

# Reward-based online learning in non-stationary environments: Adapting a P300-speller with a “Backspace” key

Emmanuel Dacé <sup>1 2</sup>    Timothée Proix <sup>2</sup>    Liva Ralaivola <sup>3</sup>

<sup>1</sup>Ecole Centrale Marseille

<sup>2</sup>INS - Aix-Marseille Université - Inserm UMR 1106

<sup>3</sup>LIF - Aix-Marseille Université - CNRS UMR 7229

July 13, 2015

# Table of contents

- 1 Introduction
- 2 Reward-based learning
- 3 Simulations
- 4 Conclusion



# Plan

- 1 Introduction
- 2 Reward-based learning
- 3 Simulations
- 4 Conclusion



# Brain-Computer Interfaces

- Embedded classifiers:
  - real-time
  - noisy (EEG)
  - subject/use case specific
  - non-stationary

⇒ Adaptive Learning

- Brain Computer Interfaces, a tool for :
  - Communication (in the absence of a motor capabilities)
  - Brain monitoring / neurofeedback
  - Motor rehabilitation
- “CO-ADAPT” project : INRIA Sophia/ INSERM Lyon/ CNRS LATP, ... (French ANR funding)
  - “co-adaptive” motor imagery
  - “co-adaptive” P300 speller



# Embedded classifiers

- Classification problem : sources  $\rightarrow$  signal  $\rightarrow$  features extraction  $\rightarrow$  classification
  - adaptive feature extraction
  - adaptive classification
- Online learning : “light” classifier update at each processing step
  - supervised online learning: stochastic gradient descent
  - unsupervised online learning: online mobile centers (K-means, EM,...)
  - reward-based online learning: stochastic classifiers + policy/value iteration
    - exploration/exploitation trade-off
    - which reward?



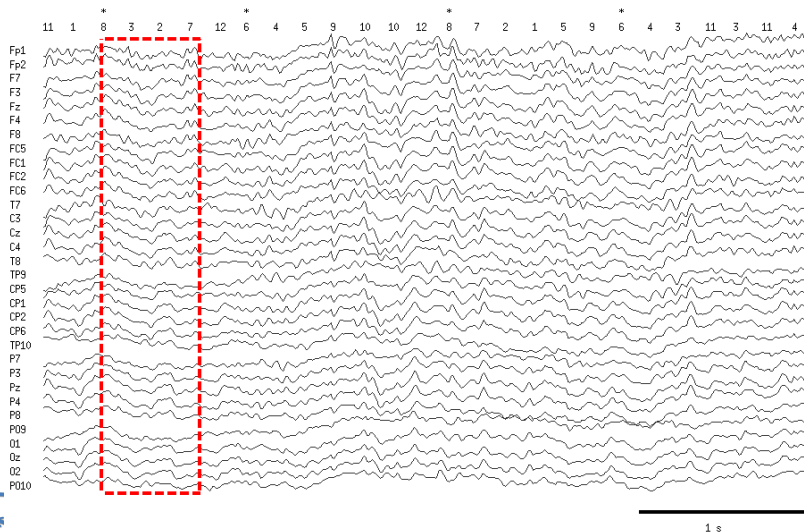
# P300 speller

- EEG :
  - 10 - 60 channels (surface electric potential - H Berger, 1929)
  - high temporal resolution / low spatial resolution
  - noisy, non-reliable,... “Evoqued potentials” technique
  - the “P300” ERP is “surprise” effect (“oddball” paradigm)
- P300-speller (Farwell and Donchin, 1988):
  - based on the “oddball” paradigm
  - 6 x 6 letters grid
  - random row/column magnification (every 150-300 ms)
  - row/column evidence build-up + argmax choice
  - low SNR / low bit rate (many flashes for one letter)
  - spelling accuracy tends to decrease in the long run

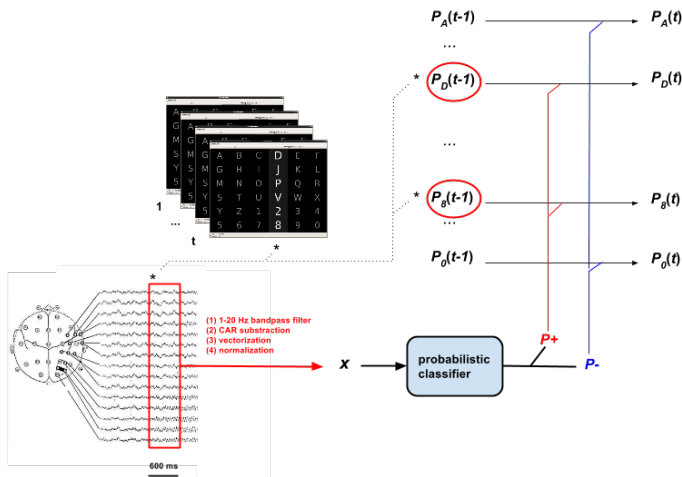
A	B	C	D	E	F
G	H	I	J	K	L
M	N	O	P	Q	R
S	T	U	V	W	X
Y	Z	1	2	3	4
5	6	7	8	9	0



# EEG data (from Inserm U1028, Lyon, France)



# Data processing pipeline





# P300-speller roadmap

- Transfer learning [Kindermans et al., 2012b, Congedo et al., 2013]:
  - Across subjects
  - generic classifier “pre-learning”
  - smart initialization
- Optimal display
- Evidence build-up [Perrin, 2012, Kindermans and Schrauwen, 2013]:
  - Probabilistic classifier : posterior estimate
  - Evidence accumulation
  - Threshold-based dynamic stopping
  - Speed-accuracy trade-off
- Online learning :
  - Toward a subject-specific classifier
  - Recover from unexpected changes
  - Different approaches :
    - EM [Li and Guan, 2006, Kindermans et al., 2012a]
    - RL [Daucé et al., 2013]



# Plan

- 1 Introduction
- 2 Reward-based learning
- 3 Simulations
- 4 Conclusion



# Reinforcement learning : multiple approaches

- general RL problem :
  - observations :  $x_1, \dots, x_T$
  - policy (random decision = exploration) :  $\pi(x, y) = P(y|x)$
  - responses :  $y_1, \dots, y_T$
  - reward :  $r_T \in \mathbb{R}$
- Task : try to improve the politics for increasing time/trial
  - Actor-critic / Q-learning :
    - temporal credit assignment (delayed reward)
    - discrete state/action space (LUT)
    - Value update vs. action
  - Multi-armed bandit :
    - instant reward
    - contextual bandit
  - Policy-gradient :
    - instant or delayed rewards
    - discrete or continuous state/actions space
    - direct update of the policy



# Rewards in classification

- Optimal  $\mathbf{w}$  unknown  $\rightarrow$  model-free, trial and error
- Online stochastic classifier :
  - read input observations set :  $\underline{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_K)$
  - give a score to every class :  $\forall k, \pi(\underline{\mathbf{x}}, k; \mathbf{w}) = \frac{\exp(\mathbf{x}_k \mathbf{w}^T)}{\sum_l \exp(\mathbf{x}_l \mathbf{w}^T)}$
  - choose the response at random (Softmax choice)
  - read the *reward*  $r$
  - update  $\mathbf{w}$
- Which reward ?
  - “error” potential after the classifier’s response:



- “BACKSPACE” key on the virtual keyboard

# Learning and forgetting : Regularized Policy gradient update

- Regularized optimization ( $\lambda$  hyperparameter):

$$\max_{\mathbf{w}} \mathcal{H} = \max_{\mathbf{w}} E(r) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Regularized gradient ascent :  $\nabla_{\mathbf{w}} \mathcal{H} = E(r \nabla_{\mathbf{w}} \ln \pi(\underline{\mathbf{x}}, y; \mathbf{w})) - \lambda \mathbf{w}$ 
  - Gradient estimator (stochastic gradient) :

$$\langle r_t \nabla_{\mathbf{w}} \ln \pi(\underline{\mathbf{x}}_t, y_t; \mathbf{w}) \rangle_{1..T}$$

- Online update (learning rate  $\eta \ll 1$ ):

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \eta (r \nabla_{\mathbf{w}} \ln \pi(\underline{\mathbf{x}}, y; \mathbf{w}) - \lambda \mathbf{w}) \\ &= (1 - \eta \lambda) \mathbf{w} + \eta r \nabla_{\mathbf{w}} \ln \pi(\underline{\mathbf{x}}, y; \mathbf{w}) \end{aligned}$$

- The old examples “fade away” as time passes  $\rightarrow$  tracking algorithm and novelty detection (Kivinen et al, 2010)

# The “oddball” update case

- Policy gradient:

$$\begin{aligned} \mathbf{g}(\underline{\mathbf{x}}, y) &= r \nabla_{\mathbf{w}} \ln \pi(\underline{\mathbf{x}}, y; \mathbf{w}) \\ &= r \left( \mathbf{x}_y - \sum_k \pi(\underline{\mathbf{x}}, k; \mathbf{w}) \mathbf{x}_k \right) \end{aligned}$$

- Update :

$$\mathbf{w}_t = (1 - \eta\lambda) \mathbf{w}_{t-1} + \eta r_t \left( \mathbf{x}_{y_t,t} - \sum_{k=1}^K \pi(\underline{\mathbf{x}}, k; \mathbf{w}_{t-1}) \mathbf{x}_{k,t} \right)$$



# Special cases

- Binary rewards :  $r^+, r^-$ ; let  $y^*$  be the “real” response :

$$E_{Y|X}(\mathbf{g}(\underline{\mathbf{x}}, y)) = (r^+ - r^-) \pi(\underline{\mathbf{x}}, y^*; \mathbf{w}) \left( \mathbf{x}_{y^*} - \sum_k \pi(\underline{\mathbf{x}}, k; \mathbf{w}) \mathbf{x}_k \right)$$



# Special cases

- Binary rewards :  $r^+, r^-$ ; let  $y^*$  be the “real” response :

$$E_{Y|X}(g(\underline{x}, y)) = \underbrace{(r^+ - r^-)\pi(\underline{x}, y^*; \mathbf{w})}_{=1?} \underbrace{\left( \mathbf{x}_{y^*} - \sum_k \pi(\underline{x}, k; \mathbf{w}) \mathbf{x}_k \right)}_{\text{Logistic gradient}}$$

•



# Special cases

- Binary rewards :  $r^+, r^-$ ; let  $y^*$  be the “real” response :

$$E_{Y|X}(\mathbf{g}(\underline{\mathbf{x}}, y)) = \underbrace{(r^+ - r^-)\pi(\underline{\mathbf{x}}, y^*; \mathbf{w})}_{=1?} \underbrace{\left( \mathbf{x}_{y^*} - \sum_k \pi(\underline{\mathbf{x}}, k; \mathbf{w}) \mathbf{x}_k \right)}_{\text{Logistic gradient}}$$

- Noisy rewards : let  $p_{\text{valid}}$  be the rate of valid rewards :

$$E_{Y|X}(\mathbf{g}(\underline{\mathbf{x}}, y)) = (2p_{\text{valid}} - 1)(r^+ - r^-)\pi(\underline{\mathbf{x}}, y^*; \mathbf{w}) \left( \mathbf{x}_{y^*} - \sum_k \pi(\underline{\mathbf{x}}, k; \mathbf{w}) \mathbf{x}_k \right)$$

# Plan

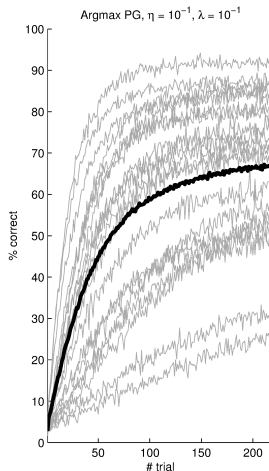
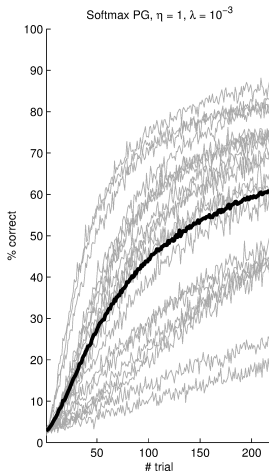
- 1 Introduction
- 2 Reward-based learning
- 3 Simulations**
- 4 Conclusion



- EEG data : 20 subjects  $\times$  32 EEG channels  $\times$  220 trials  $\times$  12 row/columns  $\times$  5 flashes per row/column
- data preprocessing. For each trial :
  - 600 ms sample after each flash (100 Hz sampling)
  - Bandpass [1-20] Hz filter
  - common average reference subtraction + channel normalization
  - 5 repetitions : average response calculation per row/column
  - vector construction :  $\mathbf{x} \in \mathbb{R}^{32 \times 60}$ ,  $\|\mathbf{x}\| = 1$
  - set construction :  $\underline{\mathbf{x}}^{\text{row}} = (\mathbf{x}_1^{\text{row}}, \dots, \mathbf{x}_6^{\text{row}})$ ,  $\underline{\mathbf{x}}^{\text{col}} = (\mathbf{x}_1^{\text{col}}, \dots, \mathbf{x}_6^{\text{col}})$
- cross-validation : for one  $(\eta, \lambda)$  couple:
  - learning from scratch :  $w_0 = 0$
  - simulated rewards :  $r \in (r^+, r^-)$  with  $r^+ = 5$ ,  $r^- = -1$ .
  - 1000 simulations  $\times$  20 subject with shuffled spelling order
  - “softmax” and “argmax” classifier variants

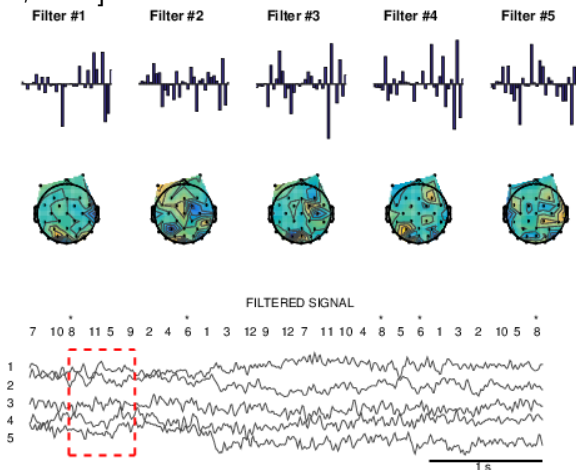


# Softmax/Argmax spelling improvement

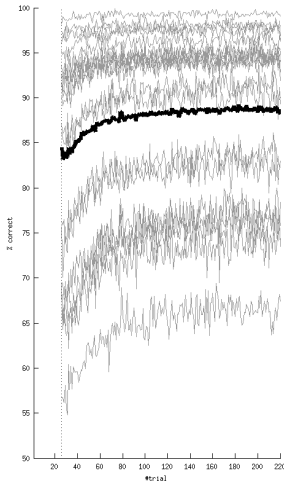


# XDAWN spatial filter

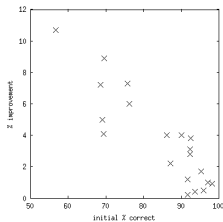
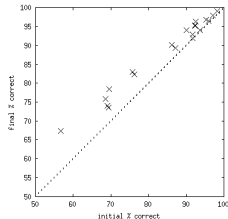
[Rivet et al., 2009]



# Classification improvement after a 25-trials training session



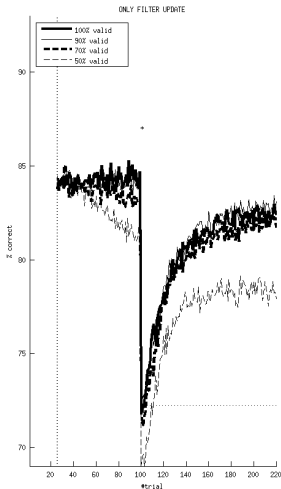
Daucé et al.



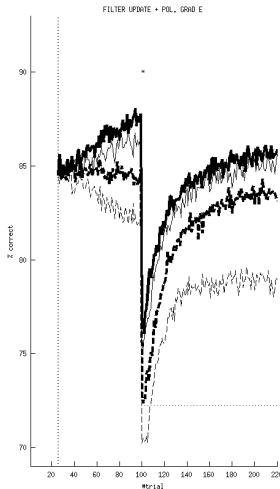
IJCNN 2015

22

# Global recovery after electrode break



Daucé et al.

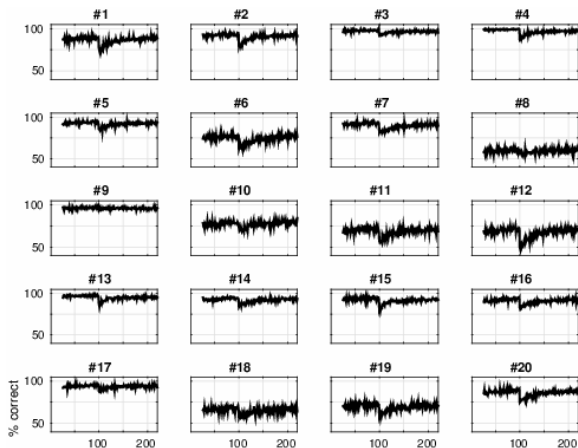


IJCNN 2015

23

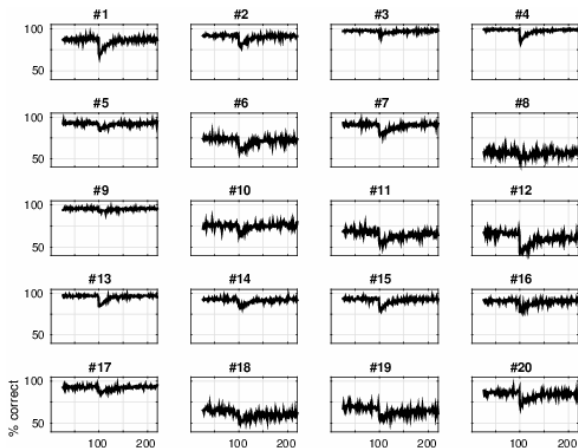


# Individual recovery, feedback noise = 10%





# Individual recovery, feedback noise = 30%



# Plan

- 1 Introduction
- 2 Reward-based learning
- 3 Simulations
- 4 Conclusion**



# Conclusion

- “Instant” reward-based policy gradient descent implements classifier training
- Binary instant rewards allow non-stochastic exploration (argmax choice instead of softmax) with MLE convergence guaranties
- Regularization allows to track environmental changes
- Is efficient in the P300-speller case, but :
  - high variability between subjects
  - reward extraction problem (difficulty shift)
  - reward non-reliability bound = 0.7?



Congedo, M., Barachant, A., and Andreev, A. (2013).

A new generation of brain-computer interface based on riemannian geometry.

*CoRR*, abs/1310.8115.



Daucé, E., Proix, T., and Ralaivola, L. (2013).

Fast classifier adaptation in non-stationary environments: a policy gradient application to the bci p300-speller “oddball” paradigm.

In *proc. of ESANN 2013*, pages 197–202.



Kindermans, P., Verstraeten, D., and Schrauwen, B. (2012a).

A bayesian model for exploiting application constraints to enable unsupervised training of a p300-based bci.

*PloS one*, 7(4).



Kindermans, P.-J. and Schrauwen, B. (2013).

Dynamic stopping in a calibration-less p300 speller.

*5th International Brain-Computer Interface Meeting*.

Dauce et al.

IJCNN 2015

28





Kindermans, P.-J., Verschore, H., Verstraeten, D., and Schrauwen, B. (2012b).

A p300 bci for the masses: Prior information enables instant unsupervised spelling.

*In Advances in Neural Information Processing Systems (NIPS 2012).*



Li, Y. and Guan, C. (2006).

An extended em algorithm for joint feature extraction and classification in brain-computer interfaces.

*Neural Computation*, 18:2730–2761.



Perrin, M. (2012).

*Coadaptation cerveau machine pour une interaction optimale :application au P300-speller.*

PhD thesis, Université Claude Bernard - Lyon 1.



Rivet, B., Souloumiac, A., Attina, V., and Gibert, G. (2009).



xDAWN algorithm to enhance evoked potentials: application to brain-computer interface.

*IEEE Transactions on Biomedical Engineering*, 56(8):2035–43.

