



## Safe transfer learning for dialogue applications

Nicolas Carrara, Romain Laroche, Jean-Léon Bouraoui, Tanguy Urvoy, Olivier Pietquin

### ► To cite this version:

Nicolas Carrara, Romain Laroche, Jean-Léon Bouraoui, Tanguy Urvoy, Olivier Pietquin. Safe transfer learning for dialogue applications. SLSP 2018 - 6th International Conference on Statistical Language and Speech Processing, Oct 2018, Mons, Belgium. hal-01928102

**HAL Id: hal-01928102**

**<https://hal.science/hal-01928102>**

Submitted on 20 Nov 2018

**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.

# Safe transfer learning for dialogue applications.

Nicolas Carrara<sup>1,2</sup>, Romain Laroche<sup>3</sup>, Jean-Léon Bouraoui<sup>2</sup>, Tanguy Urvoy<sup>2</sup>,  
and Olivier Pietquin<sup>1</sup>

<sup>1</sup> Univ. Lille, CNRS, Centrale Lille, INRIA UMR 9189 - CRISTAL, Lille, France

`name.surname@inria.fr`

<sup>2</sup> Orange Labs, Lannion, France

`name.surname@orange.com`

<sup>3</sup> Microsoft Maluuba, Montral, Canada

`name.surname@microsoft.com`

**Abstract.** In this paper, we formulate the hypothesis that the first dialogues with a new user should be handle in a very conservative way, for two reasons : avoid user dropout; gather more successful dialogues to speedup the learning of the asymptotic strategy. To this extend, we propose to transfer a safe strategy to initiate the first dialogues.

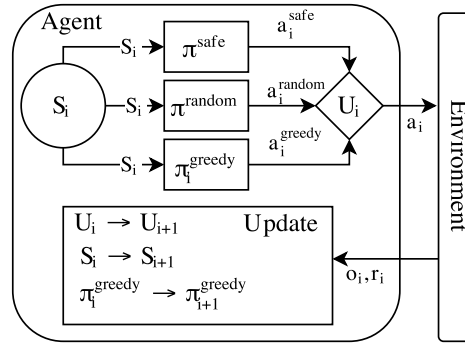
**Keywords:** Transfer Learning · Dialogue · Safety

## 1 Introduction

During its early steps of learning, a Reinforcement Learning (RL) [8] based dialogue agent does a lot of exploration that may lead to penalizing behavior. However, the first interactions between a user and a dialogue system are crucial to gain his trust. To improve jump-start performance of an RL agent, one can transfer a strategy [9, 7]. In dialogue, the strategy focuses on the success of the dialogue while minimizing its length [3, 5, 2, 4]. Rushing the dialogue may be problematic with some users and induce premature dialogue hangups. In one hand, it could lead to the lose of this user once for all. In an other hand, the lack of succeeding dialogues may affect the learning speed of the RL agent. To this extend, we introduce a novel algorithm :  $\epsilon$ -safe. It transfers a safe strategy which avoid any critical dialogue act to avoid the aforementioned problems.

## 2 $\epsilon$ -safe

$\epsilon$ -safe (Fig 1) is a  $Q$ -learning algorithm [6] where each action is decided by a randomly chosen policy among the greedy policy, an exploratory policy and the transferred safe policy. It may be an handcrafted policy, an RL policy (with large reward penalty on the catastrophic event) or even a safe-RL policy [1].

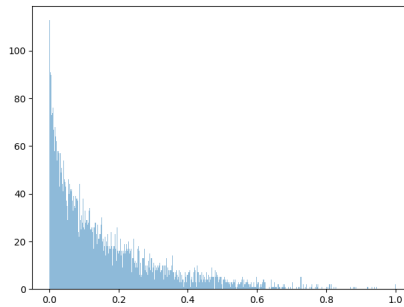
Fig. 1:  $\epsilon$ -safe algorithm.

### 3 Experiment

We test our algorithm on a simple slot-filling application. The agent ask for slot values, slot by slot in a fix order. Several acts are available :

- **ask\_next** : ask next slot (with NLU errors).
- **repeat\_oral** : repeat current slot (with NLU errors).
- **repeat\_numpad**: repeat using numeric pad (without NLU errors).
- **summarize\_inform**: summarize slots values and return the form result. If values are correct, the dialogue ends successfully, if not, the slot values are reset and the dialogue continues from the first slot.

**repeat\_numpad** is an unsafe action : the user hangups with probability  $p$ . For each new user,  $p$  is randomly generated. An histogram of  $p$  values is displayed Fig 2.

Fig. 2: Half-gaussian distribution of  $p$  values.

We define 2 handcrafted dialogue systems. The **safe** system uses **repeat\_oral** if the recognition score is below 0.5, otherwise **ask\_next**, or **summarize\_inform** after the last slot; The **unsafe** system uses **repeat\_numpad** instead. We compare two  $\epsilon$ -safe agents. **safe\_on** uses **safe** as transferred policy while **unsafe\_on** transfers the **unsafe** policy.

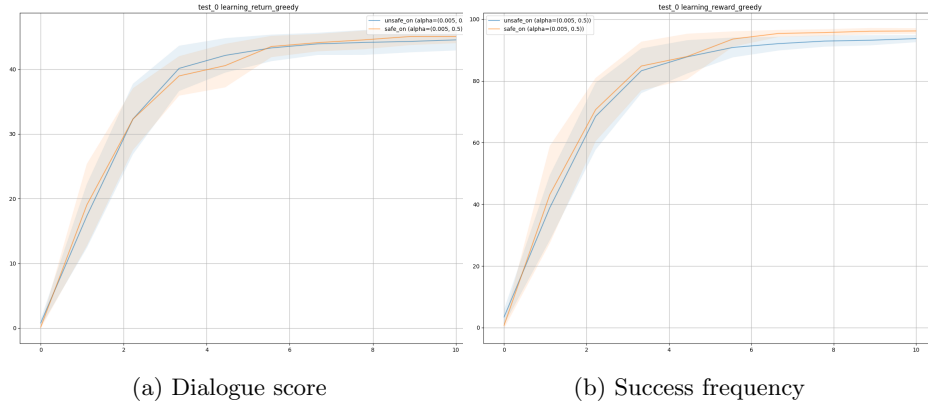


Fig. 3: Performance of the greedy policies.

We test **safe\_on** and **unsafe\_on** with 10 randomly generated users, for 1000 dialogues. We repeat the experiment 10 times. We display the performance of the greedy policies. The dialogue score (reward penalized by dialogue length) is plotted on Fig 3a while the dialogue success (reward only) is plotted on Fig 3b. We see that despite a slightly difference on the dialogue success, the dialogue score is the same. That means  $\epsilon$ -safe doesn't improve the learning speed of the greedy policy even if it is conservative enough to keep the user in the dialogue by avoiding catastrophic acts.

## 4 Future work

We believe that the hangup-model is too simple, so we plan to design the hangup-model with a Poisson distribution. We also want to replace the handcrafted policies by actual RL policies learned on source users. Finally, a real application on DSTC2 may be considered.

## References

1. Carrara, N., Laroche, R., Bouraoui, J., Urvoy, T., Pietquin, O.: A Fitted-Q Algorithm for Budgeted MDPs. Workshop on Safety, Risk and Uncertainty in Reinforcement Learning (UAI2018) (2018)
2. Carrara, N., Laroche, R., Pietquin, O.: Online learning and transfer for user adaptation in dialogue systems. In: SIGDIAL/SEMDIAL on negotiation dialog (2017)
3. Casanueva, N., Hain, T., Christensen, H., Marxer, R., Green, P.: Knowledge transfer between speakers for personalised dialogue management (2015)
4. Chandramohan, S., Geist, M., Pietquin, O.: Optimizing Spoken Dialogue Management from Data Corpora with Fitted Value Iteration. Proceedings of the International Conference on Speech Communication and Technologies (2010)
5. Genevay, A., Laroche, R.: Transfer learning for user adaptation in spoken dialogue systems. In: AAMAS (2016)
6. J. C. H. Watkins, C., Dayan, P.: Q-learning. In: Machine Learning. pp. 279–292 (1992)
7. Lazaric, A.: Transfer in Reinforcement Learning: a Framework and a Survey (2012)
8. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press Cambridge (1998)
9. Taylor, M., Stone, P.: Transfer Learning for Reinforcement Learning Domains : A Survey. JMLR (2009)