



INSTITUT NATIONAL
DES SCIENCES
APPLIQUÉES
TOULOUSE

IA FRAMEWORKS

COURSE PRESENTATION

Brendan Guillouet

November ??, 2020

Institut National des Sciences Appliquées de Toulouse

TABLE OF CONTENTS

Definition & Context

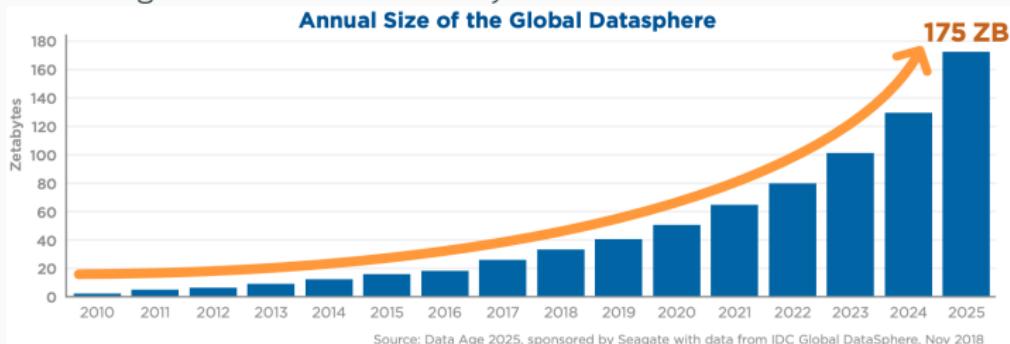
Organization & Evaluation

Github Reminder

DEFINITION & CONTEXT

BIG DATA and its 3 Vs :

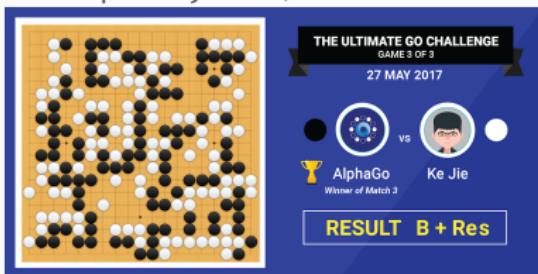
- **VOLUME** : Exponential growth.
 - According to IDC : from 33 zeta bytes in 2018 to 175 in 2025.?



- **VARIETY** : numeric, signal, images, time series, graph, text etc...
- **VELOCITY** : real time decision.
- **(VALORISATION)**.

Now ARTIFICIAL INTELLIGENCE (again)

- BIG DATA,
- COMPUTATION POWER AVAILABILITY : GPU, cloud computing.
- 'NEW' ALGORITHMS : Deep Learning, Reinforcement learning,
- that produce impressive results, (Alpha Go, Atari Game defeated, write poetry etc..).

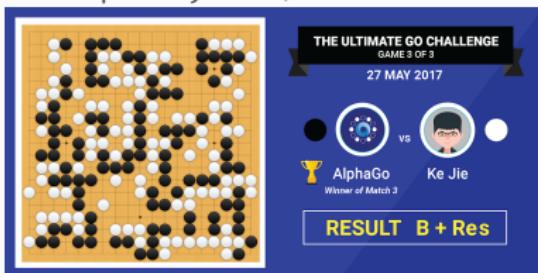


- .. and benefits some good marketing.

IA & BIG DATA

Now **ARTIFICIAL INTELLIGENCE** (again)

- **BIG DATA**,
- **COMPUTATION POWER AVAILABILITY** : GPU, cloud computing.
- **'NEW' ALGORITHMS** : Deep Learning, Reinforcement learning,
- that produce impressive results, (Alpha Go, Atari Game defeated, write poetry etc..).



- .. and benefits some good marketing.

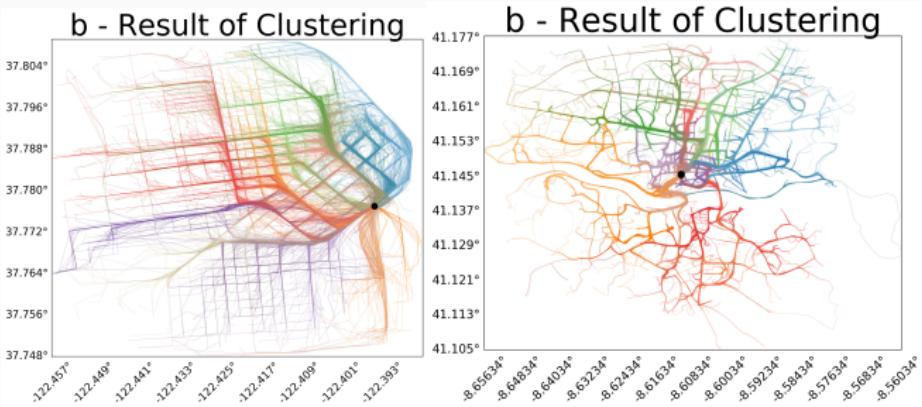
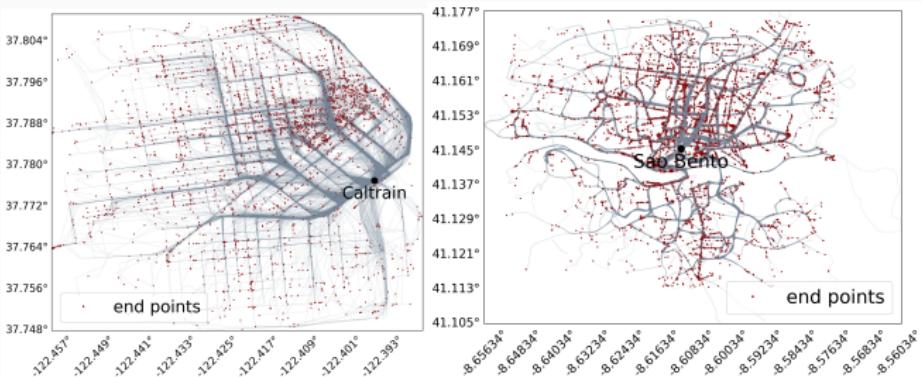
The definition is still very unclear, but IA is still NOT here.

A WIDE RANGE OF APPLICATION FIELDS...

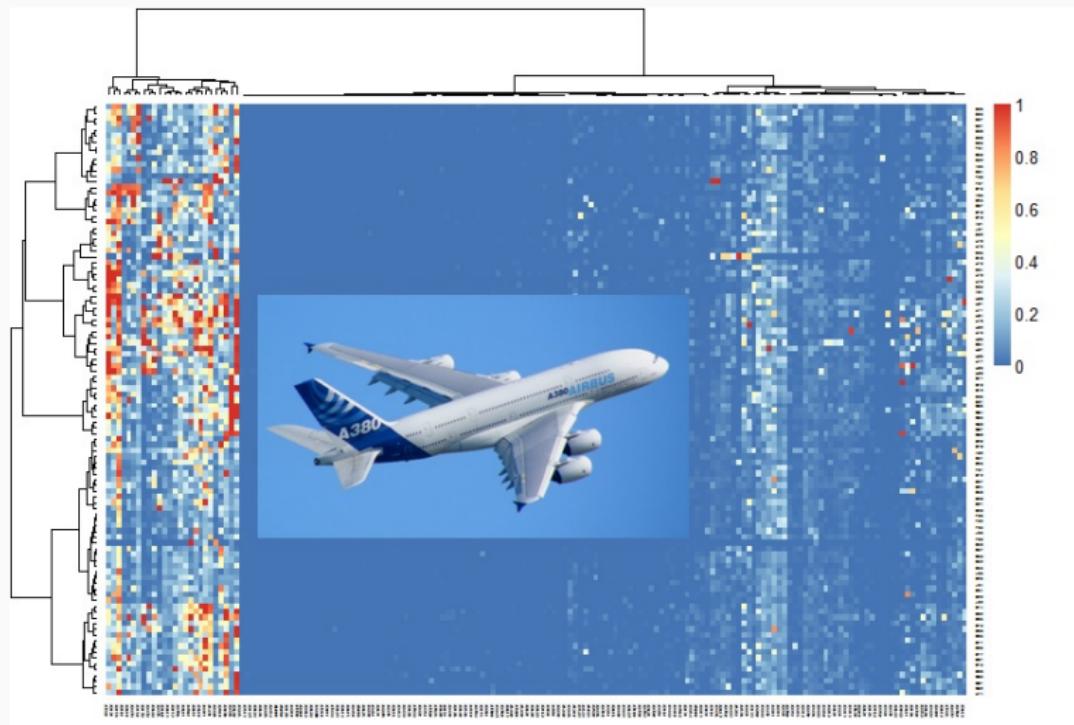
..at different levels of advancement.

- **E-COMMERCE** : Recommendation system, advertising, targeting.
- **TEXT TRANSLATION** : DeepL , Google Translate.
- **IMAGE & VIDEO** : Classification, videosurveillance etc..
- **GAME** : Chess, Go, Atari Game, StarCraft.
- **PUBLIC** : administrations, e-health (*open data*)
- **INDUSTRY** : 4.0 industry, Autonomous vehicle.
- **METEOROLOGY** : Hybridation of learning and physical model
- **CERTIFICATION**.
- **INTERNET OF THINGS.**

VARIETY : TRAJECTORY CLUSTERING FROM GPS DATA.



ANALYSIS OF FLIGHT INCIDENT MESSAGES



2020 INTERNSHIP OF INSA STUDENT

- Detection of field (agriculture) and batiment (military) from satellite images. *SSII & Startup*
- Anonymization of name on text documents. *SSII*
- Fall detection of elder people. *PME*
- Autonomous drone flight. *SSII*
- Rain prediction in Ivory Coast. *Startup*
- Detection of diabetic retinopathy from x-ray photo. *Startup-Norway*
- Estimation of fish results. *Research Center*

JOB SPECIFICATION

TASK

- Managing data Warehousing
- Ensure Data accessibility & consumption via API.
- Continuously monitoring and testing the system to ensure optimized performance.
- Extraction Transformation & Load
- Cleaning and organizing raw data.
- Programming skills.
- Using descriptive statistics to get a big-picture view of their data.
- Analyzing interesting trends found in the data.
- Creating visualizations and dashboards to help the company interpret and make decisions with the data.
- Presenting the results of a technical analysis to business clients or internal teams.
- Evaluating statistical models to determine the validity of analyses.
- Using machine learning to build better predictive algorithms.
- Testing and continuously improving the accuracy of machine learning models.
- Continuously Perform scientific Watch.
- Building data visualizations to summarize the conclusion of an advanced analysis.

JOB SPECIFICATION

TASK	DATA ENGINEER	DATA ANALYST	DATA SCIENTIST
• Managing data Warehousing	✓✓	✗	✗
• Ensure Data accessibility & consumption via API.	✓✓	✗	✗
• Continuously monitoring and testing the system to ensure optimized performance.	✓✓	✗	✗
• Extraction Transformation & Load	✓✓	✓	✓
• Cleaning and organizing raw data.	✓	✓	✓
• Programming skills.	✓	✓	✓
• Using descriptive statistics to get a big-picture view of their data.	✗	✓✓	✓
• Analyzing interesting trends found in the data.	✗	✓✓	✓
• Creating visualizations and dashboards to help the company interpret and make decisions with the data.	✗	✓✓	✓
• Presenting the results of a technical analysis to business clients or internal teams.	✗	✓✓	✓
• Evaluating statistical models to determine the validity of analyses.	✗	✓	✓✓
• Using machine learning to build better predictive algorithms.	✗	✓	✓✓
• Testing and continuously improving the accuracy of machine learning models.	✗	✗	✓✓
• Continuously Perform scientific Watch.	✗	✗	✓✓
• Building data visualizations to summarize the conclusion of an advanced analysis.	✗	✗	✓✓

JOB SPECIFICATION

TASK	You're Here ↓		
	DATA ENGINEER	DATA ANALYST	DATA SCIENTIST
• Managing data Warehousing	✓✓	✗	✗
• Ensure Data accessibility & consumption via API.	✓✓	✗	✗
• Continuously monitoring and testing the system to ensure optimized performance.	✓✓	✗	✗
• Extraction Transformation & Load	✓✓	✓	✓
• Cleaning and organizing raw data.	✓	✓	✓
• Programming skills.	✓	✓	✓
• Using descriptive statistics to get a big-picture view of their data.	✗	✓✓	✓
• Analyzing interesting trends found in the data.	✗	✓✓	✓
• Creating visualizations and dashboards to help the company interpret and make decisions with the data.	✗	✓✓	✓
• Presenting the results of a technical analysis to business clients or internal teams.	✗	✓✓	✓
• Evaluating statistical models to determine the validity of analyses.	✗	✓	✓✓
• Using machine learning to build better predictive algorithms.	✗	✓	✓✓
• Testing and continuously improving the accuracy of machine learning models.	✗	✗	✓✓
• Continuously Perform scientific Watch.	✗	✗	✓✓
• Building data visualizations to summarize the conclusion of an advanced analysis.	✗	✗	✓✓

WHAT DO WE NEED TO START?

- Local or cloud computing? (*GCP, AWS, AZURE*)
- Size of disk? Size of RAM? CPU or GPU? How many of its?
- From where and how did you get the data : local/bucket?
- How are they stored ? Csv, Data base ? Which one ?
- Distributed architecture ? (*Hadoop, Spark*)
- In which format do you expect the results ?
- How it is automate ?



DataFarming

A NEW ECONOMIC MODEL

SITUATION :

- Algorithm are available and free (python library, github etc..)
- Most used Software (Tensorflow, Spark,...) are under free licences (GNU, MIT, Apache)
- **Margin** on hardware are not really valuable

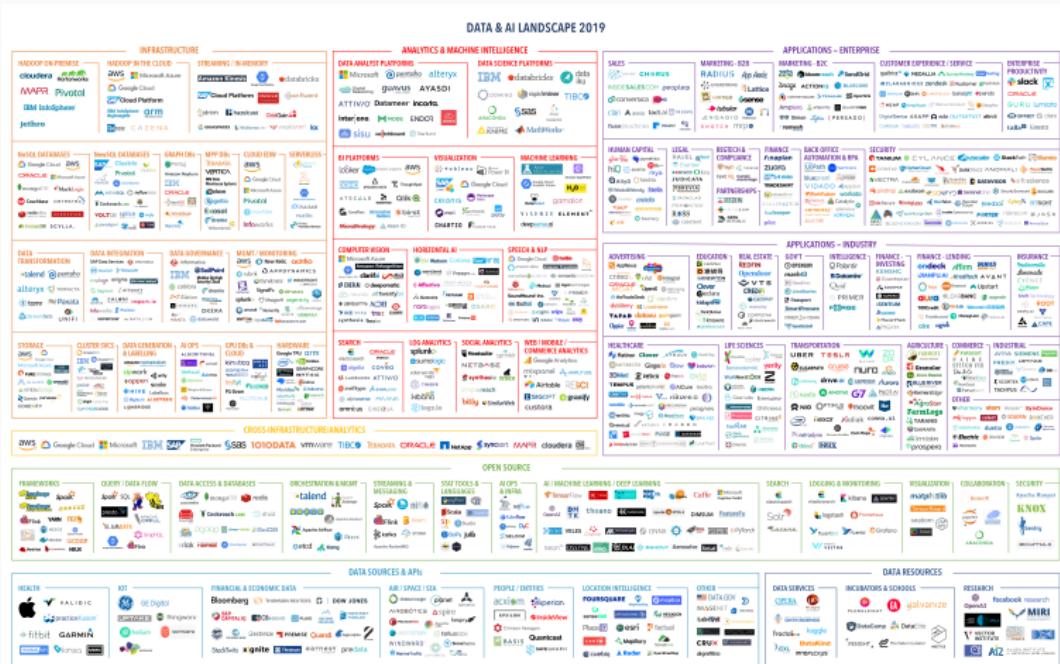
NEW STRATEGY : SELL SERVICES

- Enthought (**Canopy**), Continuum analytics (**Anaconda**), Horton Works, Cloudera (**Hadoop**, **Spark**...), Databricks (**Spark**), Oxdata (**H2O**), Revolution Analytics (**RHadoop**) – Microsoft

NEW MARKET : CLOUD COMPUTING (PICK-AND-SHOVEL PLAY)

- Actors : Amazon Web Service Platform, Microsoft Azure, Google Cloud Computing **aaaS**, Software **aaaS**, Service **aaaS**...
- A huge variety of services, Storage, Instances, ML Pipeline, DataFlow, etc....

BIG DATA & AI LANDSCAPE (TURCK, 2019)



July 16, 2019 - FINRA 2019 VERSION

© Matt Turck (@mattturck), Lisa Xu (@lisaxu92), & FirstMark (@firstmarkcap). mattturck.com/data2019

NEW (?) DATA SCIENCE

As a **STATISTICIAN**.

- Deep comprehension of algorithms
- Scientific watch

What's really new for a **DATA SCIENTIST**.

- Programming skills.
- Better understanding of all data-pipeline technology.
- Manage bigger data.
- New questions about ethics.

⇒ Learn to "self-learn".

ORGANIZATION & EVALUATION

OBJECTIVES - TOOLS

Tools you know

ML Python
Libraries



K Keras

gensim

surprise



Python
Environment



Viz' Python Libraries



seaborn



Framework
& Tool



Google Cloud Platform



git

OBJECTIVES - TOOLS

Tools you'll discover

ML Python
Libraries



K Keras

gensim

surprise



Python
Environment



Viz' Python Libraries



seaborn



Framework
& Tool



docker



Google Cloud Platform



git

OBJECTIVES - TOOLS

Tools you'll discover

ML Python
Libraries



K Keras

gensim

surprise



Python
Environment



Viz' Python Libraries



seaborn



Framework
& Tool



git

Question :

- What are the best tools to used? In which circumstances?
- Is it valuable to use CPU/GPU? Own laptop/cloud resources?
Distributed architecture?

OBJECTIVES - DATA & ALGORITHM

MACHINE LEARNING & HIGH DIMENSIONAL STATISTICS

DATA

- Structured Data
- Time Series & Signal
- Image

ALGORITHM

- Machine Learning : *SVM, RF, Regression.*
- Deep Learning : *MLP,CNN.*

OBJECTIVES - DATA & ALGORITHM

MACHINE LEARNING & HIGH DIMENSIONAL STATISTICS

DATA

- Structured Data
- Time Series & Signal
- Image

ALGORITHM

- Machine Learning : SVM, RF, Regression.
- Deep Learning : MLP, CNN.

ARTIFICIAL INTELLIGENCE FRAMEWORKS

DATA

- Text (*Cdiscount*)
- Video Game (*PacMan-like*)
- Note Data (*MovieLens*)

ALGORITHM

- Text cleaning : TFIDF, stemming.
- More Deep Learning : RNN, Words Embedding.
- Reinforcement learning : Policy Gradient / Q-learning.
- Recommendation System : NMF.

SIX - LABS

- **SESSION 1 - ??- ??-20**
 - Github Reminder, Python environment and Python Script.
 - Introduction to Google Cloud Computing.
- **SESSION 2 - ??- ??-20**
 - Containerization with Docker.
 - Natural language processing : Text Cleaning + Text Vectorization.
- **SESSION 3 - ??- ??-20**
 - Natural language processing : Words Embedding.
 - Natural language processing : Recurrent Network.
- **SESSION 4 - ??- ??-20**
 - Reinforcement learning : PG Gradient.
 - Reinforcement learning : Deep Q-learning.
- **SESSION 5 - ??- ??-20**
 - Recommendation System.
 - Data Visualization.
- **SESSION 6 - ??- ??-20**
 - "Free Time".

EVALUATION - OBJECTIVE

You will be evaluated on your capacity of acting like a Data Scientist,
i.e.

1. Understand an algorithm you haven't seen during course.
2. Being able to explain it clearly.
3. Make it run on an dataset to evaluate its performances.
4. Make it run on the appropriate tools (SPark? Cloud? GPU?)
5. Share it and make your results easily reproducible (Git - docker, conda environment.).

Notation

- Oral presentation - (40%)
- Project - (60%)

EVALUATION

ORAL PRESENTATION - (40%)

- Quality of the presentation. 25%
- In-Deep explanation of the chosen algorithm. 25%
- Choice of the tools-infrastructure used. 25%
- Results you obtained. 25%

Date : January ??, 2020.

EVALUATION

PROJECT - (60%) :

- The git should contain a clear markdown Readme, which describes : (33%)
 - Which result you achieved ? In which computation time ? On which engine ?
 - What do I have to install to be able to reproduce the code ?
 - Which command do I have to run to reproduce the results ?
- The code has to be easily reproducible. (33%)
 - Packages required has to be well described. (a *requirements.txt* files is the best)
 - Conda command or docker command can be furnish
- The code should be clear and easily readable. (33%)
 - Final results can be run in a script and not a notebook.
 - Only final code can be found in this script.

Deadline : January ?? 2020.

EVALUATION

⇒ SOME PROPOSITION OF SUBJECTS

- Image : Generative Adversarial Network.
- Image : Object segmentation & Localization
- Image : Siamese Network.
- One-Shot Learning.
- NLP : Recurrent Neural Network (Text generation)
- NLP : Words Embedding Comparison.
(Glove/Fasttext/Bert/Elmo)
- RL : Actor Critic.
- RL : Application on real data (Pacman).
- Google Cloud : Comparison of AutoML with other services.

⇒ YOU'RE FREE TO CHOOSE THE INFRASTRUCTURE YOU WANT.

⇒ GROUP OF 4 TO 5.

GITHUB REMINDER

GITHUB - CLONE

wikistat / AI-Frameworks

Unwatch 3 Star 22 Fork 11

Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security Insights Settings

Science des Données Saison 5: Technologies pour l'apprentissage automatique / statistique de données massives et l'Intelligence Artificielle Edit

apprentissage-machine spark-mllib scikitlearn-machine-learning data-science use-cases Manage topics

179 commits 1 branch 0 packages 0 releases 2 contributors MIT

Branch: master New pull request Create new file Upload files Find file Clone or download

File	Description	Last Commit
bguilouet Small update on readme		3 months ago
CloudComputing	MAH TP and slides GCE	3 months ago
NaturalLanguageProcessing	NLP TP ok	3 months ago
PySpark	Type on spark notebook	3 months ago
RecomendationSystem	Recommendation system	3 months ago
ReinforcementLearning	Recommendation system	3 months ago
slides	Recommendation system	3 months ago
.gitignore	Recommendation system	3 months ago
LICENSE	Initial commit	4 years ago
README.md	Small update on readme	1 hour ago

Clone with HTTPS Use SSH
Use Git or checkout with SVN using the web URL.
<https://github.com/wikistat/AI-Frameworks>

Open in Desktop Download ZIP 3 months ago

20/25

GITHUB - GET NEW COMMIT

LOCAL LAST COMMIT

```
~/INSA/AI-Frameworks master* ↓  
INSA > git log  
commit e8211da86b578d69641c8892dc969d1a44013d7d (HEAD -> master)  
Author: bguillouet <brendan.guillouet@gmail.com>  
Date:   Tue Feb 25 15:46:47 2020 +0100  
  
Commit with new organisation for 2020/2021 (WIP)
```

GITHUB

- ⌚ Commits on Mar 24, 2020
 - Small update on readme bguillouet committed 37 minutes ago🔗 bc96981 ↻
- ⌚ Commits on Feb 25, 2020
 - Commit with new organisation for 2020/2021 (WIP) bguillouet committed 28 days ago🔗 e8211da ↻

GITHUB - PULL

```
~/INSA/AI-Frameworks master* ↓ 10s
INSA > git status
On branch master
Your branch is behind 'origin/master' by 1 commit, and can be fast-forwarded.
  (use "git pull" to update your local branch)

Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git checkout -- <file>..." to discard changes in working directory)

        modified:   README.md

Untracked files:
  (use "git add <file>..." to include in what will be committed)

        new_file.txt

no changes added to commit (use "git add" and/or "git commit -a")

~/INSA/AI-Frameworks master* ↓
INSA > git pull origin master
From github.com:wikistat/AI-Frameworks
 * branch            master      -> FETCH_HEAD
Updating e8211da..bc96981
error: Your local changes to the following files would be overwritten by merge:
      README.md
Please commit your changes or stash them before you merge.
Aborting
```

GITHUB - COMMIT & PULL

```
~/INSA/AI-Frameworks master* ↓
INSA > git add README.md

~/INSA/AI-Frameworks master* ↓
INSA > git commit -m "my commit"
[master 7735af0] my commit
 1 file changed, 31 insertions(+), 10 deletions(-)

~/INSA/AI-Frameworks master* ↑↑
INSA > git pull origin master
From github.com:wikistat/AI-Frameworks
 * branch            master      -> FETCH_HEAD
Auto-merging README.md
CONFLICT (content): Merge conflict in README.md
Automatic merge failed; fix conflicts and then commit the result.
```

GITHUB - CREATE NEW BRANCH

```
~/INSA/AI-Frameworks master* ↓
INSA > git status
On branch master
Your branch is behind 'origin/master' by 1 commit, and can be fast-forwarded.
  (use "git pull" to update your local branch)

Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git checkout -- <file>..." to discard changes in working directory)

        modified:   README.md

Untracked files:
  (use "git add <file>..." to include in what will be committed)

        new_filte.txt

no changes added to commit (use "git add" and/or "git commit -a")

~/INSA/AI-Frameworks master* ↓
INSA > git checkout -b TP1
M      README.md
Switched to a new branch 'TP1'

~/INSA/AI-Frameworks TP1*
INSA > git add README.md

~/INSA/AI-Frameworks TP1*
INSA > git commit -m "my commit"
[TP1 6a37c81] my commit
```

GITHUB - PULL ON MASTER

```
~/INSA/AI-Frameworks TP1*
INSA > git checkout master
Switched to branch 'master'
Your branch is behind 'origin/master' by 1 commit, and can be fast-forwarded.
  (use "git pull" to update your local branch)

~/INSA/AI-Frameworks master* ↓
INSA > git pull origin master
From github.com:wikistat/AI-Frameworks
 * branch            master      -> FETCH_HEAD
Updating e8211da..bc96981
Fast-forward
 README.md | 11 ++++++----
 1 file changed, 8 insertions(+), 3 deletions(-)

~/INSA/AI-Frameworks master*
INSA > git log
commit bc96981b978d6f12693d64aa2e4cbdf8da945160 (HEAD -> master, origin/master, origin/HEAD)
Author: bguillouet <brendan.guillouet@gmail.com>
Date:   Tue Mar 24 16:05:06 2020 +0100

    Small update on readme

commit e8211da86b578d69641c8892dc969d1a44013d7d
Author: bguillouet <brendan.guillouet@gmail.com>
Date:   Tue Feb 25 15:46:47 2020 +0100

    Commit with new organisation for 2020/2021 (WIP)
```