

Reproductibilité en apprentissage automatique

David Beauchemin

Département d'informatique et de génie logiciel,
Université Laval

david.beauchemin.5@ulaval.ca

30 octobre 2020



Objectifs de la présentation

- Sensibiliser sur les enjeux de la reproductibilité.
- Inciter l'intégration des solutions permettant une meilleure reproductibilité dans vos solutions d'affaires ou académiques.

Mes qualifications



DAVID BEAU- CHEMIN

Candidat au doctorat

Département d'informatique et de génie logiciel

- Introduit (informellement) à la recherche reproductible en 2016 (RMarkdown et Git).
- Participation à REPROLANG, visant la reproductibilité d'articles scientifiques ayant mené à la publication d'un article en 2020 (où nous abordons les notions présenté ici) [?].
- Membre actif dans le développement de solution d'intégration facilitant la reproductibilité (Poutyne, MLFlow callback).

1 Introduction

2 Les barrières à la répliquabilité

3 Ok, mais comment ?

4 La suite

C'est quoi la reproductibilité ?

La reproductibilité est le principe qu'on ne peut tirer de conclusions que d'un événement bien décrit, qui est apparu plusieurs fois, provoqué par des **personnes différentes**.

Toutefois, on utilise souvent ce terme pour spécifiquement désigné la **réplicabilité**. Soit la réplication (reproduction) des résultats d'un articles dans des environnements pas (toujours) différents [?, ?].

- Être capable de répliquer les résultats d'un article/ d'un projet,
- à partir du même jeux de données ou un jeux de données différents (mais proche),
- en utilisant la procédure d'entraînement de l'article ou en utilisant notre procédure d'entraînement et
- en utilisant le code du projet.

Pourquoi s'y intéressé ?

- 70 % des chercheurs en science on échoué dans leur tentative de reproduire un article d'un autre chercheur,
- 50 % n'on pas réussi à reproduire leur **propre** expérimentations [?].

Pourquoi s'y intéressé ?

L'informatique ne fais pas exception à cela malgré la simplicité (théorique) de réplication des résultats. Selon une étude, sur 255 articles près de 40 % n'était pas répliquable [?].

Pourquoi s'y intéressé ?

La répliquabilité du code et d'un article facilite la réutilisation pour d'autres projets de recherche **et** le transfert vers l'industrie.

1 Introduction

2 Les barrières à la répliquabilité

3 Ok, mais comment ?

4 La suite

- Non disponibilité du jeux de données ou version (pas clair) du jeux de données,
- mauvaise spécification ou sous-spécification du modèle ou de la procédure de formation,
- manque de disponibilité du code nécessaire pour exécuter les expériences, ou erreurs dans le code,
- configuration du modèle déficiente [?]¹.

1. Liste sélective de ceux présenté dans l'article.

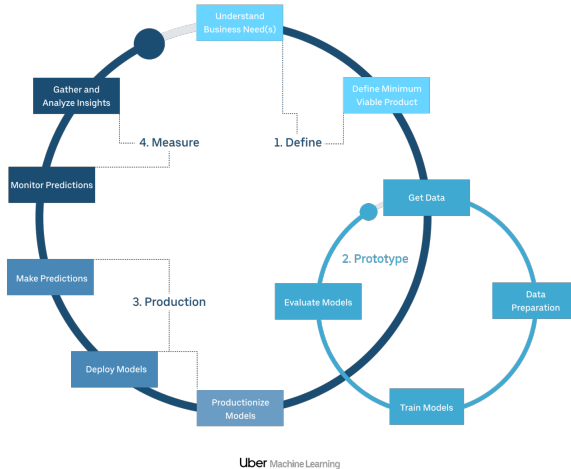


FIGURE – From Uber Engineering²

2. <https://eng.uber.com/scaling-michelangelo/>

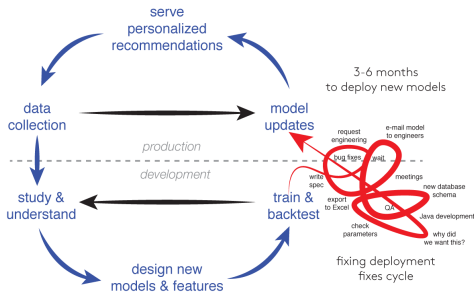


FIGURE – The need for Agile machine learning³

3. <https://johann.schleier-smith.com/blog/2015/08/09/need-for-agile-machine-learning.html>

Plan

1 Introduction

2 Les barrières à la répliquabilité

3 Ok, mais comment ?

4 La suite

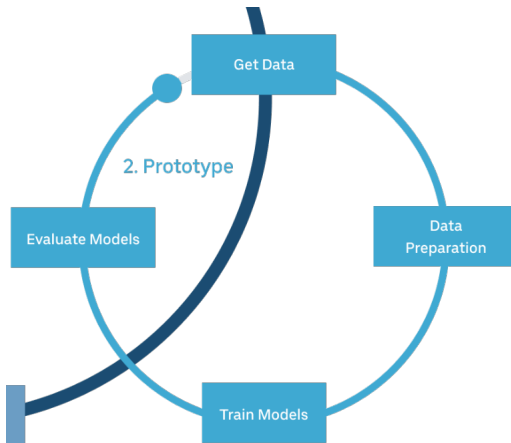


FIGURE – From Uber Engineering⁴

4. <https://eng.uber.com/scaling-michelangelo/>

- Avec quelle version des données on travail ?
- Comment gérer **facilement** plusieurs version des données ?
- Comment définir **facilement** les étapes de pré-processing des données ?

Il nous faut des *data pipelines*, des tuyaux que nous pouvons raccorder **facilement** à nos modèles pour l'entraînement et la mise en production, par exemple, Data Version Control (DVC) ⁵.

5. <https://dvc.org/>

- Avec qu'elle version du code on travail ?
- Comment savoir **rapidement** qu'elle est la différence d'implémentation entre deux versions du modèles ?
- Comment gérer **facilement** les embranchements d'expérimentations ?

Il nous faut un outil nous permettant de visualiser la différence entre des fichiers de code et nous permettant d'avoir plusieurs version du code pouvant exister **en même temps**, par exemple Git⁶.

6. <https://git-scm.com/>

- Ne pas réinventer la roue.
- Simplifier l'écriture de code pour développer des modèles.
- Qui facilite l'entraînement (GPU, multi-GPU/CPU).

Il nous faut des outils nous permettant de simplifier le développement de nos modèles, par exemple, Poutyne [?], PyTorch Lightning [?], Scikit-Learn [?], Gensim [?] et Allen NLP [?].

Entraînement, configuration et résultats

- Avec quelle version du code, du modèle et des données avons-nous fait cet entraînement ?
- Quels sont les résultats ?
- Comment visualiser **rapidement** les résultats et les paramètres de configuration ?

Il nous faut des outils nous permettant de *logger* les paramètres d'entraînement et les résultats, par exemple, MLFlow [?] et Sacred [?].

- Comment créer des tableaux de résultats **facilement** (pas à la *mitaine*) ?
- Comment s'assurer **facilement** que les résultats sont à jours ?
- Comment visualiser **rapidement** les résultats et les paramètres de configuration ?

Il nous faut des outils nous permettant de créer des tableaux de résultats à même les résultats, soit de diminuer le plus possible le travail manuel, par exemple, Python2latex⁷ et Markdown⁸.

7. <https://github.com/jsleb333/python2latex>

8. <https://fr.wikipedia.org/wiki/Markdown>

- Comment s'assurer que nos modèles fonctionnent sur d'autres environnements ?
- Comment faciliter la réutilisation de notre code ?

Docker !⁹

9. <https://www.docker.com/>

Plan

1 Introduction

2 Les barrières à la répliquabilité

3 Ok, mais comment ?

4 La suite

Développer des processus rigoureux (par essai et erreur) et ne pas prendre tout ce qui a été discuter ici comme l'unique solution.

- Clean code¹⁰
- Dagger de projet¹¹

10. <https://www.oreilly.com/library/view/clean-code-a/9780136083238/>

11. [https:](https://github.com/davebulaval/cookiecutter-machine-learning-research)

[//github.com/davebulaval/cookiecutter-machine-learning-research](https://github.com/davebulaval/cookiecutter-machine-learning-research)

References I
