



**HAL**  
open science

## Évaluation de systèmes apprenant tout au long de la vie

Yevhenii Prokopalo, Sylvain Meignier, Olivier Galibert, Loïc Barrault,  
Anthony Larcher

### ► To cite this version:

Yevhenii Prokopalo, Sylvain Meignier, Olivier Galibert, Loïc Barrault, Anthony Larcher. Évaluation de systèmes apprenant tout au long de la vie. 6e conférence conjointe Journées d'Études sur la Parole (JEP, 33e édition), Traitement Automatique des Langues Naturelles (TALN, 27e édition), Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RÉCITAL, 22e édition), 2020, Nancy, France. pp.516-524. hal-02798580v2

**HAL Id: hal-02798580**

**<https://hal.science/hal-02798580v2>**

Submitted on 18 Jun 2020 (v2), last revised 23 Jun 2020 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Évaluation de systèmes apprenant tout au long de la vie

Yevhenii Prokopalo<sup>1</sup> Sylvain Meignier<sup>1</sup>

Olivier Galibert<sup>2</sup> Loïc Barrault<sup>3</sup> Anthony Larcher<sup>1</sup>

(1) LIUM, 72 Le Mans, France, (2) LNE, 78 Trappes, France, (3) University of Sheffield, UK

yevheniiprokopalo@univ-lemans.fr

## RÉSUMÉ

---

Aujourd'hui les systèmes intelligents obtiennent d'excellentes performances dans de nombreux domaines lorsqu'ils sont entraînés par des experts en apprentissage automatique. Lorsque ces systèmes sont mis en production, leurs performances se dégradent au cours du temps du fait de l'évolution de leur environnement réel. Une adaptation de leur modèle par des experts en apprentissage automatique est possible mais très coûteuse alors que les sociétés utilisant ces systèmes disposent d'experts du domaine qui pourraient accompagner ces systèmes dans un apprentissage *tout au long de la vie*. Dans cet article nous proposons un cadre d'évaluation générique pour des systèmes apprenant tout au long de la vie (SATLV). Nous proposons d'évaluer l'apprentissage assisté par l'humain (actif ou interactif) et l'apprentissage au cours du temps.

## ABSTRACT

---

### Evaluation of lifelong learning systems

Current intelligent systems need the expensive support of machine learning experts to sustain their performance level when used on a daily basis. To reduce this cost, i.e. remaining free from any machine learning expert, it is reasonable to implement lifelong (or continuous) learning intelligent systems that will continuously adapt their model when facing changing execution conditions. In this work, the systems are allowed to refer to human domain experts who can provide the system with relevant knowledge about the task. Nowadays, the fast growth of lifelong learning systems development rises the question of their evaluation. In this article we propose a generic evaluation methodology for the specific case of lifelong learning systems. Two steps will be considered. First, the evaluation of human-assisted learning (including active and/or interactive learning) outside the context of lifelong learning. Second, the system evaluation across time, with propositions of how a lifelong learning intelligent system should be evaluated when including human assisted learning or not.

---

**MOTS-CLÉS :** Apprentissage automatique, évaluation, apprentissage tout au long de la vie.

**KEYWORDS:** Machine learning, Lifelong learning, Evaluation.

---

## 1 Introduction

Les systèmes intelligents utilisent une représentation du monde, un modèle, dont l'apprentissage est réalisé en laboratoire par des experts en apprentissage automatique (EAA) (Bishop, 2006). Le rôle de ces experts est triple : (1) ils sélectionnent et annotent les données d'apprentissage ; (2) ils déterminent les méta-paramètres inhérents à la structure du système en utilisant des données de développement ;

(3) ils évaluent et comparent les systèmes afin de déterminer lequel mettre en production sur un jeu de données de test. Dans cet article, nous appellerons "données initiales" l'ensemble de ces trois jeux de données. Une fois en production, ces systèmes sont confrontés aux données réelles. Le cycle de vie d'un tel système est illustré sur la figure 1-A (Chen & Liu, 2016). Si il est raisonnable de considérer que les données d'entrée du système lors de sa mise en production sont proches des données initiales, il arrive la plupart du temps qu'elles s'en éloignent avec le temps. Cet effet est bien connu et entraîne de sévères dégradations de performances (Quionero-Candela *et al.*, 2009) qui ne peuvent être résolues que grâce à l'intervention d'un expert qui entraînera un nouveau modèle. Afin d'éviter le recours coûteux à cet expert, la communauté scientifique se mobilise pour développer des systèmes qui apprennent seuls au cours du temps, sans l'aide d'un expert en apprentissage automatique mais en adaptant continuellement leur modèle à l'aide des données disponibles et la possible intervention d'un expert humain du domaine. Ces systèmes sont appelés des systèmes apprenant tout au long de la vie (SATLV). Les SATLV diffèrent des systèmes à apprentissage statique car ils disposent de modules d'adaptation qui doivent permettre de maintenir leurs performances au cours du temps en apprenant des nouvelles données d'entrée. Le cycle de vie d'un système SATLV est décrit par la figure 1-B.

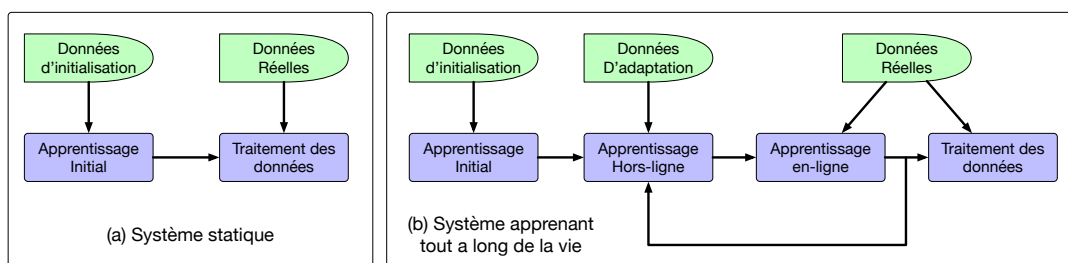


FIGURE 1: Comparaison des cycles de vie d'un système à apprentissage statique et d'un système apprenant tout au long de la vie.

Les SATLV peuvent adapter leur modèle selon deux modes : en-ligne et hors-ligne. Dans ce travail, nous définirons l'adaptation en-ligne comme l'action d'un système qui met à jour ses connaissances (le modèle) en apprenant des données d'entrée qu'il doit traiter. En d'autres termes, il est possible pour le système d'utiliser toutes les données d'entrée qui lui sont fournies lors de son exploitation. Dans ce cas, le SATLV reçoit les mêmes données que n'importe quel système à apprentissage statique.

L'adaptation hors-ligne permet à un SATLV d'utiliser l'ensemble des données qui lui sont accessibles pour améliorer son modèle. Ces données incluent entre autre les données initiales, que le système peut décider de réutiliser, des données additionnelles qui peuvent être fournies par l'expert humain du domaine, des données collectées automatiquement par le SATLV et même des données réelles que le SATLV a dû traiter par le passé.

Lors d'une adaptation hors-ligne, le SATLV n'a pas à générer d'hypothèse (de sortie) et peut bénéficier de ressources de calcul importantes. Dans la suite de cet article, nous supposons que les données d'adaptation en-ligne sont fournies sans aucune annotation tandis que les données d'adaptation hors-ligne peuvent inclure des données annotées et des données non-annotées.

Les processus d'adaptation en-ligne et hors-ligne peuvent utiliser des méthodes d'adaptation non-supervisées ou assistées par l'humain. L'apprentissage assisté par l'humain inclut l'apprentissage actif (AA ; pour lequel le système prend l'initiative d'engager un échange avec l'humain) et de l'apprentissage interactif (AI ; dans lequel l'humain est à l'initiative de cet échange).

L'apprentissage non-supervisé permet au système d'adapter son modèle en exploitant les données

d'entrée sans annotation. L'AA autorise le système à demander à l'humain de corriger une partie de l'hypothèse qu'il a générée et l'AI permet à l'humain de renvoyer au système des corrections de son hypothèse courante afin que le système intègre ces corrections et génère une nouvelle hypothèse.

Les SATLV, comme les systèmes à apprentissage statique doivent être évalués avant leur mise en production. Cette évaluation prend souvent la forme d'une comparaison entre systèmes respectant un protocole strict et contrôlé. Notons que le présent article ne traite que de l'évaluation des systèmes automatiques et non de leur développement.

L'évaluation des systèmes à apprentissage statique suit une méthodologie bien établie et requiert un jeu de données d'évaluation, une métrique objective et un protocole d'évaluation. Pour de très nombreuses tâches, ces éléments ne sont pas disponibles pour l'évaluation de systèmes SATLV et la dérivation des modalités d'évaluation des systèmes statiques pour permettre l'évaluation de SATLV effectuant les mêmes tâches n'est pas simple. L'évaluation des SATLV nécessite de pouvoir évaluer les performances du système au cours du temps ainsi que les performances des modules d'AA et d'AI individuellement. Cette évaluation doit également prendre en compte la chronologie des données. En effet, si les systèmes automatiques traitent la plupart du temps les données actuelles, il est possible que pour certaines tâches spécifiques, ce système ait à traiter des données issues du passé. Les performances de ce système sur des données actuelles ou des données du passé peuvent avoir une incidence différente pour l'utilisateur. Nous introduisons dans la section 3.1 une politique-utilisateur que la métrique d'évaluation doit prendre en compte pour évaluer la capacité du SATLV à s'adapter selon la politique dictée par l'utilisateur.

Le travaux existant pour l'évaluation de tels systèmes font souvent état d'un manque de protocoles et de données nécessaires et proposent de simuler la chronologie ou l'évolution des données en permutant les données de corpus existant (Kemker *et al.*, 2018; Lomonaco & Maltoni, 2017). Il apparaît également que les performances de tels systèmes sont évaluées à plusieurs instant dans le temps ou sur des données présentant de nouveaux événements mais qu'aucune évaluation de leur comportement temporel n'est rapportée (Parisi *et al.*, 2019). Dans cet article<sup>1</sup>, nous proposons des métriques et protocoles pour l'évaluation de l'apprentissage assisté par l'humain pour les cas d'apprentissage hors-ligne et en-ligne (sections 2.1 et 2.2) (Prokopalo *et al.*, 2020). Nous proposons également une métrique pour évaluer la capacité des SATLV à adapter leur modèle pour respecter la politique fixée par l'utilisateur au cours du temps. Une métrique d'évaluation de l'apprentissage assisté par l'humain est introduite dans la section 2. Conscients que les métriques d'évaluation et les protocoles sont toujours limités et ne permettent pas d'évaluer tous les aspects de systèmes complexes, nous discutons des limitations de nos propositions dans la section 4

L'ensemble des propositions de cet article considèrent que les systèmes à apprentissage statiques sont évalués par des métriques similaires à des taux d'erreurs (le plus faible le meilleur) mais il est important de noter que les solutions proposées peuvent être appliquées à toute métrique scalaire sans perdre de leurs généralités.

## 2 Évaluation de l'apprentissage assisté par l'humain

Pour supprimer la nécessité de faire intervenir un expert en apprentissage automatique, il est raisonnable de faire appel à un expert humain du domaine (EHD). Dans cette section, nous ne considérons

---

1. ce travail a fait l'objet d'une publication à LREC 2020

pas l'évaluation au cours du temps mais seulement l'apprentissage assisté par l'humain. L'interaction entre l'EHD et le SATLV peut avoir lieu hors-ligne ou en-ligne.

Hors-ligne, l'EHD interagit avec le SATLV sur les données disponibles sans considérer de contrainte de temps ou de calcul. Dans ce contexte, les SATLV sont libres d'apprendre de n'importe quelles données non-annotées et ont l'opportunité de poser à l'EHD un nombre limité de questions à propos de ces données. Durant cet apprentissage actif, le système doit poser les questions qui maximisent la généralisation des réponses apportées par l'utilisateur. La section 2.1 décrit une façon d'évaluer l'apprentissage actif.

En-ligne, l'EHD est partie prenante de la chaîne de production, ses interactions avec le SATLV ont lieu entre l'arrivée des données à traiter et l'envoi de l'hypothèse finale à l'utilisateur. Dans ce travail, nous restreignons le cadre d'interaction de l'EHD aux seules données à traiter (sans considérer que l'EHD peut apporter une information sur des données additionnelles). Étant donné un ensemble  $X$  de données à traiter, le SATLV peut poser des questions à l'EHD. Il est également possible pour l'EHD de surveiller le traitement automatique et de fournir des informations au système au cours du processus de traitement.

Afin d'évaluer les processus d'adaptation hors-ligne et en-ligne, l'EHD doit être simulé par un système qui assure la reproductibilité des tests.

## 2.1 Apprentissage hors ligne assisté par l'humain

L'efficacité de l'apprentissage assisté par l'humain réside dans la capacité du système à exploiter au mieux les interactions avec l'humain pour en minimiser le coût. De nombreux articles dans la littérature conseillent de mesurer la qualité de l'apprentissage assisté par l'humain en fonction du coût des interactions et des performances du système (Krogh & Vedelsby, 1995; Siddhant & Lipton, 2018; Drugman *et al.*, 2019; Beluch *et al.*, 2018; Pérez-Dattari *et al.*, 2018; Celemin & Ruiz-del Solar, 2019). Cette mesure nous semble pertinente, ainsi nous ne proposons pas de nouveauté dans ce domaine et nous contentons ici de décrire cette métrique. Étant donné un modèle initial, un processus itératif est initié par le système avec pour but de réduire le taux d'erreurs obtenu sur un jeu de données. Le système peut poser des questions relatives à n'importe quel document auquel il a accès et ses performances sont évalués sur un jeu de données de test qui ne peut pas être utilisé pour adapter son modèle. Dans ce contexte, il n'est pas certain que le taux d'erreurs diminue et atteigne zéro au cours de l'apprentissage. Pour cette raison, un coût maximum d'interaction est fixé. Une fois atteint ce coût, l'apprentissage actif est interrompu. Il est important de noter que lors d'un apprentissage hors-ligne, le système n'a pas à traiter de données en particulier et l'apprentissage assisté par l'humain se trouve donc réduit à l'apprentissage actif. Le cas où le EHD initie l'interaction peut être vu comme de l'apprentissage faiblement supervisé ou supervisé.

## 2.2 Apprentissage en ligne assisté par l'humain

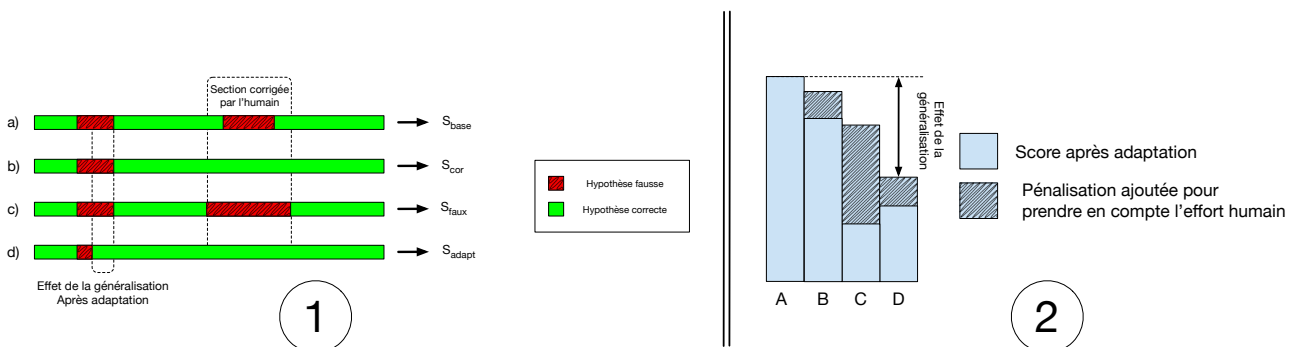
Lors d'un apprentissage en-ligne, le système doit traiter une séquence de données réelles et renvoyer les hypothèses correspondantes. En fonction du scénario, deux types d'interactions sont possibles : soit le SATVL initie une session d'apprentissage actif avant de renvoyer l'hypothèse finale, soit l'EHD initie une session d'apprentissage interactif en corrigeant itérativement les hypothèses générées par le système.

Dans le cas d'un apprentissage interactif, il est raisonnable de considérer que l'EHD est en mesure de

fournir des corrections jusqu'à obtenir une hypothèse entièrement correcte. L'efficacité de l'apprentissage assisté par l'humain peut alors être mesurée par le coût de supervision minimal présenté dans (Geoffrois, 2016).

Différents systèmes peuvent autoriser différentes modalités d'interaction. Il est donc difficile d'estimer le coût de l'interaction humaine de façon homogène et comparable entre systèmes. (Broux *et al.*, 2018). Dans l'idéal, une mesure de performance produirait une comparaison équitable de systèmes offrant des modalités d'interaction diverses. Nous proposons pour ceci de mesurer le coût d'interaction sous forme d'une métrique scalaire donnée dans la même unité que le score de performance du système (taux d'erreur par exemple). Cette pénalité peut s'appliquer à l'apprentissage actif ou interactif. Calculer une pénalité dans la même unité que la performance du système permet de calculer un unique score pénalisé qui reflète à la fois la performance finale du système et l'effort humain nécessaire pour l'obtenir.

Une première option consiste à calculer le coût d'interaction comme la quantité d'information donnée au système par l'utilisateur. Cependant, cette stratégie ne fonctionne pas pour des métriques non linéaires comme par exemple BLEU pour la traduction automatique (Papineni *et al.*, 2002). Nous proposons de calculer la pénalité comme la part du score correspondant aux données corrigées par l'humain. Afin de calculer cette quantité, nous calculons deux scores : un score corrigé,  $S_{cor}$ , et un score dégradé,  $S_{faux}$ , dont le calcul est décrit par la figure 2-1.



Une hypothèse (a), produite par un système automatique contenant une partie correcte (vert) et des erreurs (rouge) obtient un score  $S_{base}$ . Lors de l'apprentissage assisté par l'humain, une section de l'hypothèse (partiellement fautive) est corrigée. Pour calculer l'effort fourni par l'humain, un nouveau score  $S_{cor}$  est calculé immédiatement après application de la correction (b). Nous introduisons un score  $S_{faux}$  qui est calculé en remplaçant la section corrigée par une hypothèse entièrement fautive (c). Après avoir reçu une information provenant de l'humain, le système adapte son modèle et produit une nouvelle hypothèse qui obtient un score  $S_{adapt}$ . Un système capable de généraliser l'information reçue devrait améliorer l'hypothèse finale sur des portions où l'utilisateur n'est pas intervenu manuellement (d).

Effet de la pénalisation dans différents cas. Le score A,  $S_{base}$ , est obtenu sur l'hypothèse initiale, avant toute intervention de l'humain. Les trois autres colonnes représentent des scores pénalisés,  $S_{pen}$ , obtenus en utilisant différentes méthodes d'apprentissage assisté par l'humain. La partie uniforme des colonnes B, C et D représente le score finale après adaptation,  $S_{adapt}$ . La partie hachurée représente la pénalité. Le score (B) illustre le cas où le système exploite peu l'information apportée par l'humain. Dans le cas (C), le système obtient un gain significatif au prix d'un effort humain important. Dans le cas (D), le système obtient un taux d'erreur final plus élevé que pour (C) mais l'effort humain nécessaire est beaucoup plus faible. Notons que la différence entre  $S_{base}$  et  $S_{pen}$  (indiqué pour le score D) correspond au gain obtenu par le système qui a généralisé les corrections appliquées par l'humain.

FIGURE 2: Pénalisation de l'apprentissage assisté par l'humain.

Un système produit une première hypothèse (Figure 2-1-A) et obtient un score  $S_{base}$  avant toute interaction avec l'EHD. L'EHD fournit alors des informations en corrigeant une partie de l'hypothèse courante. La partie des données corrigée est indiquée sur la figure 2-1-B et cette correction permet d'obtenir un score  $S_{cor}$ . En fonction de la tâche, il est possible qu'une partie de l'hypothèse qui est corrigée ait été en partie correcte. (Dans le cas de la traduction il peut s'agir de la correction d'une phrase qui contenait déjà certains mots corrects). La différence entre  $S_{base}$  et  $S_{cor}$  correspond à l'amélioration apportée par la seule correction mais ne reflète pas le coût de la correction. Nous calculons alors un nouveau score  $S_{faux}$  qui est obtenu en remplaçant l'intégralité de la partie corrigée de l'hypothèse par une hypothèse erronée. (Figure 2-1-C) La différence  $S_{faux} - S_{cor}$  mesure la part

du score qui correspond au coût réel de la correction. Finalement, l'hypothèse corrigée est renvoyée au système qui adapte son modèle et re-génère une nouvelle hypothèse qui obtiendra le score  $S_{adapt}$ . Le scores pénalisé,  $S_{pen}$ , est obtenu de la façon suivante :

$$S_{pen} = S_{adapt} + (S_{faux} - S_{cor}) \quad (1)$$

Notons qu'un système qui ne prendrait pas la correction en compte ou demanderait une information déjà correcte serait pénalisé deux fois : une fois par la valeur élevée de  $S_{adapt}$  et une fois par le second terme de l'équation 1.

L'effet de cette pénalisation est illustré par la figure 2-2. Le score  $S_{base}$ , illustré par le score A est obtenu avant toute interaction humaine. À partir de cette hypothèse, différentes stratégies d'adaptation pourraient mener à différents scores pénalisés, illustrés par les barres B, C et D de la figure 1. Les parties uniformes de ces barres illustrent le score  $S_{adapt}$  tandis que les parties hachurées illustrent la pénalité appliquée à chaque version du système. Le système idéal obtiendra un score pénalisé très faible tout en nécessitant une petite quantité d'interactions.

### 3 Évaluation au cours du temps

Une fois le modèle initial,  $M_0$ , d'un système SATLV entraîné et testé en laboratoire, le système est mis en production et traite une séquence de données,  $[X_{l_i}]_{i=1...T}$ , lui parvenant au cours du temps. Cette séquence est ici nommée *données au cours du temps* et  $l_i$  est le temps auquel  $X_{l_i}$  a été produit. Le SATLV peut adapter son modèle après avoir traité des données  $X_{l_i}$  pour obtenir une nouvelle version  $M_i$ . Les performances du SATLV sont évaluées sur les données  $X_{l_i}$  et la séquence de scores obtenue est  $[S(M_i, l_i)]_{i=1...T}$ .

#### 3.1 Prise en compte d'une politique fixée par utilisateur

Les données que les SATLV doivent traiter évoluent au cours du temps mais il arrive qu'un système rencontre des données semblables à des données vues dans le passé. L'adaptation du modèle aux données actuelles peut entraîner « l'oubli catastrophique » des données du passé au sens de (French, 1999). Le coût de cet oubli pour l'utilisateur dépend de l'application visée et il est essentiel que l'utilisateur puisse l'indiquer au SATLV afin d'adapter la politique d'adaptation du SATLV. Différentes politiques sont illustrées sur la figure 3

Pour l'évaluateur du système, il est essentiel d'évaluer la capacité du SATLV à adapter son modèle de façon à répondre au mieux aux exigences de l'utilisateur. Afin d'évaluer la capacité du SATLV à suivre une politique d'adaptation donnée, une nouvelle séquence de données,  $[Y_{t_j}]_{j=1...N}$ , est utilisée. La date de création de ces données doit couvrir la même période que les données sur lesquelles le système a été adapté. Il faut que :  $l_1 < t_1 < t_N < l_T$ . Ainsi, chaque version  $M_i$  du modèle peut être évaluée sur l'ensemble de la séquence de test  $Y_{t_j}$  pour obtenir une séquence de scores  $[s(M_i, t_j)]_{j=1...N}$ . Idéalement, ces séries temporelles de données devraient être strictement alternées mais de nombreux critères sur la nature des données et les événements rencontrés au cours du temps doivent être pris en compte afin de produire des études objectives. Notons que cette séquence de scores est calculée pour une unique version,  $M_i$ , du modèle. Les performances de ce système sur les données antérieures à  $l_i$  montrerons à quel point le système a oublié le passé. La figure 3-1 illustre la

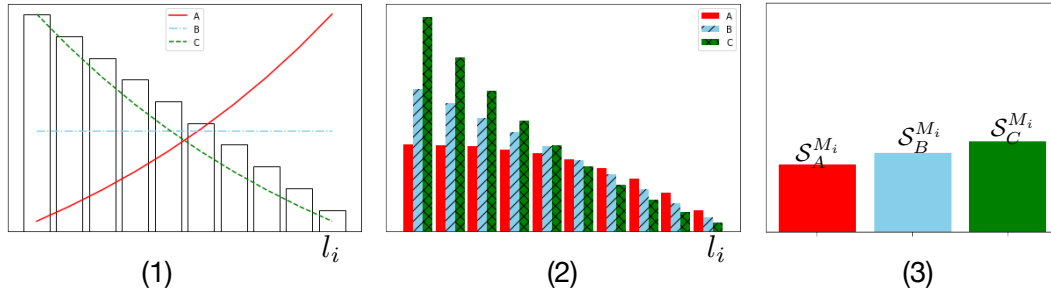


FIGURE 3: Scores obtenus par un système apprenant au cours du temps et 3 politiques d'adaptation différentes (1). Scores pondérés au cours du temps (2). Score final résultat de la somme pondérée des scores au cours du temps pour les 3 politiques (3).

séquence de scores obtenue par la version  $M_i$  du système dont le modèle a été adapté jusqu'à  $l_i$ . Cet exemple factice illustre le cas d'un système optimisé pour les données les plus récentes mais dont les performances sur les données du passé sont moins bonnes. La même figure montre trois politiques d'adaptation : (A) l'utilisateur désire que le système privilégie les données actuelles ; (B) les données du passé sont aussi importantes que les données actuelles ; (C) l'utilisateur privilégie les données du passé. Évaluer la politique d'adaptation pour ce système consiste à pondérer les scores obtenus au cours du temps par la politique définie par l'utilisateur (Figure 3-2) La somme de ces scores pondérés est un indicateur des performances du système mais également de sa capacité à satisfaire la politique d'adaptation,  $\mathcal{P}$ , fixée.

$$S_{\mathcal{P}}^{M_i} = \sum_{t_j \in [l_1; l_i]} s(M_i, t_j) \cdot \mathcal{P}_{t_j} \quad (2)$$

On observe sur la figure 3-3, que le système factice obtient un meilleur score pour la politique qui privilégie les données actuelles, comme attendu.

### 3.2 Évaluation au cours du temps de l'apprentissage assisté par l'humain

Dans ce travail nous considérons que l'adaptation hors-ligne au cours du temps représente un coût fixe pour le système car elle nécessite la présence permanente d'un expert humain, mais ne modifie pas la disponibilité du système. En revanche, l'adaptation en-ligne influe sur le temps de réponse du système qui effectue des cycles d'adaptation entre la réception des données et la communication de l'hypothèse au client. Ainsi nous proposons d'intégrer la coût de l'adaptation en-ligne en remplaçant, dans l'équation 2, le score  $s(M_i, t_j)$  par sa version pénalisée proposée dans la section 2.2. Le score pénalisé et pondéré par la politique d'adaptation devient alors :

$$S_{\mathcal{P}}^{M_i} = \sum_{t_j \in [l_1; l_i]} S_{pen}(M_{i,j}, t_j) \cdot \mathcal{P}_{t_j} \quad (3)$$

## 4 Discussion

Dans cet article, nous proposons trois contributions pour évaluer les systèmes automatiques apprenant tout au long de la vie (SATLV). Nous proposons une nouvelle mesure du coût de l'interaction humaine dans l'apprentissage assisté par l'humain, donnée dans la même unité que la mesure de performance du système. Nous introduisons le concept de politique d'adaptation afin d'évaluer la capacité des



SATLV à respecter une politique définie par l'utilisateur. Enfin, nous avons combiné ces éléments pour proposer un cadre complet d'évaluation des systèmes SATLV. Notre travail considère des métriques de type « taux d'erreur » mais peut également être appliqué à des métriques scalaires définies sur un intervalle semi-ouvert.

La pénalité proposée pour sanctionner le coût de l'interaction humaine dans l'adaptation du modèle est exprimée dans la même unité que la mesure de performance du système afin d'être directement combinée et d'offrir un indicateur de performance unique. Cette pénalité permet d'évaluer la capacité du système à généraliser l'information reçue de l'expert du domaine. Dans certains cas, le calcul du score intermédiaire,  $S_{faux}$ , pose des questions complexes. Par exemple dans le cas de la traduction automatique, la définition de la traduction la plus fautive est très complexe car il est toujours possible de dégrader le score en insérant des mots. Une borne raisonnable consiste à limiter la taille de l'hypothèse fautive à la taille de la référence. Ce calcul sera appliqué dans une tâche de l'évaluation internationale WLT 2020 dédié aux SATLV. Le calcul de la pénalité considère que l'information fournie par l'humain peut être directement appliquée à l'hypothèse pour calculer un nouveau score,  $S_{cor}$ . Cette hypothèse n'est pas toujours réaliste car elle dépend des interactions autorisées.

Afin de garantir la reproductibilité et l'équité de l'évaluation de l'adaptation assistée par l'humain, il est nécessaire d'implémenter une simulation d'expert humain. Différentes implémentations sont possibles et tous les systèmes comparés devront interagir avec le même simulateur. Idéalement, plusieurs simulateurs pourraient être utilisés afin de confronter les systèmes à différents cas de figure. Le développement de ces simulateurs est un sujet de recherche en soi qui n'est pas l'objet de cet article.

Notre deuxième contribution est l'introduction d'une politique d'adaptation. Cette politique permet à l'utilisateur d'introduire une connaissance du domaine relative à l'évolution des données au cours du temps (par exemple le cas où un système rencontrera des données qui varient de façon cyclique). La définition d'une telle fonction par un utilisateur non-expert en apprentissage automatique nécessite une interface utilisateur intuitive et bien définie qui peut être difficile à évaluer en elle-même.

Dans ce travail, nous avons tenté de définir un cadre général d'évaluation qui pourrait être transposé à un grand nombre de tâches et de mesures de performance. Comme c'est le cas pour de nombreuses métriques existantes et utilisées, nos propositions ne répondent pas à tous les besoins et doivent être combinées à d'autres afin d'évaluer tous les aspects de systèmes complexes par définition.

Enfin les métriques et protocoles définis dans ce travail ont été appliqués au cas de la segmentation en locuteur et de la traduction automatique et seront utilisés lors d'évaluations internationales organisées en 2020 : ALLIES/ALbayzin 2020 et WMT2020. À la fin de ces évaluations, les métriques, protocoles et données collectées pour ces évaluations seront distribuées gratuitement pour les fins de recherche.

## Remerciements

Ce travail a été financé par le projet Chist-ERA ALLIES (ARN-17-CHR2-0004-01)<sup>2</sup>

---

2. <https://lium.univ-lemans.fr/allies/>

## Références

- BELUCH W. H., GENEWEIN T., NÜRNBERGER A. & KÖHLER J. M. (2018). The power of ensembles for active learning in image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, p. 9368–9377.
- BISHOP C. M. (2006). *Pattern recognition and machine learning*. springer.
- BROUX P.-A., DOUKHAN D., PETITRENAUD S., MEIGNIER S. & CARRIVE J. (2018). Computer-assisted speaker diarization : How to evaluate human corrections. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*.
- CELEMIN C. & RUIZ-DEL SOLAR J. (2019). An interactive framework for learning continuous actions policies based on corrective feedback. *Journal of Intelligent & Robotic Systems*, **95**(1), 77–97.
- CHEN Z. & LIU B. (2016). Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, **10**(3), 1–145.
- DRUGMAN T., PYLKKONEN J. & KNESER R. (2019). Active and semi-supervised learning in asr : Benefits on the acoustic and language models. *arXiv preprint arXiv :1903.02852*.
- FRENCH R. M. (1999). Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, **3**, 128–135.
- GEOFFROIS E. (2016). Evaluating interactive system adaptation. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, p. 256–260.
- KEMKER R., MCCLURE M., ABITINO A., HAYES T. L. & KANAN C. (2018). Measuring catastrophic forgetting in neural networks. In *Thirty-second AAAI conference on artificial intelligence*.
- KROGH A. & VEDELSBY J. (1995). Neural network ensembles, cross validation, and active learning. In *Advances in neural information processing systems*, p. 231–238.
- LOMONACO V. & MALTONI D. (2017). Core50 : a new dataset and benchmark for continuous object recognition. In *Conference on Robot Learning*, p. 17–26.
- PAPINENI K., ROUKOS S., WARD T. & ZHU W.-J. (2002). BLEU : A method for automatic evaluation of machine translation. In *Proceedings of the Annual Meeting on Association for Computational Linguistics (ACL) : ACL*.
- PARISI G. I., KEMKER R., PART J. L., KANAN C. & WERMTER S. (2019). Continual lifelong learning with neural networks : A review. *Neural Networks*.
- PÉREZ-DATTARI R., CELEMIN C., RUIZ-DEL SOLAR J. & KOBER J. (2018). Interactive learning with corrective feedback for policies based on deep neural networks. *arXiv preprint arXiv :1810.00466*.
- PROKOPALO Y., MEIGNIER S., GALIBERT O., BARRAULT L. & LARCHER A. (2020). Evaluation of lifelong learning systems. In *International Conference on Language Resources and Evaluation (LREC)*.
- QUONERO-CANDELA J., SUGIYAMA M., SCHWAIGHOFER A. & LAWRENCE N. D. (2009). *Dataset Shift in Machine Learning*. The MIT Press.
- SIDDHANT A. & LIPTON Z. C. (2018). Deep bayesian active learning for natural language processing : Results of a large-scale empirical study. *arXiv preprint arXiv :1808.05697*.