



Choisir son notebook



Qu'est ce qu'un notebook ?

Outil de prédilection des *Data Scientists*, et dont l'usage se démocratise auprès des *Data Analystes*, le notebook les aide à concevoir des programmes dans l'ordre qui leur semble le plus adapté (programmation lettrée).

Il possède 3 grandes fonctionnalités :

- le mode *Markdown* qui procure les fonctionnalités d'un traitement de texte (Formules LaTeX, Structuration, typographie, hyperliens, etc.) via une syntaxe très simple (langage de balisage léger), pour un rendu en html ;
- le mode *Code*, qui permet d'interpréter des langages, de se connecter à des SGBD, à des moteurs de calcul ou encore à des interfaces système ;
- l'affichage de rendu visuels.

Typiquement un notebook est un client qui s'utilise à partir d'un navigateur internet et qui nécessite le lancement d'un service à partir d'un serveur.

Le notebook n'a pas vocation à servir à développer des programmes complexes, (pour cet usage on préférera un IDE), mais se révèle très utile au *Data Scientist* pour explorer les données et comparer les algorithmes.

On pourra également l'utiliser pour la présentation de travaux, il s'agit d'un outil immédiatement opérationnel et très intuitif.



Jupyter VS Zeppelin



Jupyter et Zeppelin sont actuellement les deux notebooks les plus utilisés, Qu'est-ce qui distingue l'un de l'autre et quels sont leurs avantages comparatifs ?



Généralités



Version actuelle
sortie le 28 septembre 2016 **2.4.0**

Version actuelle
sortie le 15 octobre 2016 **0.6.2**

fork du projet IPython Depuis **2014**

Depuis **2015**

Licence **BSD Modifiée**

Licence **Apache**



Partenaires



Sponsor

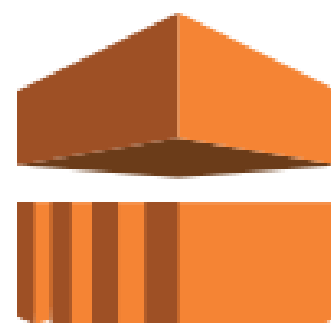


Partenaires institutionnels

Distributions qui l'intègrent



Sociétés et institutions qui l'utilisent



Communauté



Documentation

Très complète

Très sommaire

<http://jupyter-notebook.readthedocs.io/en/latest/>

<https://zeppelin.apache.org/docs/latest/>

Support

Communauté 3 fois plus active

Activité relativement faible

<http://stackoverflow.com/questions/tagged/jupyter>

<http://stackoverflow.com/questions/tagged/apache-zeppelin>



Implémentation



Mode service

Très simple
Installer anaconda, et lancer la commande
jupyter notebook dans un shell.

Très simple,
Télécharger le .tgz sur le site de zeppelin, décompresser,
lancer la commande ./zeppelin-daemon.sh start dans un
shell à partir du dossier bin et taper localhost:8080 dans un
navigateur internet

Mode hub

Il est possible d'utiliser Jupyter et Zeppelin en mode Hub au moyen de **JupyterHub** et **ZeppelinHub**. Chaque utilisateur d'une même société pourra ainsi créer son propre compte sur le hub via un portail commun. Ceci lui permettra d'avoir accès à un outil en silo qui hébergera tous ses notebooks.

JupyterHub

Possibilité d'installer Jupyterhub soi-même
(composant d'Anaconda dont les sources sont disponibles)

ZeppelinHub

Licence propriétaire
(solution hébergée par ZEPL Inc)

Possibilité d'utiliser les notebooks en mode
collaboratif
(gestion des droits et des utilisateurs)

Possibilité de publier les notebooks
(en mode public ou privé)



Interopérabilité



- Interpréteurs et connecteurs ne sont pas natifs
- Un notebook ne peut interpréter qu'un seul langage

Pour lancer les notebooks dans d'autres langages, il est nécessaire d'installer des noyaux (*kernels*) additionnels. La liste des kernels (<https://github.com/ipython/ipython/wiki/IPython-kernels-for-other-languages/>) comprend, les **Interpréteurs de langage** (liste non exhaustive) :

- C#
- Erlang
- ghc
- Go
- gnuplot
- julia
- Lua
- Matlab
- nodejs
- Octave
- Perl
- PHP
- Processing
- R
- ruby
- Scala
- Scilab
- Torch

Les connecteurs :

- SGBD distribués : Redis
- Moteurs de calculs distribués : Scala Spark
- Interface système :
- bash
- windows

- Interpréteurs et connecteurs natifs
- Un notebook peut interpréter plusieurs langages à la fois

Pour activer tous les interpréteurs et langages disponibles, il suffit de lancer ./install-interpreter.sh -a depuis le dossier bin. Les **interpréteurs de langages** disponibles sont :

- Python
- Scala
- R

Les connecteurs :

- SGBD distribués
 - bigquery
 - cassandra
 - hbase
 - hive
 - Jdbc
 - lens
 - postgresql
- Moteurs de calculs distribués
 - Spark
 - Fink
 - Ignite
- Interface système
 - Shell
 - Alluxio
 - HDFS file interpreter



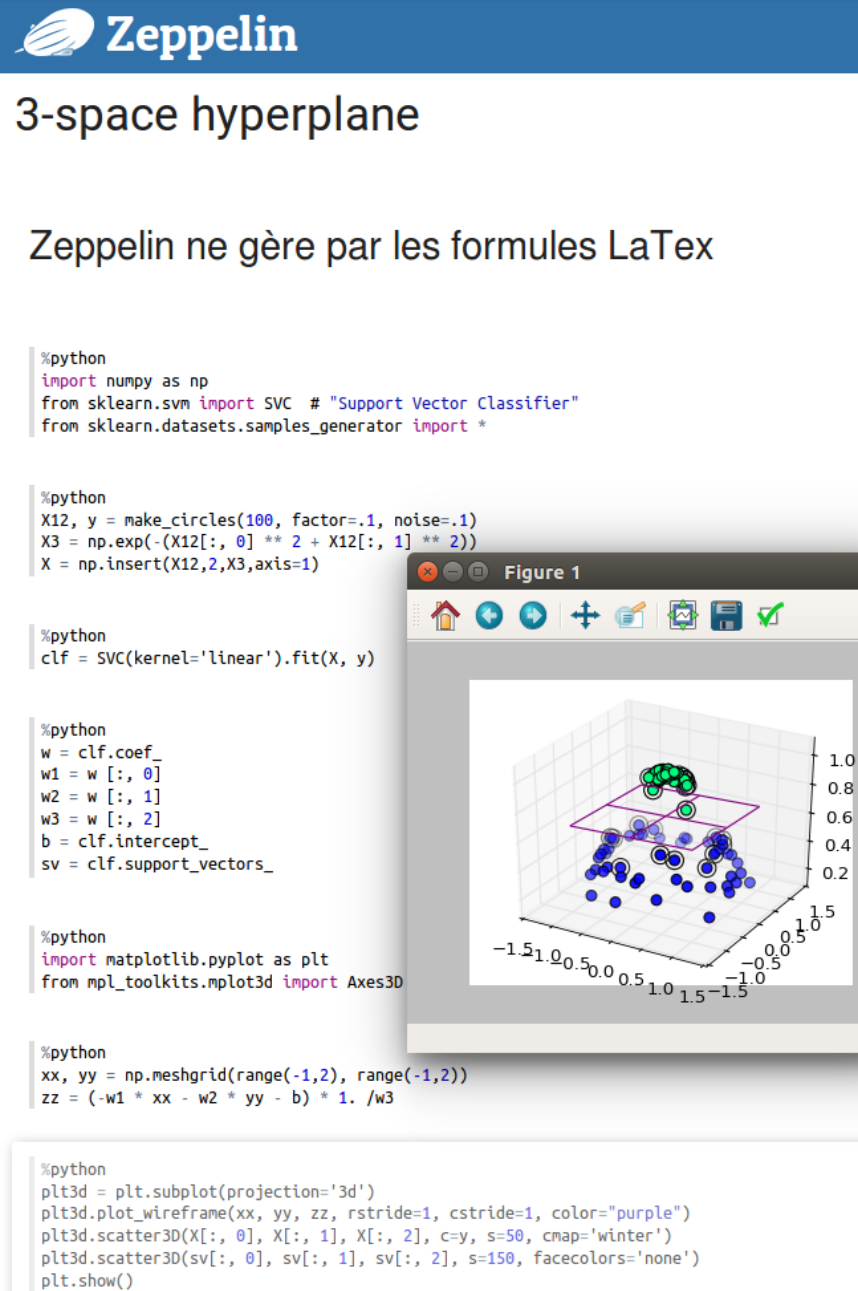
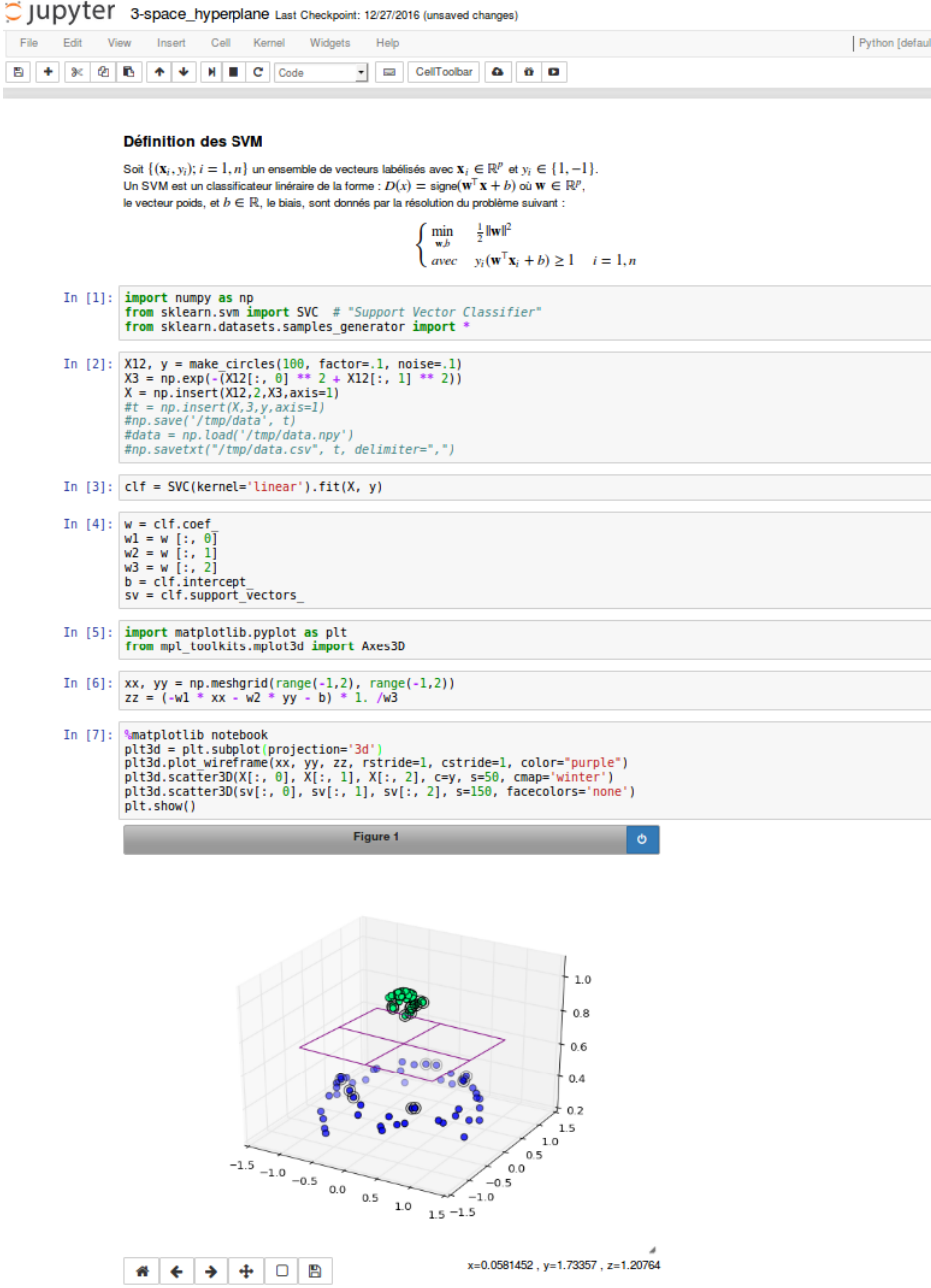
Interface



Ergonomie

- Supporte l'**auto-complétion** (TAB)
- Les **graphiques** peuvent s'afficher de manière **interactive** (images) à l'intérieur des *notebook*.
- Formules **LaTeX** gérées par le mode *Markdown*

- Auto-complétion** très limitée (Ct1 + .)
- Les **graphiques** s'ouvrent dans une nouvelle fenêtre et sont **interactifs**.
- Chaque **fenêtre** du *notebook* est **dimensionnable**.
- Supporte **Angular**



Formats



Sauvegarde

Ipynb

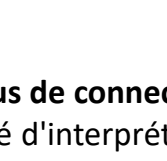
JSON

Export

- Notebook (ipynb)
- PDF via LaTeX
- Presentation
- Python
- HTML
- Markdown
- JSON



Récapitulatif



Avantages comparatifs

Projet plus mature
Plus utilisé
Documentation plus complète
Beaucoup plus d'interpréteurs de langages
Plus de formats d'export
Auto-complétion
JupyterHub est une solution open source qu'on peut installer soi-même.
Formules LaTeX

Plus de connecteurs
Possibilité d'interpréter plusieurs langages dans un même notebook
Chaque fenêtre du notebook est dimensionnable
Possibilité d'avoir des graphiques dynamiques avec Angular
ZeppelinHub permet d'avoir des notebooks collaboratifs et publiables

Verdict

Data science

Data Engineering