

Programming by Feedback

Riad Akrou, Marc Schoenauer, Jean-Christophe Souplet, Michèle Sebag
TAO, CNRS – INRIA – LRI, Université Paris-Sud, France



Ages of programming

- 1970s Specifications
- 1990s Examples
- 2010s Interactive Learning and Optimization

Languages & thm proving
Pattern recognition & ML

- Visual rendering
- Information retrieval
- Robotics

Brochu et al. 2010
Joachims et al., 2012
Knox et al. 2010, Akrou et al., 2012; Wilson et al., 2012; Saxena et al. 2013

Programming by Feedback: Active Computer, Judging Expert



Knowledge-constrained

Computation, memory-constrained

Algorithm: Iterate

- 1 Computer presents the expert with a pair of solutions y_{t_1}, y_{t_2}
- 2 Expert emits preferences $y_{t_1} \succ y_{t_2}$
- 3 Computer approximates Expert's utility function
- 4 Computer searches for next solution pair *with best expected posterior utility*

Formally

\mathcal{X} Search space, solution space (controllers in RL)

\mathcal{Y} Evaluation space (trajectories)

True utility function U^* (with unknown w^* in W):

$$U : \mathcal{Y} \mapsto \mathbb{R}, U(y) = \langle w^*, y \rangle$$

Modelling the expert's competence: Noise model

$$\delta \sim U[0, M]$$

Given preference margin $z = \langle w^*, y - y' \rangle$

$$P(y \prec y' \mid w^*, \delta) = \begin{cases} 0 & \text{if } z < -\delta \\ 1 & \text{if } z > \delta \\ \frac{1+z}{2} & \text{otherwise} \end{cases}$$

Learning the expert's utility function

find θ_t posterior on W

Proposition. Given evidence $\mathcal{U}_t = \{y_0, y_1, \dots; (y_{i_1} \succ y_{i_2}), i = 1 \dots t\}$,

$$\begin{aligned} \theta_t(w) &\propto \prod_{i=1,t} P(y_{i_1} \succ y_{i_2} \mid w) \\ &= \prod_{i=1,t} \left(\frac{1}{2} + \frac{w_i}{2M} \left(1 + \log \frac{M}{|w_i|} \right) \right) \end{aligned}$$

with $w_i = \langle w, y_{i_1} - y_{i_2} \rangle$, capped to $[-M, M]$.

Most informative demonstrations (y, y') ?

Expected utility of selection:

$$\begin{aligned} EUS(y, y') &= \mathbb{E}_{\theta_t}[\langle w, y - y' \rangle > 0] \cdot U(\theta_t^+, y) \\ &\quad + \mathbb{E}_{\theta_t}[\langle w, y - y' \rangle < 0] \cdot U(\theta_t^-, y') \end{aligned}$$

Expected posterior utility:

$$\begin{aligned} EPU(y, y') &= \mathbb{E}_{\theta_t}[\langle w, y - y' \rangle > 0] \cdot \max_y U(\theta_t^+, y) \\ &\quad + \mathbb{E}_{\theta_t}[\langle w, y - y' \rangle < 0] \cdot \max_y U(\theta_t^-, y) \\ &= \mathbb{E}_{\theta_t}[\langle w, y - y' \rangle > 0] \cdot U(\theta_t^+, y^*) \\ &\quad + \mathbb{E}_{\theta_t}[\langle w, y - y' \rangle < 0] \cdot U(\theta_t^-, y'^*) \end{aligned}$$

Therefore

$$\text{Find } \arg\max EUS(y, y')$$

Viappiani & Boutillier 10

Optimization in the demonstration space

Proposition. $EUS^{\text{noiseless}}(y, y') - L \leq EUS^{\text{noise}}(y, y') \leq EUS^{\text{noiseless}}(y, y')$

Proposition. $EUS_t^{\text{noiseless}} - L \leq EPU_t^{\text{noise}} \leq EUS_t^{\text{noiseless}} + L$

Optimization in the solution space

- Find $\arg\max EUS(y_t^*, y)$ decreases cognitive burden

- Given the mapping Φ : Solution \mapsto Demonstration space,

$$\mathbb{E}_{\Phi}[EUS^{NL}(\Phi(x), y_t^*)] \geq EUS^{NL}(\bar{y}, y_t^*)$$

- Draw $w_0 \sim \theta_t$ and let $x_1 = \arg\max \langle w_0, \bar{y} \rangle$
Iteratively, find $x_{i+1} = \arg\max \langle \mathbb{E}_{\theta_i}[w], \bar{y} \rangle$, with θ_i posterior with $\bar{y}_i > \bar{y}_t^*$.

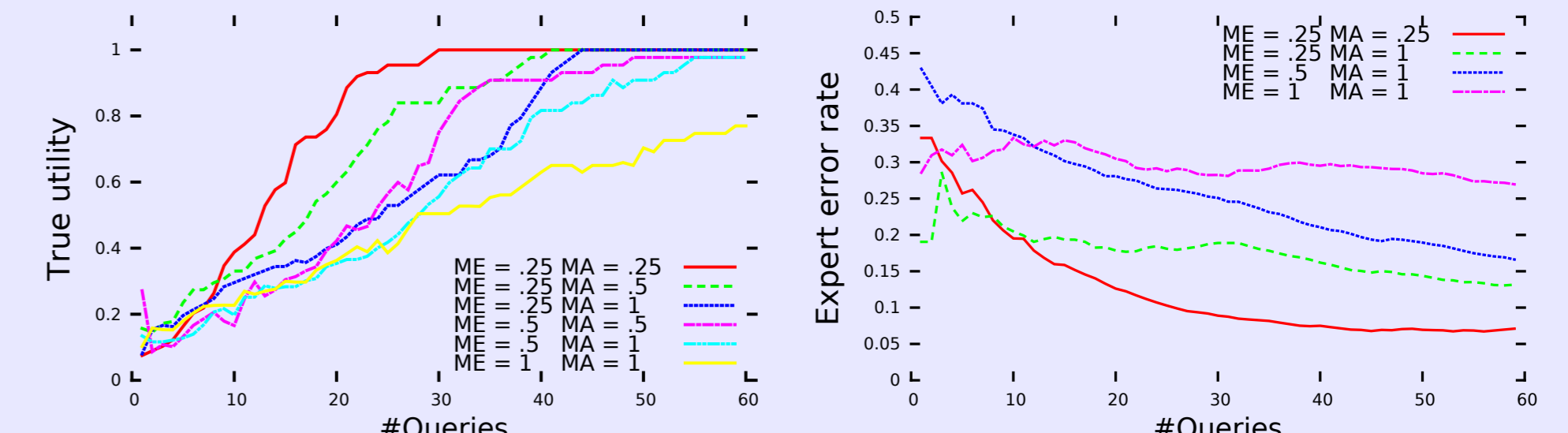
Proposition. The sequence monotonically converges toward a local optimum of $EUS^{\text{noiseless}}$.

Experimental study

Grid world: Discrete Case, no Generative Model

25 states, 5 actions, horizon 300, 50% transition motionless

	...	1/4	1/2	1
1/64			1/4	1/2
1/128	1/64			1
1/256	1/128	1/64		



True w^* on gridworld

True utility of x_t

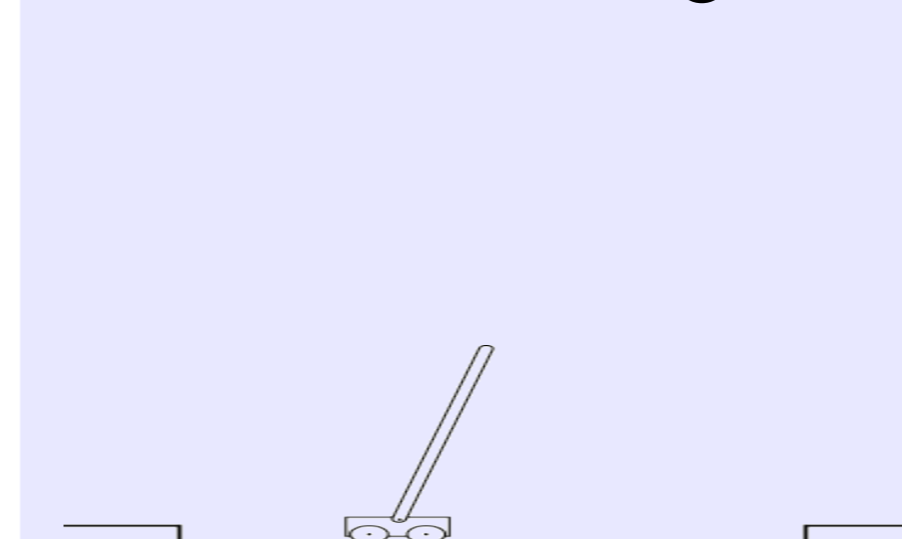
user's mistakes

Sensitivity study wrt user's competence (M_E) and computer trust (M_A):
a cumulative (dis)advantage phenomenon

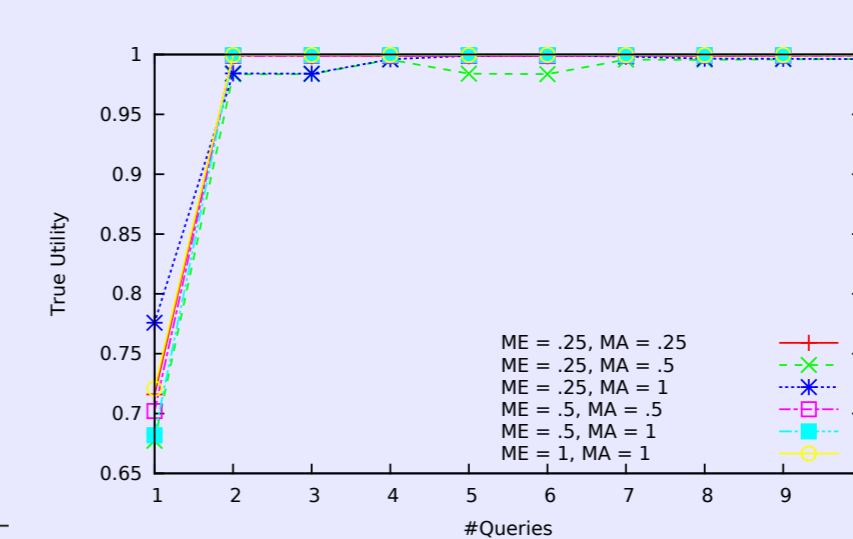
The number of (emulated) user mistakes *increases* as the computer underestimates the user's competence. For low M_A , the computer learns faster, submits more relevant demonstrations to the user, thus priming a virtuous educational process.

The Cartpole: Continuous Case, no Generative Model

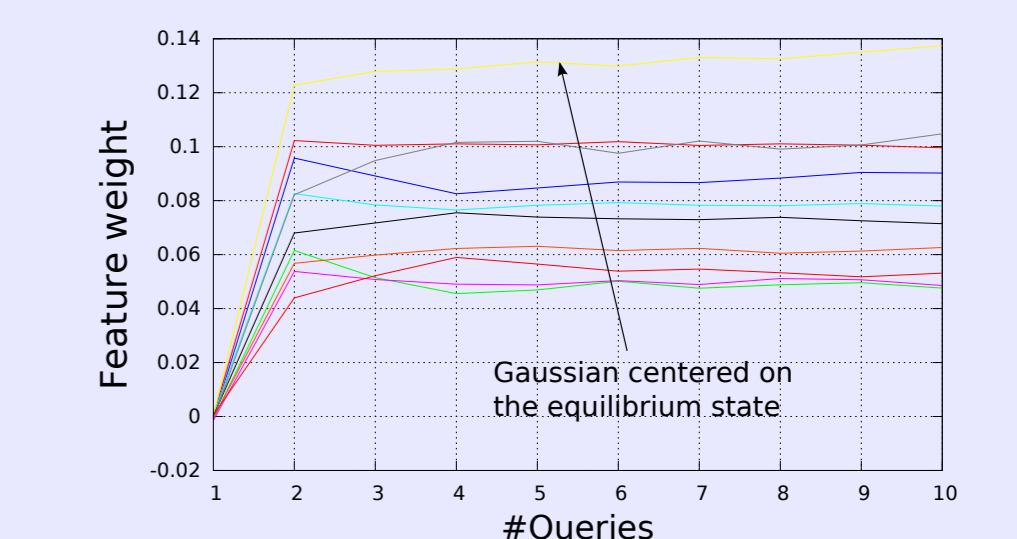
State space \mathbb{R}^2 (the angle and angular velocity of the pendulum), 3 actions; demonstration length 3,000.



Cartpole



True utility of x_t



Estimated utility of features

Demonstration space $\mathcal{Y} = \mathbb{R}^9$ (feature = Gaussian in state space). Simulated user's feedback: best demonstration is the longest one (+ noise). True utility: fraction of the demonstration in equilibrium.

Two interactions required on average to solve the cartpole problem, irrespective of the noise model hyper-parameters

The Bicycle: Continuous Case, with Generative Model

State space \mathbb{R}^4 , action space

\mathbb{R}^2 , demonstration length $\leq 30,000$. Solution

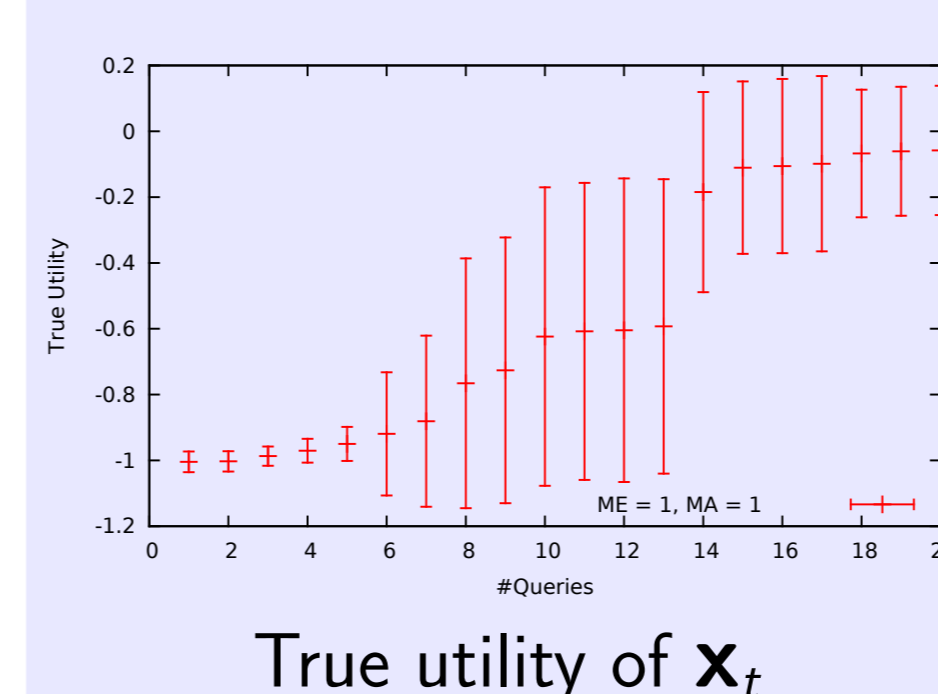
space $\mathcal{X} \subseteq \mathbb{R}^{210}$ (weight vector of a 1-layer

feedforward NN with 4 input, 29 hidden neurons

and 2 output). Optimization component: CMA-ES

black box optimization (Hansen et al., 2001)

as LSPI fails with the estimated utility function.

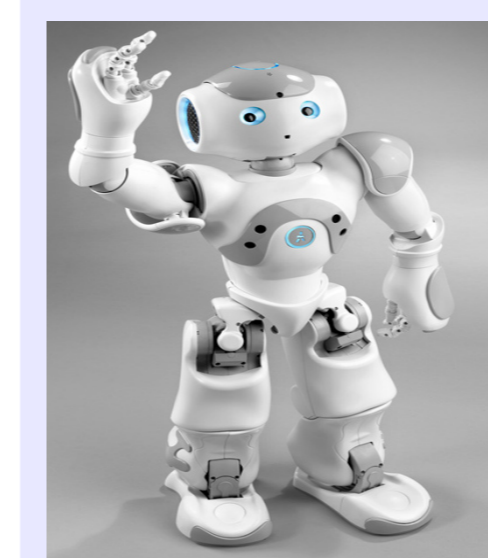


True utility of x_t

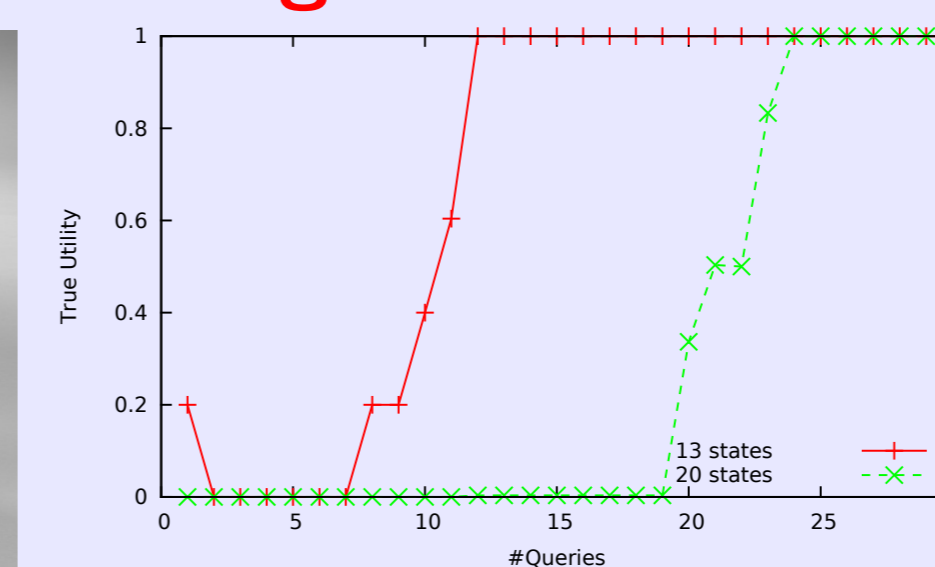
15 interactions required on average to solve the bicycle problem for the low noise setting ($M_E = M_A = 1$).

Improves on the state of the art: circa 20 queries required with discrete action space in Wilson et al. 2012; explained from the more compact search space (V as opposed to Q).

The Nao: Training in-situ



The Nao robot



True utility of x_t

Goal: reaching a given state. Transition matrix estimated from 1,000 random (s, a, s') triplets. Demonstration length 10, initial state is fixed.

Discussion and perspectives

- Feasibility of the **Programming by Feedback** paradigm.



One could carry through the organization of an intelligent machine with only two interfering inputs, one for pleasure or reward, and the other for pain or punishment. 1950

- Importance of noise: all experts (and us) made mistakes. The computer must trust the expert to a limited extent; but computer distrust increases the expert mistakes.
- Next: Identifying the sub-behaviors responsible for the expert's like/dislikes, taking inspiration from Wilson et al. 2012.
- Next: Accounting for the variance of the demonstrations associated to a solution as a secondary criterion (or switching to multiple-objective optimization).