

Plan du cours

- Introduction au data mining (1h)
- Apprentissage supervisé et non supervisé
- Focus sur l'apprentissage supervisé + mesures de performance et de validation (2h)
- Focus sur l'apprentissage non supervisé + mesures de performance et de validation (2h)
- Focus sur l'analyse de texte (2h)



Plan du jour

- Quelques définitions (datascience, apprentissage, machine learning, bigdata...etc.)...pour dire la même chose ?
- Quelques use case
- But du cours
- Le processus d'analyse de données
- L'équation d'un bon data scientist
- Visualisation des données
- ➤ Introduction à Python

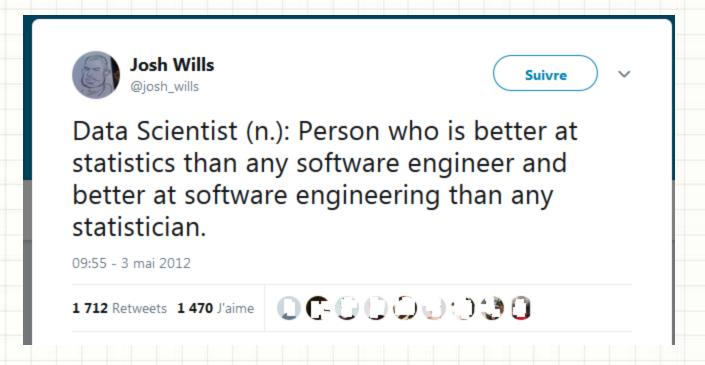
Quelques définitions

- La statistique est l'étude d'un phénomène par la collecte de données, leur analyse, leur traitement, l'interprétation des résultats et leur présentation afin de rendre les données compréhensibles par tous.
- Le « Machine Learning » (ou apprentissage artificiel, ou apprentissage statistiques) est un sous domaine de l'informatique qui donne au ordinateurs la possibilité d'apprendre sans être explicitement programmés.
- La « Data science » est un domaine interdisciplinaire qui utilise des méthodes scientifiques, des processus, des algorithmes et des systèmes pour extraire des connaissances à partir de données sous différentes formes, structurées ou non. ... est un « concept pour unifier les statistiques, l'analyse des données, l'apprentissage automatique et leurs méthodes associées »
- Le « Big data » incluent généralement des ensembles de données dont la taille dépasse la capacité des outils logiciels couramment utilisés pour capturer, organiser, gérer et traiter les données dans un délai raisonnable.

Quelques définitions... mon point de vue

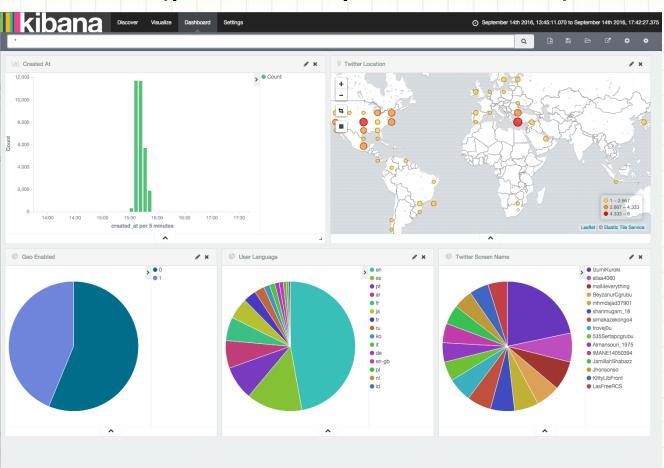
- Les termes sont tous connexes, c'est une évolution des statistiques qui est le terme le plus ancien.
 - 1^{er} évolution vers le machine learning lorsque les informaticiens ont pris un peu pour eux le sujet. Le machine learning sont des statistiques de point de vue informatique!... ou en d'autre termes, comment faire apprendre les ordinateurs sans qu'ils soient explicitement programmés, en utilisant des concepts statistiques.
 - 2ème « évolution » vers le bigdata, pour adresser l'explosion du volume de données à traiter
 - d'où les technologies de bigdata : hadoop, map/reduce, spark, NoSQL, etc.

Quelques définitions... et la data science alors ?

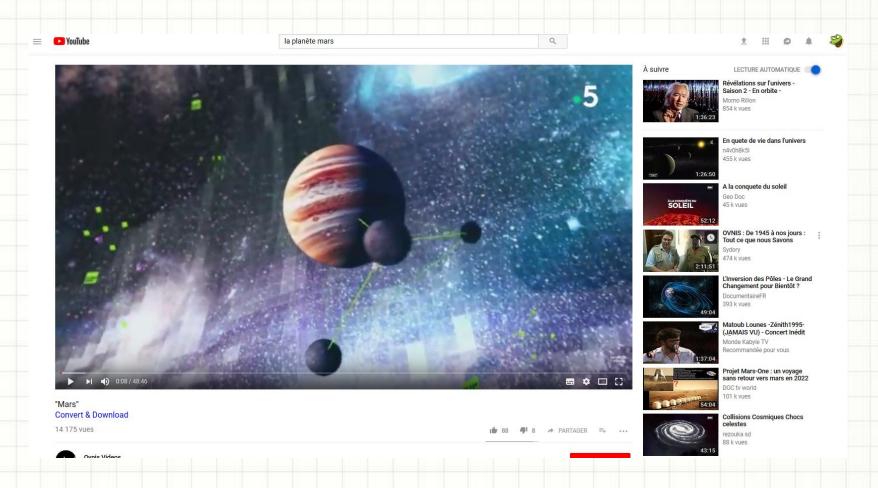


Quelques use case

Dashboard (peut être personnalisé)



Moteur de recommandation

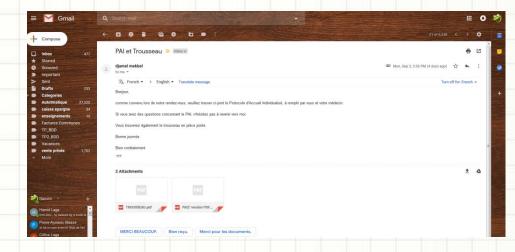


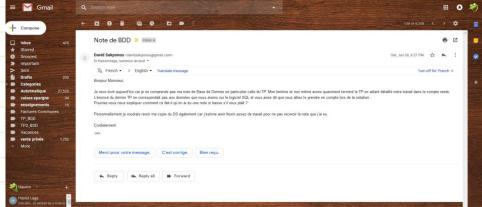
Moteur de recommandation (publicité)



- Moteur de recherche
 - Indexation d'un volume gigantesque de données: structurée, semi-structurées, et non-structurées

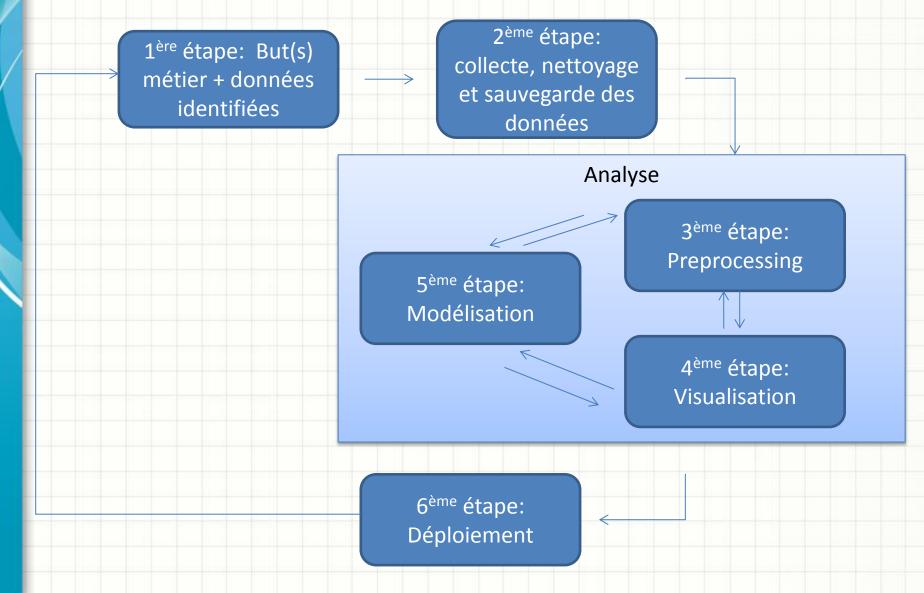
- Fonction de prédiction intégrée dans un logiciel
 - détecteur de spams
 - classification
 automatique de
 documents
 - Détection d'intention dans un texte, et rédaction de réponse automatique





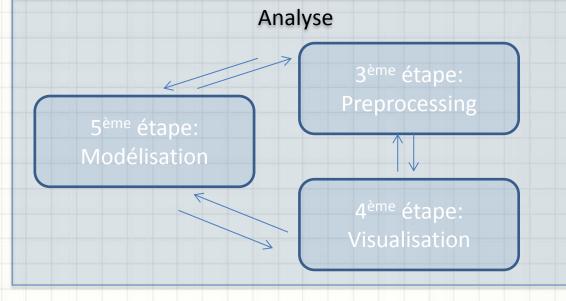
But du cours

- Nous n'allons pas faire du bigdata
- Je suis plus ingénieur © que statisticien 😊
- Par conséquent, nous allons faire du machine learning (ou apprentissage machine en français), en se focalisant plus sur la pratique que la théorie

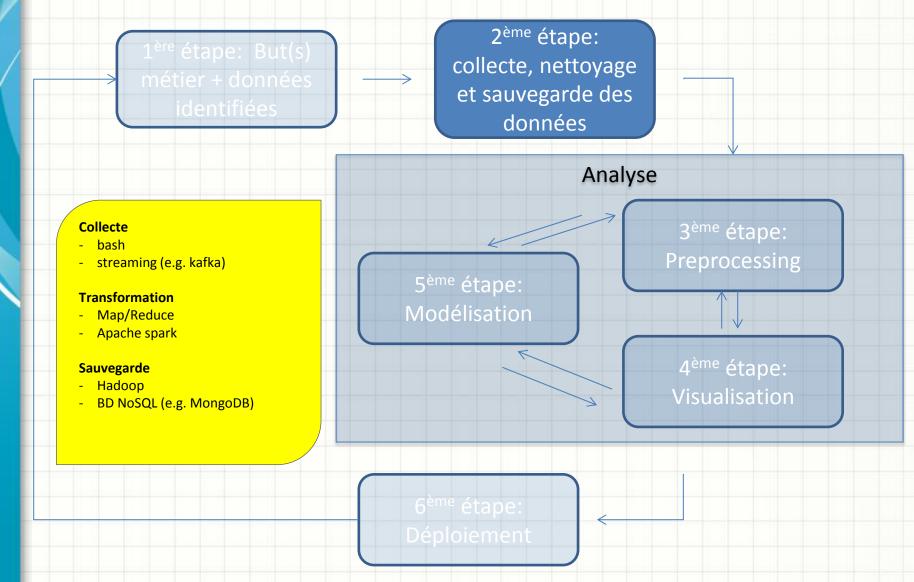


1^{ère} étape: But(s) métier + données identifiées 2^{ème} étape: collecte, nettoyage et sauvegarde des données

- détection de spams / exemples de spams et de mails normal
- Regroupement de profiles client / liste de profiles client
- Estimation du prix de l'immobilier / liste de biens avec leurs caractéristiques associés au prix



6^{ème} étape: Déploiemen



1^{ere} étape: But(s) métier + données identifiées 2^{ème} étape: collecte, nettoyage et sauvegarde des données

Preprocessing

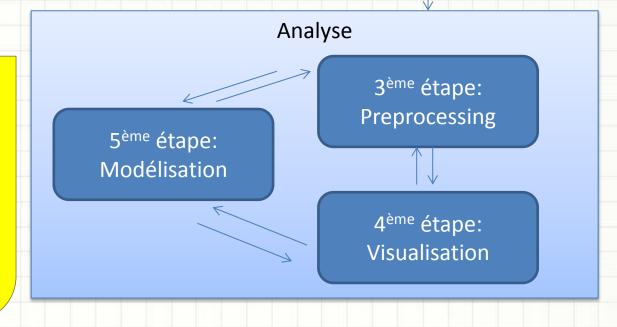
- Enlever les stopwords (mot vide en Français)
- Traiter les valeurs null
- Oversampling/Undersampling
- Réuction de dimension

Visualisation

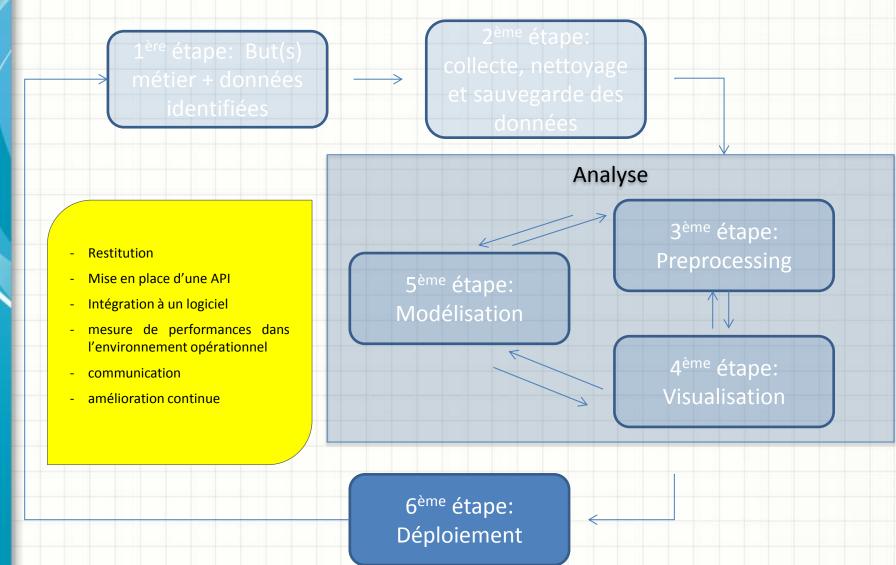
- Détection des catégories sousreprésentées
- Détection de vocabulaire

Modélisation

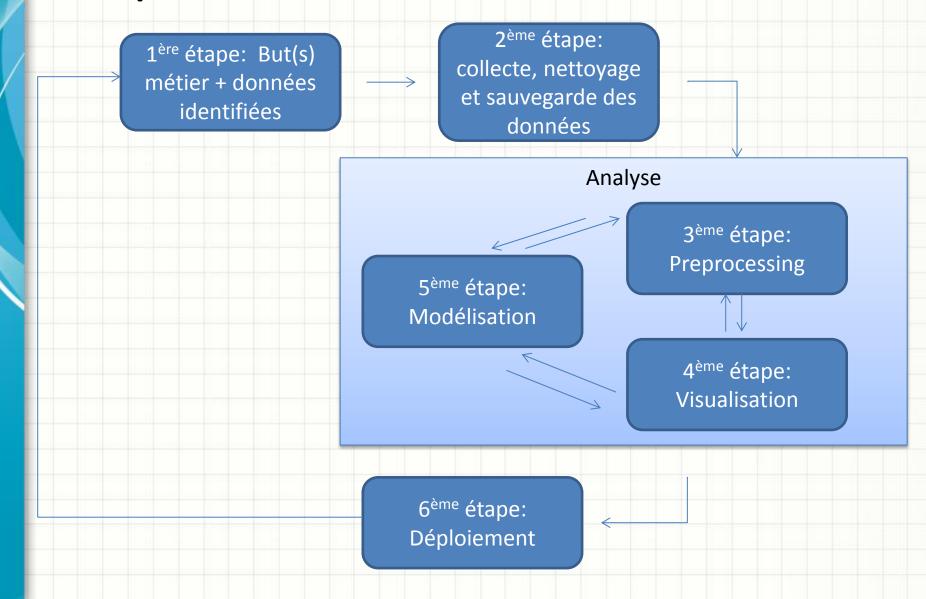
- Tests de différents algorithmes avec différents paramètres
- Analyse des performances



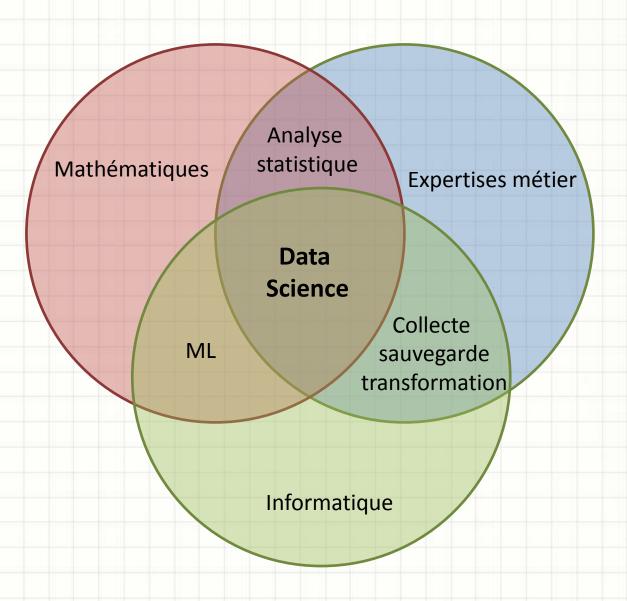
6^{ème} étape:



L'équation d'un bon data scientist



L'équation d'un bon data scientist

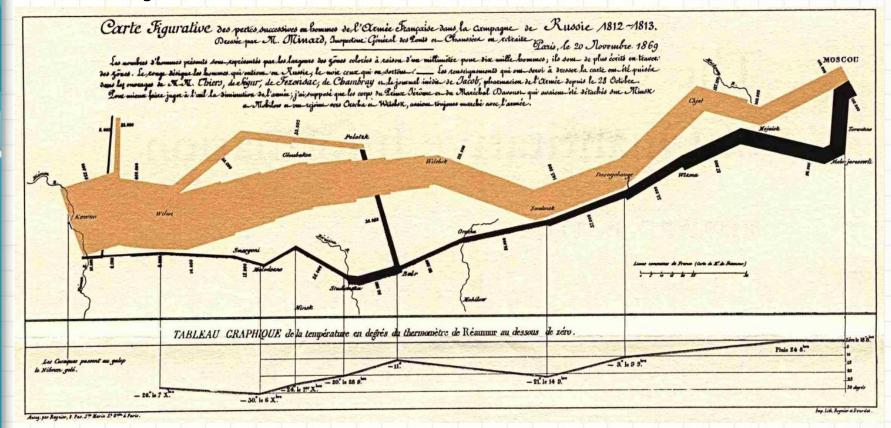


L'équation d'un bon data scientist

 Josh Wills: "Person who is better at statistics than any software engineer and better at software engineering than any statistician."

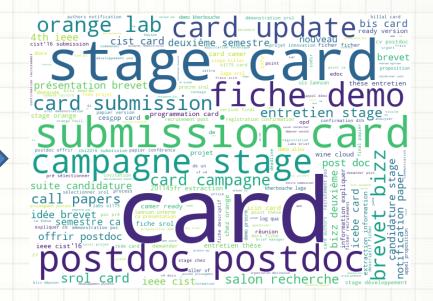
- Utilisée à la fois pour
 - L'étape d'exploration de jeux de données et de prétraitement
 - L'étape de restitution des résultats

 Charles Joseph Minard, ingénieur civil français, 1969

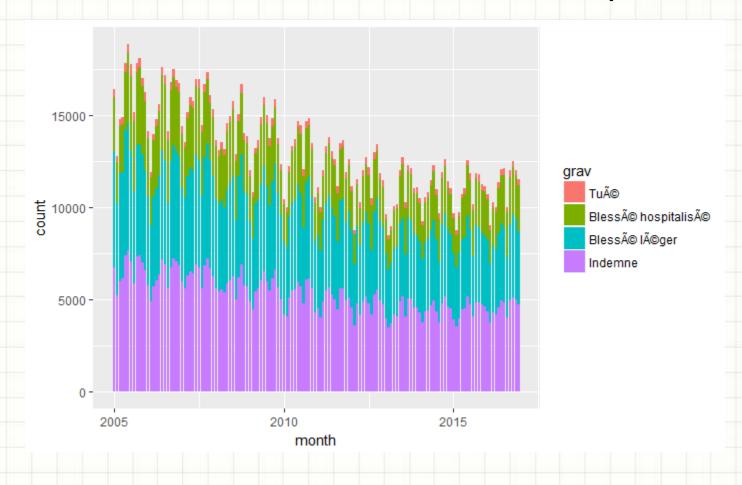


Prétraitement des données, exemples

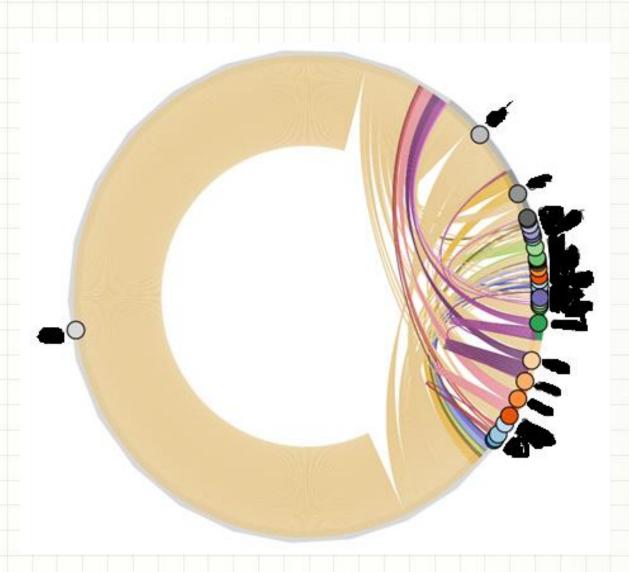




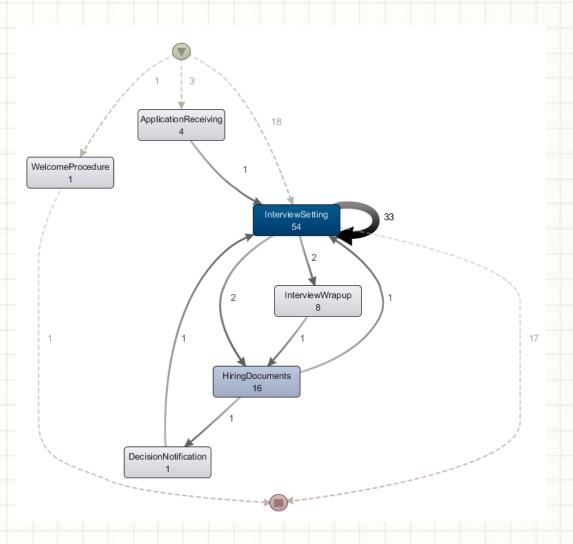
Prétraitement des données, exemples



 Restitution, exemples



Restitution,
 exemples



Les langages / outils de la datascience

- R
- Python
- Julia
- •

Introduction à python

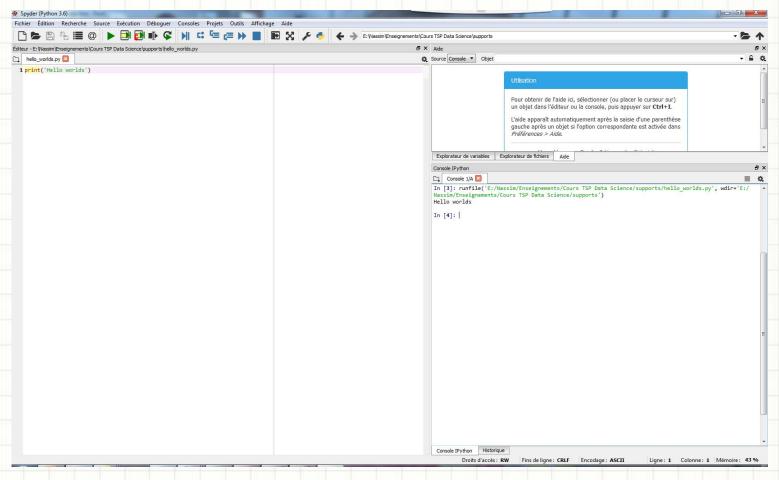
- Langage interprété
 - Plus lent, mais portable assez facilement d'une machine à une autre
 - Nécessite un interpréteur auquel on soumet un fichier .py
- Comme tout langage évolué il intègre :
 - des types de données simples et évolués, des boucles, des conditions, des procédures et des fonctions
 - la possibilité de programmer en Orienté Objet avec les notions de classes, d'instances, et d'héritage
- Une communauté très active, donnant lieu à de nombreuse librairies très intéressantes
 - Numpy, Scipy, scikit-learn, pandas, etc.

Console interactive

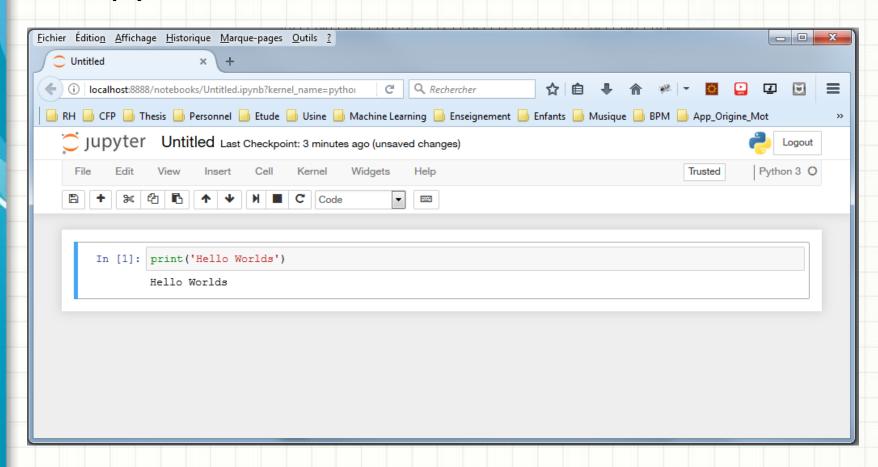
 Création d'un fichier .py et exécution à partir de la ligne de commande



spyder

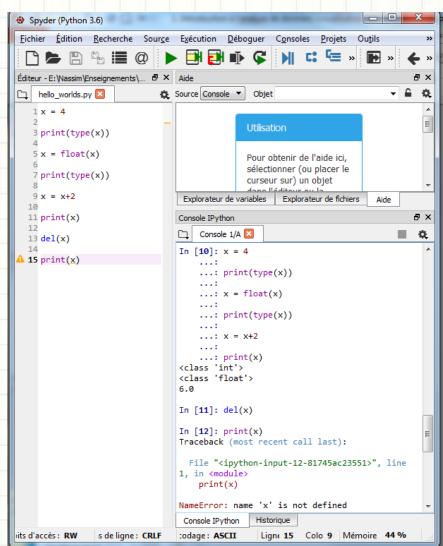


Jupyter



Introduction à python – les opérations de base

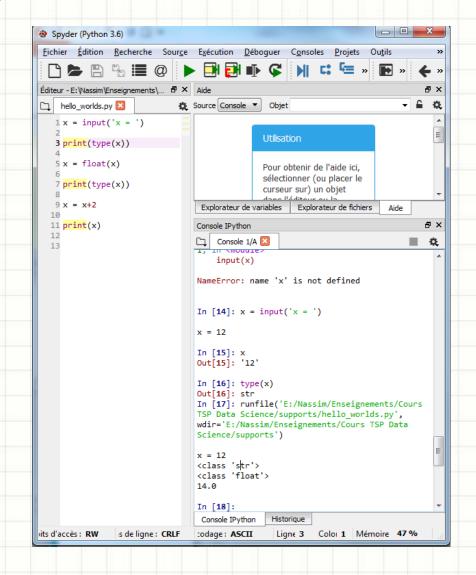
- Affectation
- Typage & Transtypage
- Connaître le type d'un objet/d'une variable
- Calcul
- Libération de la mémoire



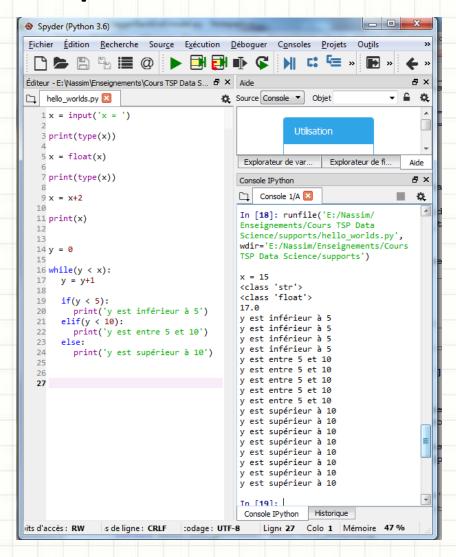
Introduction à python – les types

Types primitifs	Types complexes
int	list
float (double)	tuple : N-uplet
str	dict
Bool (True, False)	set: liste d'éléments uniques
	file
	None
	exception

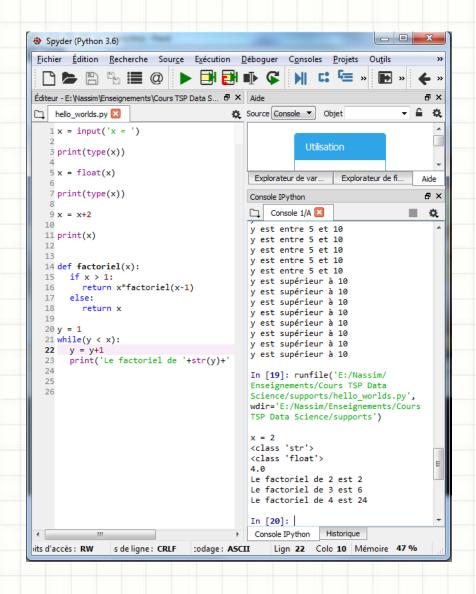
Introduction à python – les entrées/sorties



Introduction à python – structures algorithmiques



Introduction à python – les fonctions



Introduction à python – les listes

- C'est une collection d'éléments, potentiellement identiques, potentiellement de types différents
- Les éléments sont indexés numériquement x[0] récupère le premier élément de la liste
- On connait la taille de la liste grace à la fonction len(x)
- La fonction del(x[i]) supprime l'élément d'index i de la liste x
- x.append(element) ajoute l'élément en fin de liste

```
27 ys = [2, 3, 5, 10, 20, 20]
28 for y in ys:
29 print('Le factoriel de '+str(y)+' est '+str(factoriel(y)))
```

Python et l'analyse de données

- Librairies
 - Numpy
 - SciPy
 - Pandas
 - Scikit-learn
 - matplotlib (visualisation)
 - ggplot (visualisation)
 - Seaborn (visualisation)

- Une librairie permettant de gérer les données de 2 dimensions, à des fin d'analyse
- Concept issues des outils de statistiques telle que : R, matlab, et SAS
- Deux types de données clefs
 - Series
 - Dataframe

- Series
 - Structure de données d'une dimension, indexable
 - Plusieurs manière de créer une Series

- Series
 - Structure de données d'une dimension, indexable
 - Plusieurs manière de créer une Series

- Series
 - Sélection d'une valeur par index

```
In [21]: #Accès aux contenu s['p1']
Out[21]: 'John'
```

Sélection de plusieurs valeurs

```
In [23]: #Accès aux contenu
s[['p1', 'p3']]
Out[23]: p1 John
p3 Jack
dtype: object
```



- Series
 - Sélection par filtre

- Opérations groupées, potentiellement conditionnées

- dataframe
 - Une structure de données à deux dimensions
 - Contient des lignes et des colonnes
 - indexée au niveau des lignes et des colonnes
 - peut être considérée comme un dictionnaire de Series, avec un index partagé

Dataframe

- Plusieurs manière de créer une dataframe
- à partir d'un dictionnaire

- Dataframe
 - Plusieurs manière de créer une dataframe
 - à partir fichier csv

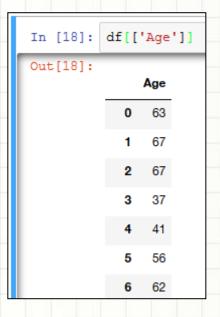
[10]: df = df	<pre>if = pandas.read_csv('Heart.csv') if</pre>															
[10]:	,	Jnnamed: 0	Age	Sex	ChestPain	RestBP	Chol	Fbs	RestECG	MaxHR	ExAng	Oldpeak	Slope	Ca	Thal	AHD
0)	1	63	1	typical	145	233	1	2	150	0	2.3	3	0.0	fixed	No
1		2	67	1	asymptomatic	160	286	0	2	108	1	1.5	2	3.0	normal	Yes
2	2	3	67	1	asymptomatic	120	229	0	2	129	1	2.6	2	2.0	reversable	Yes
3		4	37	1	nonanginal	130	250	0	0	187	0	3.5	3	0.0	normal	No
4		5	41	0	nontypical	130	204	0	2	172	0	1.4	1	0.0	normal	No
5	,	6	56	1	nontypical	120	236	0	0	178	0	0.8	1	0.0	normal	No
6	i	7	62	0	asymptomatic	140	268	0	2	160	0	3.6	3	2.0	normal	Yes
7	,	8	57	0	asymptomatic	120	354	0	0	163	1	0.6	1	0.0	normal	No
8	1	9	63	1	asymptomatic	130	254	0	2	147	0	1.4	2	1.0	reversable	Yes
9	,	10	53	1	asymptomatic	140	203	1	2	155	1	3.1	3	0.0	reversable	Yes

Dataframe

- Plusieurs manière de créer une dataframe
- à partir d'une base de données, à travers une requête SQL
- ...etc.

```
In [13]: import pandas as pd
         from sqlalchemy import *
         from sglalchemy.orm import sessionmaker
         from sqlalchemy.types import TIMESTAMP
         from sglalchemy.exc import IntegrityError
         from sqlalchemy.sql import select
         from sqlalchemy import desc
         #Loading data
         database cnf = 'mysql+pymysql://root:@localhost:3306/examen'
         engine = create engine(database cnf, echo=False)
         conn = engine.raw connection()
         df = pd.read sql(sql='SELECT * FROM `compte`', con=conn)
         C:\ProgramData\Anaconda3\lib\site-packages\pymysql\cursors.py
         E9)' for column 'VARIABLE VALUE' at row 480")
           result = self. query(query)
Out[13]:
            id id client num agence solde
                               1 22000
                               1 7000
```

- Les opérations sur les dataframes
 - Accès par colonnes
 - Accès par ligne



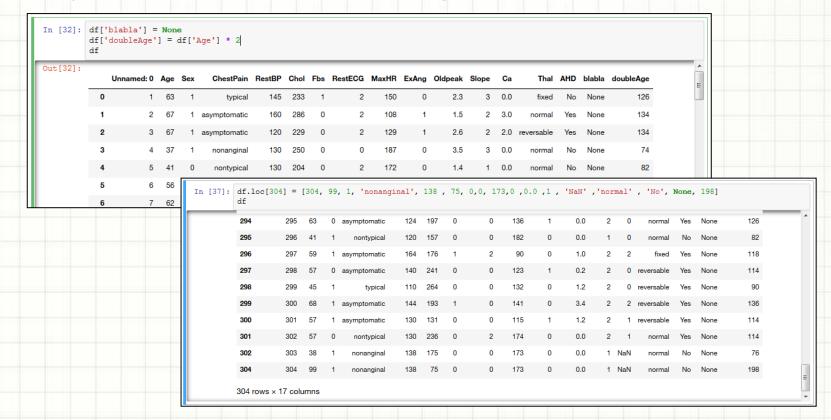
```
In [27]: df.loc[[0]]

Out[27]:

Unnamed: 0 Age Sex ChestPain RestBP Chol Fbs RestECG MaxHR ExAng Oldpeak Slope Ca Thal AHD

0 1 63 1 typical 145 233 1 2 150 0 2.3 3 0.0 fixed No
```

- Les opérations sur les dataframes
 - Ajout de colonnes et de lignes

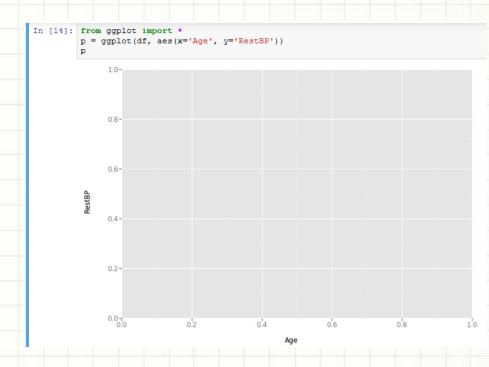


- D'autres opérations avancées:
 - description statistiques: moyenne (mean), min, max, sum, describe, ...etc
- Visualisation
 - ggplot

- Grammar Graphics Plot
- Créé pour permettre de créer « facilement » des graphe pour visualiser les données et leurs caractèristiques
- Écrit initialement pour R, mais utilisé dans Python
- Les pré-requis
 - matplotlib
 - numpy
 - scipy
 - statmodels
- Installation
 - pip install ggplot

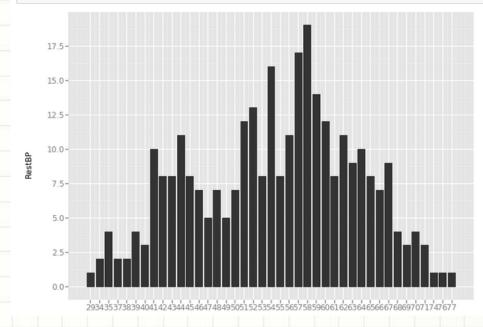
- Création d'un graphique avec python ggplot principe :
 - 1. Spécification des données et des axes
 - 2. Description des axes et leur apparence
 - 3. Ajout de couche sur le graphe

 Spécification des données et des axes



- Spécification des données et des axes
- Description des axes et leur apparence

```
from ggplot import *
p = ggplot(df, aes(x='Age', y='RestBP'))
p = p + geom_bar()
p
```



- Spécification des données et des axes
- Description des axes et leur apparence
- Ajout de couche sur le graphe

- Spécification des données et des axes
- Description des axes et leur apparence
- Ajout de couche sur le graphe

```
from ggplot import *
p = ggplot(aes(x='Age', y='MaxHR'), data=df)
C:\ProgramData\Anaconda3\lib\site-packages\ggplot\stats\stat_smooth.py:77: FutureWarning: sort(
  smoothed_data = smoothed_data.sort('x')
```