

Training a Gaming Agent on Brainwaves

Bartolomé Francisco, Moreno Juan, Navas Natalia, Vitali José,
Ramele Rodrigo, Member, IEEE, Santos Juan Miguel

Abstract—Error-related potential (ErrP) are a particular type of ERP elicited by a person attending a recognizable error. The purpose of this study is to determine if these signals can be used to train a Reinforcement Learning (RL) algorithm to learn an optimal policy. A game scenario is used to trigger the feedback response embedded in Electroencephalographic (EEG) signals of an observation human critic (OHC) that observes an agent playing a game. ErrP signals are captured using a Brain-Computer Interface (BCI) system. The experimental process consists of an individual observing a simple game scenario while their brain signals are captured. The game consists of a grid, where the object has to reach a desired target in the fewest amount of steps. Initially the object moves randomly within the grid, and a RL algorithm is trained for it to learn the optimal policy, in order to know how to reach the target in the least amount of steps. Results show that there is an effective transfer of information and that the agent successfully learns to solve the game efficiently. Initially when the game is simulated with random movements it takes the agent an average of 97 steps to reach the objective, whereas when the trained Q-Table is used to reach the objective, the agent can reach the goal in the optimal number of steps, which is 8 steps. The algorithm can be trained with information from different OHCs, since the resulting dataset is composed of rewards for each corresponding step, independent of who they were generated from. This shows that the error classification accuracy (approximately 0.67) is good enough for the algorithm to learn. The reward function only penalizes wrong steps, which means that type II error (not properly identifying a wrong movement) does not affect the accuracy, they only make the learning process slower. This study shows that: (i) the structure of a simple grid-based game that can elicit the ErrP signal component; (ii) the verification that low classification accuracy of just above chance level that produces noisy rewards is enough to allow an agent to learn the optimal policy; (iii) collaborative rewards from multiple observational human critics can compensate the lack of accuracy or the limited scope of transfer learning schemes.

Index Terms—ErrP, BCI, EEG, RL, Agent, AI

I. Introduction

THE effectiveness of today's human-machine interaction and artificial intelligence is limited by a communication bottleneck, as humans are required to translate high-level concepts into a machine-mandated sequence of instructions [?], [?]. Hence, new interaction methods are required to increase the communication bandwidth between computers and humans or to produce alternative communications systems to increase the efficiency of this channel. In this respect, video games have been widely used as test tools to assess new means of interactions [?], [?]. Video gaming agents are computer programs that can

sense the computer game environment, process information, and react accordingly within the environment. They are used in the context of testing and evaluating artificial intelligence algorithms that aim to win the game or to behave like a real user player [?]. In this work, the feedback obtained from an observational human critic (OHC) in the form of electroencephalographic (EEG) signals is used to evaluate the operational performance of a gaming agent. Observational human critics are silent subjects observing a computer gaming agent playing the game.

The feasibility of a distinct non-biological communication channel between the Central Nervous System (CNS) and a computer device has been previously proven with Brain Computer Interfaces (BCI) or Brain Machine Interfaces (BMI). [?]. BCI systems provide a new input modality that can be used in the context of a computer game [?], [?]. This advancement is relevant in the context of the accessibility for video games [?] and the growing area of e-sports [?].

In this study, gaming agents are trained using only signal components called Error-related Potentials (ErrP) that can be identified in the observer's brain signals. These types of signals can be found on EEG traces and are elicited when subjects are aware of the presence of an unexpected outcome, which they identify as an error. The analysis of ErrP signals is currently an extensive area of research in the neuroscience community [?]. Error-related Potentials can be detected by signal processing and machine learning techniques [?] and are also used in Brain-Computer Interfaces to implement or enhance artificial communication channels [?].

Given the scenario, Reinforcement Learning (RL) [?] stands out as a natural method to train the agent. Reinforcement Learning refers to an algorithmic learning strategy inspired on how biological agents learn by exploring their environment while getting negative or positive feedback rewards. The method aims to maximize positive rewards while minimizing negative feedback. Thus, the learning problem is posed as a stochastic optimization strategy [?]. Recently, this technique has been used extensively in the context of advances in artificial intelligence [?]. The influence of DeepBrain's AlphaGo project cannot be neglected, since it was the first to reach a very high proficiency when it won the complex game Go against several world champions [?].

Previous research has explored the usage of RL with reward signals based on brain activity, recorded by an EEG-based BCI system during task execution. The papers [?], [?], [?] have successfully demonstrated that a robot can be controlled with brain signals from a person who

R. Ramele and J.M.Santos are with the Department of Computer Engineering, Instituto Tecnológico de Buenos Aires(ITBA), Argentina, e-mail: rramele@itba.edu.ar.

Manuscript received December 9, 2019; revised August 25, 2020.

is observing a robot solve a task. Moreover, a growing number of studies have demonstrated the feasibility of using ErrPs as rewards for RL schemes such as to enhance robotic behaviour [?], to assess air traffic controller's decisions [?] or to categorize actions as errors [?]. Other approaches have used these signals as important feedback for human-robot interaction or to implement shared-control strategies [?]. Additionally, ErrPs have also been used in the context of games as an additional feedback channel that can be explored to improve gaming experience [?], [?].

Therefore, we aim to use the information extracted from brainwaves to enhance the performance of a gaming agent. The three contributions are (1) a simple game that can elicit the ErrP potential, (2) results that confirm that even when ErrP classification accuracy is low and produces a noisy reward signal, enough information is generated for an agent to learn the optimal policy and solve a simple game and (3) collaborative rewards from multiple observational human observers can compensate the lack of classification accuracy or the inefficacy of transfer learning procedures for brainwaves signals.

In Section II the general layout of the cognitive game is described. Sections II-A and II-B outline the cognitive game procedure used to obtain rewards in the form of ErrP components. Section II-F describes the gaming agent learning procedure. Lastly, results and conclusions are exposed in Sections III and IV.

II. Materials and Methods

The experimental procedure is summarized in Figure 1. The proposed system has two distinct parts. This first part consists of the collection of brainwave signals from a person that is watching an agent play a game. The agent knows the game rules but not how to win it. The second part, the gaming agent learning phase, is where the agent can learn the winning strategy using the person's feedback to improve its own performance.

A. Brainwave Session

The retrieval of the OHC's brain activity, called the brainwave session, is one of the most critical parts of the study. Subjects are recruited voluntarily and given a form with questions regarding their health (previous health issues and particular visual sensitivity), habits (sleeping hours, caffeine and alcohol consumption), and a consent petition to collect the required data. The brainwave sessions are performed with 8 subjects, 5 males and 3 females, with an average age of 25.12 years, a standard deviation of 1.54 years, and a range of 22-28 years. All subjects have normal vision, are right-handed and no history of neurological disorders.

After the form is filled out, a short description of the procedure is given to each subject. They are only told that the objective of the agent is to reach the goal and the four movements that the agent can make. When this concludes, the subject is introduced to the

wireless digital EEG device (g.Nautilus, g.Tec, Austria) that she/he has to wear during the brainwave session. It has eight electrodes (g.LADYbird, g.Tec, Austria) on the positions Fz, Cz, Pz, Oz, P3, P4, PO7, and PO8, identified according to the 10-20 International System, with a reference set to the right ear lobe and ground set as the AFz position. The electrode contact points are adjusted applying conductive gel until the impedance values displayed by the program g.NeedAccess (g.Tec, Austria) are within the desired range. This process takes between 10 to 15 minutes. After this step, the subject is instructed to close their eyes, make eye movements and muscle chew in order to check the program and guarantee that the live channel values are accurate.

Once the headset is correctly applied, the OpenVibe Acquisition Server program, from the OpenVibe platform [?], is launched and configured with a sampling rate of 250 Hz. A 50 Hz notch filter is applied to filter out power line noise. An additional bandpass filter between 0.5 Hz and 60 Hz is applied. Data are handled and processed with the OpenVibe Designer, from the same platform, using 8 channels for the brain data (one channel per electrode) and an additional channel to record the stimulus, which corresponds to a game movement performed by the agent. After everything is connected, the subject is seated in a comfortable chair in front of a computer screen. The brightness of the screen is set to the maximum setting to avoid any visual inconvenience in which the subject can not distinguish the components of the game that appear on the screen.

The Acquisition Server receives and synchronizes the signal data from the headset and any event information from the game, and transfers it to the OpenVibe Designer application. When the subject is ready, the Game Manager and the OpenVibe Designer programs are launched and configured to communicate with the previously mentioned Acquisition Server. A brainwave session consists of several matches, each one being a gameplay. In the end, the sequence of game movements and the signal data generated for each match are saved for offline processing ¹.

B. Cognitive Game Procedure

The game parsimoniously consists of a 5x5 grid of grey circular spots with a black background. A blue spot indicates the current position of the agent and a green spot represents the goal, as shown in Figure 2. The agent's objective is to reach the goal. The circular spot representing the goal remains static at the bottom-right position of the grid, while the one representing the position of the agent starts at the upper-left position of the grid and moves in each iteration. When the agent reaches the goal, the position where the agent and the goal are located turns red, showing that the match has ended. There are four possible movements that the agent can perform: it can go upwards, downwards, towards the

¹The brainwave dataset has been published on the IEEE DataPort initiative [?].

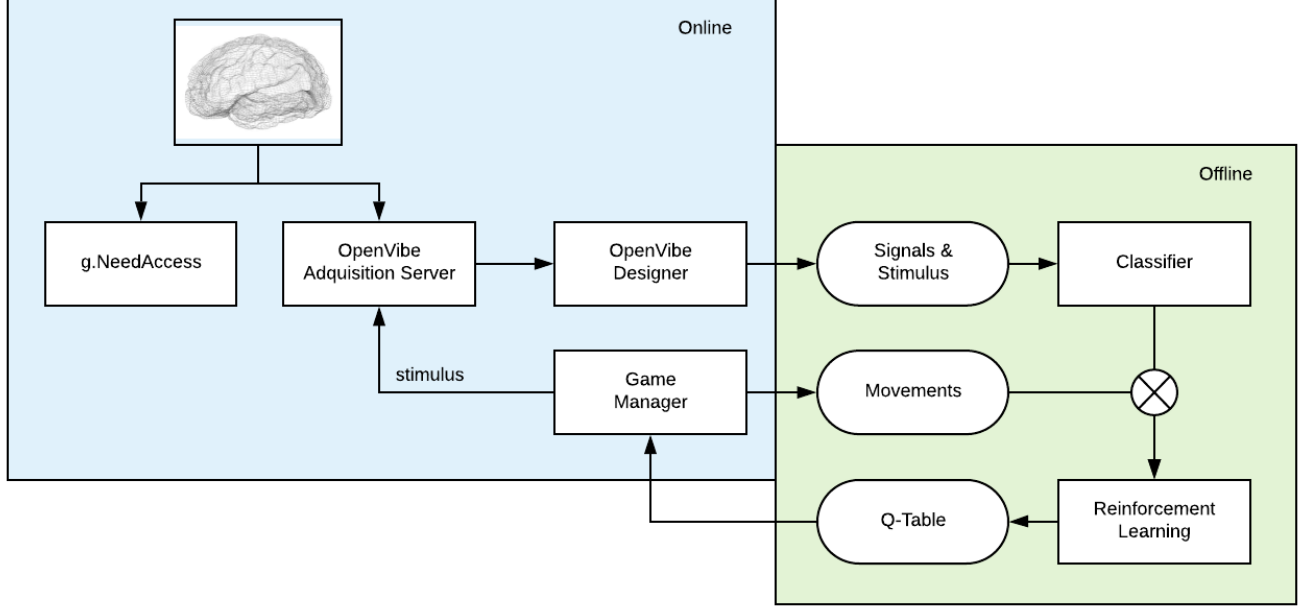


Fig. 1: Overview of the experimental procedure. Brainwaves are obtained by the OpenVibe Acquisition Server. The Game Manager is responsible for generating the game screen, the game mechanics, and the game movements performed by the gaming agent. It is also connected to the Acquisition Server to send stimulus information. The captured information is stored by the OpenVibe Designer. Offline, EEG signals are classified and they are linked to each game movement calculated by the Game Manager to determine proper rewards for each action. This information is used by a Reinforcement Learning algorithm that iteratively trains a Q-Table in order to improve the performance of the agent that plays the game.

left and the right, and those movements are bounded to avoid the agent from leaving the grid. The movement direction is selected randomly and is executed once every 2 seconds. After each gameplay, there is a pause of 5 seconds until the next match starts. Each time an agent moves, the Game Manager program sends an event marker to the Acquisition Server. This is considered a stimulus to the observational human critic. The game is designed as to be evident whenever there is an error (i.e. the agent moves away from the objective) so the subject can notice it immediately after the stimulus is presented, possibly triggering the expected cognitive response, which can be imprinted as an ErrP component within the EEG stream.

C. Signal Processing, Segmentation and Classification

To aid the detection of the ErrP response, an offline processing pipeline and classifier is constructed to identify whether the action taken by the agent is an error or not, from the human observer's point of view. It is developed in Python using the "MNE" software platform [?], which is a package designed specifically for processing EEG and Magnetoencephalography data, and built upon the machine learning library Scikit-Learn [?].

This pipeline consists of the offline processing of the collected signals used to train a classifier that can decide whether an error potential is triggered. Firstly, the output

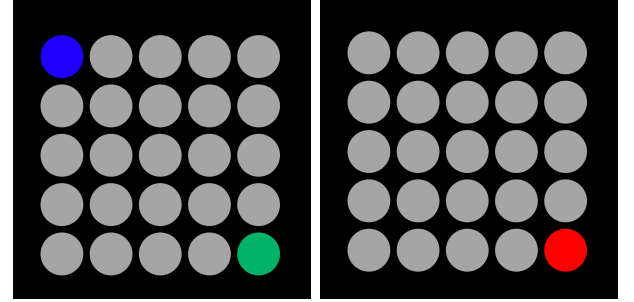


Fig. 2: Grid system representation used in the cognitive game. The blue spot represents the initial location while the green spot represents the target location. Once the agent reaches the target spot, its color turns red to indicate the end of the play.

of a brainwave session is read and an additional band-pass filter of 0.1-20.0 Hz is applied to the signal. Samples that correspond to the start of an event are tagged using the data from the stimulus channel.

After the raw data are loaded and tagged, epochs are extracted. Epochs consist of all the sample points that take place during the 2 seconds from the start of the event, 2 seconds corresponding to the time it takes for each action to take place, resulting in 500 samples per channel, as the sampling frequency is 250 Hz. Thus, each epoch is

composed of a matrix 500 x 8 channels.

Samples that do not correspond to an epoch (located beyond the 2 seconds frame after the onset of the event) are not used. Also, epochs referring to the start or finish of each match are excluded.

In this way, the raw data of a brainwave session is processed into an array of matches where each element is an array of epochs tagged with a number specifying the prediction of the classifier, i.e. if the epoch corresponds to an action that made the agent move further from the goal (hit) or an action that made the agent moves closer to the goal (no-hit). The ErrP is expected to be found in hits. To get the data ready for classification, the stimulus channel is removed to classify the signals using only the EEG data. Each epoch is regularized using a MinMaxScaler, i.e. subtracting the minimum value in the epoch and dividing by the signal peak-to-peak amplitude [?]. The eight channels are concatenated using the MNE Vectorizer function, which transforms the data matrix into a single array sample. Lastly, this data are used by the classification module as information to train and test a classifier. Five different classification algorithms are used: Logistic Regression, Multilayer Perceptron with a hidden layer of 100 neurons (i.e. default values for the Scikit-Learn MLPClassifier), Random Forest, KNeighbours with k=3 and finally a linear kernel Support Vector Classifier (i.e. SVM) [?].

D. Reinforcement Learning

Each match consists of a list of game movement configurations and the associated epochs obtained from OHC's brainwaves. The set of matches of each OHC is split into training and testing. Training matches are used to train the classifier to identify the ErrP signal. After a classifier is trained, the epochs extracted from the test matches are classified as hit or no-hit. A reward for each movement in the match is generated based on the prediction from the classifier for that movement. The reward can either be -1 when the event is classified as a hit or 0 when it is classified as a no-hit. The accuracy of these rewards depends on the performance of the classifier. The list of game movements and their associated reward information is used to train the agent by a variant of Reinforcement Learning called Q-Learning algorithm.

E. Q-Learning

Q-Learning [?] is a form of model-free reinforcement learning, where an agent tries an action at a particular state and evaluates its consequences in terms of the reward or penalty it receives. To represent rewards, a matrix $Q(s, a)$ is used, where rows correspond to all the possible states, and columns represent all possible actions. This matrix is known as the Q-Table. The algorithm proceeds by randomly choosing what action to do and iteratively

updating the Q-Table based on the received reward r by the following equation

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{\tilde{a}} Q(\tilde{s}, \tilde{a}) - Q(s, a)] \quad (1)$$

where s is the current state, a the action, α the learning rate and γ the discount factor, a value between 0 and 1 that determines the importance of long term results versus immediate rewards. Hence, $Q(s, a)$ is the expected value of the sum of discounted rewards that the agent will receive if in the s state, it takes the action a according to this policy. Once the environment has been extensively explored and the Q-Table has been optimized, the action chosen for a given state is the one that maximizes the expected reward according to the Q-Table matrix.

The algorithm is developed in Python and uses the OpenAI Gym toolkit [?]. Gym is a toolkit for developing and comparing reinforcement learning algorithms. It makes no assumptions about the structure of an agent, and is compatible with any numerical computation library, such as TensorFlow or Theano [?].

F. Gaming Agent Learning Procedure

The gaming agent learning procedure uses the testing matches from brainwave sessions produced during the cognitive game procedure phase, and their components are schematized in Figure 1.

This phase is divided into a sequence of run sessions and a gaming agent training match. During the run session, the agent plays 200 matches guided by a specific Q-Table with a 5% chance of randomly selecting a movement, to reduce deadlocks and loops. Following the run session, the agent performs a single gaming agent training match. The gaming agent starts first with a Q-Table initialized with zeros, so the initial policy for the agent is randomized. For the agent to learn from the feedback generated by the OHC, movement actions are determined by the reply of the agent's actions that were taken during one brainwave session match, in an offline reinforcement learning scheme [?]. This allows the Q-Table to be built based on the OHC's feedback from the movements the agent took, which were executed pseudo-randomly during the brainwave session. The previously mentioned feedback is not explicit as it comes from the interpreted brain signal data. This implies that the reward is determined by the OHC's brain activity.

Hence, following the iterative procedure based on Equation 1, the Q-Table is updated in each gaming agent training match. After the algorithm finishes replicating all the steps from the brainwave session match, the Q-Table is stored and used by the agent in the next run session.

III. Results

Figure 3 shows the binary classification accuracy obtained for the eight OHCs using five different classification algorithms and using a 10-fold cross-validation procedure.

The best overall performance is obtained using Logistic Regression. In addition, the classification accuracy obtained averaging 5 epochs is shown as well. Although the signal averaging procedure improves the Signal-To-Noise-Ratio (SNR) of the ErrP response, it reduces the number of data samples, producing a clear improvement only for the OHCs that contain more samples (1, 3 and 7).

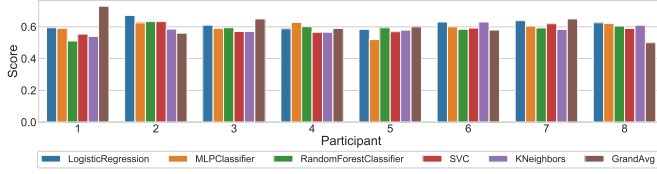


Fig. 3: Binary single trial classification score using five different classifiers while recognizing ErrP potentials for the eight OHCs. The classification score using an ensemble average of 5 epochs is shown as well (GrandAvg). Chance Level is 0.5.

On the other hand, Figure 4 shows the average amount of steps it takes for the agent to reach the goal for each OHC, as the Q-Table is progressively trained using the reward information obtained from the prediction of the trained classifier. Each point corresponds to a run session where the average number of steps for the agent to reach the goal using a specific Q-Table is specified, for 200 repetitions. The first point, at the x-value 0, represents the number of steps the agent takes to reach the goal with an untrained Q-Table, where movements are decided randomly. The next point corresponds to the amount of steps it takes to reach the goal using a policy derived from a Q-Table trained after one brainwave session match, and so on.

The results show that as the Q-Table is progressively trained the average amount of steps decreases, meaning that the agent learns. However, the rate at which it learns varies per OHC, depending on the classification accuracy of the extracted brainwaves. For example results for OHC 1 show faster learning than those of OHC 8 (Figure 4).

In the case for OHC 5 and 6, the reward information obtained from the brainwaves is not enough to train the agent effectively. Figures 4 for OHC 5 and 6 show no apparent learning, as the amount of steps to reach the goal doesn't decrease when trained. These results are also consistent with their classification ROC curves, shown in Figures 5 obtained for both OHCs, where the area under the curve are close to chance level. Both OHCs have less recorded data from the sessions in comparison to the rest of the OHCs. This variation in performance for different OHCs has been studied extensively in BCI [?]. Besides low data samples, there are other reasons affecting the classification accuracy: cognitive reasons (i.e. the OHC not paying extensively attention to the game dynamics), very low SNR of the ErrP component or even the BCI-illiteracy phenomena where the specific OHC's signals do not contain the expected component response [?].

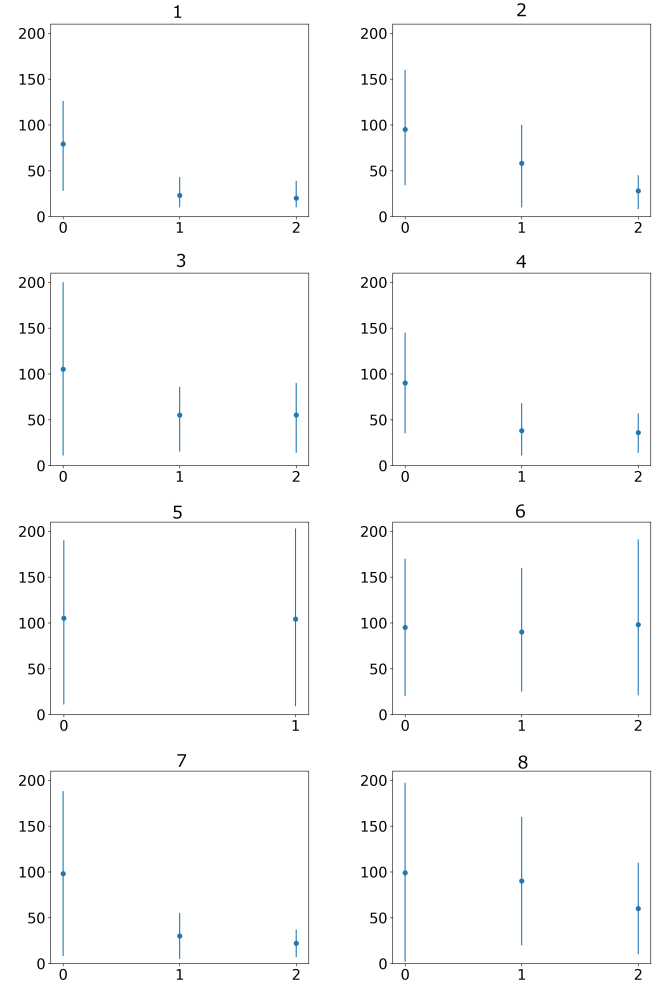


Fig. 4: Average number of steps for the agent to reach the goal when trained with rewards generated from brainwaves from OHCs 1-8. Y axis show the averaged number of steps for a run session, while x axis show the number of game matches used to cumulative train the Q-Table.

Figure 6 shows the result of an agent successively trained with brainwave session matches where the EEG is generated with random signals. In this case, random EEG signals are generated using OpenVibe Acquisition Server signal generator for all channels, as if they were produced by an OHC who doesn't pay attention to the game. The agent learns nothing, and regardless of the number of matches that are used to learn the Q-Table, the number of steps required to reach the goal does not decrease. This pattern is also obtained when the game matches from OHCs 5 and 6 are used, showing that the reward labeling predicted by the trained classifier for those cases worked like a random classifier.

Electroencephalographic signals have high inter-subject variability [?]. This is evidenced in Figure 7 where the agent training is performed with rewards obtained by classifying epochs from one Tester OHC and a classifier that was trained using the brainwaves from a different Trainer OHC. The figure shows the cumulative variation

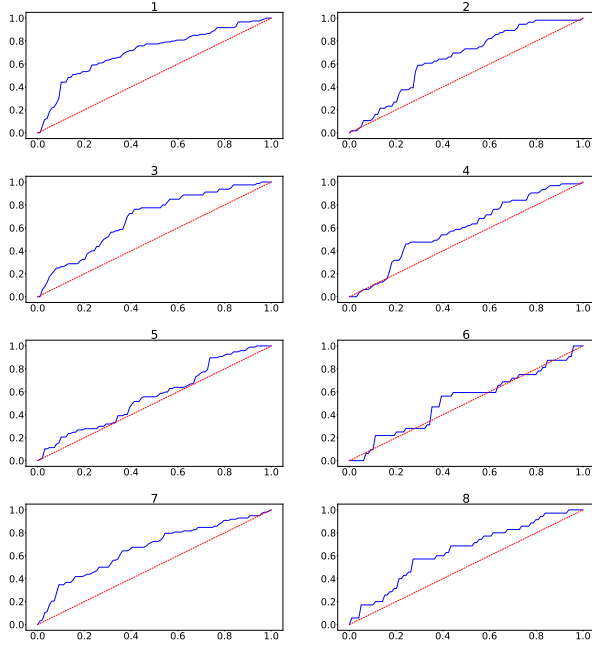


Fig. 5: ROC Curves for OHCs 1-8. True positive rate is on the vertical axis and false positive rate on the horizontal axis.

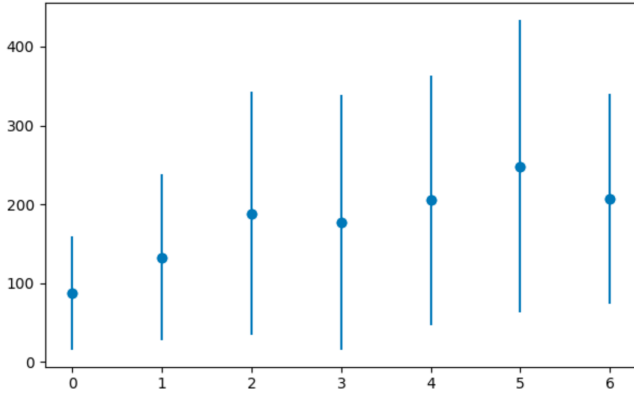


Fig. 6: Average number of steps for the agent to reach the goal when trained with a classifier produced from sham EEG signals. X axis show the number of gaming agent training matches used to train the Q-Table.

for all run sessions on the average number of steps required to reach the goal after training the agent with all the available matches from the brainwave session. Enhancements are shown as negative values. Only the diagonal of the heatmap matrix shows a clear improvement in terms of the reduction of the required number of steps to reach the goal (averaged per 200 runs) which corresponds to the same information for each OHC shown in Figure 4. For the transfer learning experiment [?], no performance gain is evidenced, the agents learn nothing and this implies that the reward information is useless.

Finally, Figure 8 shows the result of training an agent with cumulative brainwave session matches from OHCs 1, 2, 3, 4, 7, and 8. It can be seen that the overall performance

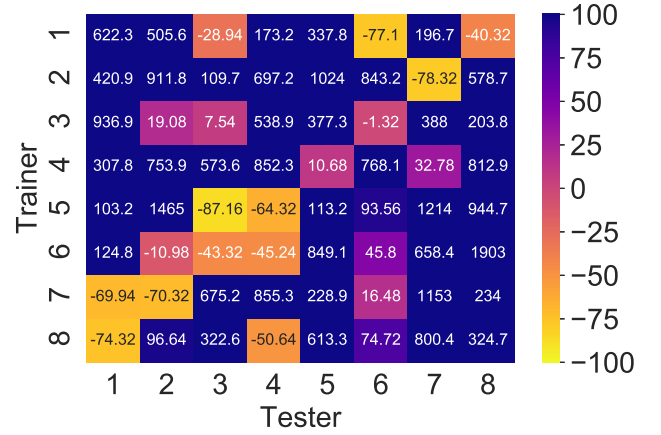


Fig. 7: Heatmap for the transfer learning experiment. Values represent the reduction in the average number of steps required to reach the goal. Negative values represent net improvements.

of the agent improves as long as there is information to produce rewards, regardless of the fact that they were generated from classifiers trained with different OHC's signals.

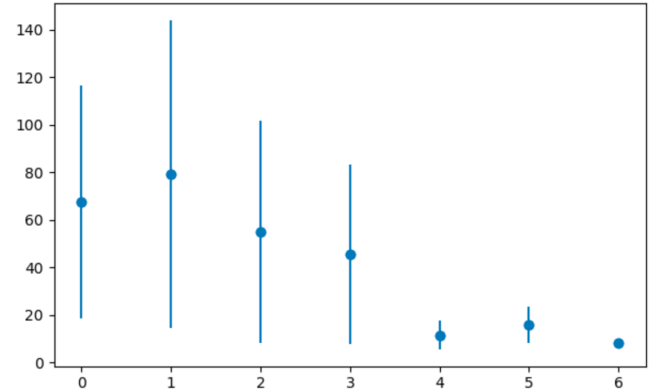


Fig. 8: Average steps using Q-Table trained with brainwave session matches from six different OHCs. X axis show the progressive number of gaming agent training matches used to train the Q-Table. These matches correspond to all subjects, excluding subjects 5 and 6 which individually show no significant learning progress.

IV. Conclusion

This work aims to propose a simple game that can use the ErrP component to train a gaming agent using a RL model. The collected data shows that ErrP signals can in fact be classified and used to train an agent effectively.

This proposal tries to keep the system as simple as possible, emphasizing information flow from the subjective error perception of the human critic. Rewards are generated using the signal processing and classification pipeline, and the Q-Table updates, enhancing the performance of the gaming agent.

One additional aspect to remark is the robustness of the learning strategy based on Q-Learning [?], [?]. The obtained accuracy to discriminate ErrPs is low. However, even with such low accuracy values, the RL algorithm was able to extract meaningful information from rewards that were helpful to improve, and often maximize, the agent's performance. Additionally, one important aspect of the classification results is the low percentage of false positives (Figure 9), showing a high specificity. On the other hand, the percentage of false negatives is generally higher. However, even though this implies that the agent misses frequently when a wrong action takes place, this is not hindering the overall performance and the agent is still learning. Though scarce, accurate rewards are very useful for the RL algorithm.

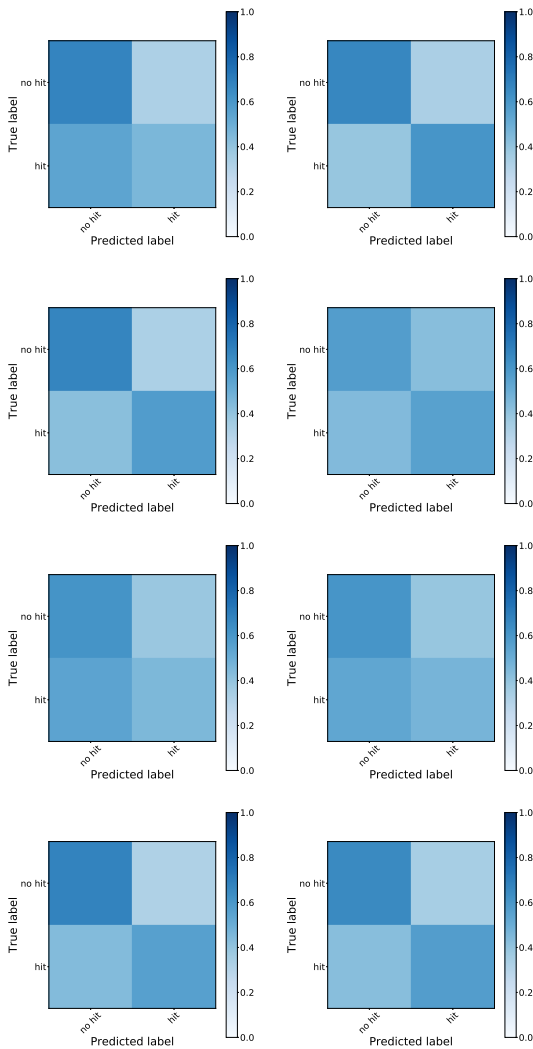


Fig. 9: Confusion Matrix for OHCs 1-8. Darker colors show higher values. It can be seen the lower percentage of false positives (upper right corner of each chart).

At the same time, effective agent training depends on the OHC's training data. Results confirm the futility or complexity of using Transfer Learning [?]: training a classifier with data obtained from one OHC, and using

the same classifier to identify ErrPs for another OHC does not increase the performance of the agent. Despite that, the rewards generated from different subject's classifiers can be used to train the same Q-Table to improve its performance, which may lead to strategies where the overall performance is enhanced based on the information from different human critics at the same time. There seems to be an agreement in terms of the subjective interpretation of what may be an appropriate movement to reach the goal.

The simple setup of the grid-based game allows further experimentation, using the reduction on the number of average steps to reach the goal as a validation of the achieved information transfer. It will be of research interest to verify if the smooth progression towards the end alters the shape of the ErrP response, how the ErrP response is triggered in relation with different shapes and colors of the board markers [?], or if there is a differential ErrP signal component in relation to up, down, left and right movements. In addition, the outcome of manipulating the stimulus could be further studied as well as the influence on the results if incentives are given to participants.

Further work will be conducted in order to increase the complexity of the game to allow the possibility that the target position is dynamically changed. Although we found that the best performing classifier is Logistic Regression, there is room for improvement. The classifier could be enhanced to recognize the Error Potential [?] more effectively or could be pre-trained to allow higher accuracy [?].

Acknowledgment

The authors would like to thank the Laboratory Centro de Inteligencia Computacional and to ITBA University.

Funding

This work was supported by the grant ITBACyT-15 issued by ITBA University.