliographie

R et Python jouent actuellement - tous deux - un rôle essentiel dans l'enseignement du Machine Learning.

Parce que

- Répondent parfaitement aux objectifs pédagogiques
- Permettent de développer des compétences de programmation
- Répondent aux attentes en matière de traitements – IA, Machine Learning – dans les entreprises (cf. les stages SISE)
- Et sont (de ce fait) explicitement spécifiées dans les offres d'emploi

Et...

Développer une expertise élevée dans les deux outils / langages n'implique pas un coût pédagogique additionnel fort.

Bibliographie - Webographie

Goebel M., Gruenwald L., « [A survey of data mining and knowledge discovery software tools](#) », ACM SIGKDD Explorations, 1(1), June 1999.

Un des premiers articles populaires ayant posé les bases de la comparaison de logiciels de data mining.

Ateliers du Master SISE, « [Logiciels de Data Science](#) », oct. 2016 ; « [Outils de la Data Science](#) », oct. 2017 ; « [Outils de la Dataviz](#) », oct. 2018.

Présentation, études de cas sous forme de scénarios de TD, tutoriels sur Youtube.

Piatetsky G., « Python leads the 11 top Data Science, Machine Learning platforms: Trends and Analysis », [KDnuggets 2019 Software Poll Results](#), May 2019.

Enquête annuelle, évolutions, comparaisons avec les années précédentes.

Piatetsky G., « Gainers, Losers, and Trends in Gartner 2019 Magic Quadrant for Data Science and Machine Learning Platforms », [KDnuggets](#), February 2019.

Enquête annuelle, évolutions, comparaisons avec les années précédentes.