



Solving weak transduction with expectation maximization

Yuri A. Ivanov*, Bruce M. Blumberg

MIT Media Laboratory, 20 Ames St., Cambridge, MA 02139, USA

Abstract

In this paper we describe an algorithm designed for learning perceptual organization of an autonomous agent. The learning algorithm performs incremental clustering of a perceptual input under reward. The distribution of the input samples is modeled by a Gaussian mixture density, which serves as a state space for the policy learning algorithm. The agent learns to select actions in response to the presented stimuli simultaneously with estimating the parameters of the input mixture density. The feedback from the environment is given to the agent in the form of a scalar value, or a *reward*, which represents the utility of a particular clustering configuration for the action selection. The setting of the learning task makes it impossible to use supervised or partially supervised techniques to estimate the parameters of the input density. The paper introduces the notion of *weak transduction* and shows a solution to it using an expectation maximization-based framework. © 2002 Published by Elsevier Science B.V.

Keywords: EM algorithm; Weak transduction; Weakly labeled data

1. Introduction

Designing the perceptual system for autonomous agents is often a difficult task. The autonomy of the agent implies that the behavior of the agent is continuously modified as the agent collects new experiences and observes outcomes of the actions that it decides to take. A typical setting of the reinforcement learning task formulates a reward function, or a function according to which the world outside the agent's body rewards or punishes the agent for taking actions.

Since the reward function is rarely known, the agent has to experiment with the world to approximate this function from a just few samples and find a *policy*, which lets the agent select its actions such that the average reward intake is maximized.

To give a concrete example, imagine the situation where an agent is trained by a human trainer to respond

to a set of voice commands. After hearing an utterance the agent performs an action. The trainer would like to have the agent respond correctly by providing (possibly noisy) rewards and punishments after observing the actions that it performs in response. In this scenario the agent needs to learn two things: (a) parameterized equivalence classes in the space of utterances; and (b) what action to select upon observing a sample from one of these classes. The first problem is that of *clustering under reward*, while the second is the typical *policy learning* task of reinforcement learning.

There exist a number of algorithms permitting learning of policies of action selection given that the perceptual system provides a good set of features. But how can such features be efficiently estimated while the policy learning is taking place? This paper focuses on the solution to the problem by embedding a simple associative search algorithm into an expectation maximization (EM) paradigm.

The task of *on-line* estimation of the perceptual organization and the policy of action selection is cast

* Corresponding author.

E-mail address: yivanov@media.mit.edu (Y.A. Ivanov).

56 here as a problem of a *multi-armed bandit with hidden*
 57 *state* and is solved iteratively within the EM frame-
 58 work. The hidden state is represented by a parame-
 59 terized probability distribution over states tied to the
 60 reward. The parameterization is formally justified, al-
 61 lowing for smooth blending between likelihood- and
 62 reward-based costs.

63 In addition, in this paper we introduces the notion
 64 of *weak transduction* in order to properly place the
 65 solution among the existing techniques used for unsu-
 66 pervised and transductive problems.

67 The paper proceeds as follows: after introducing
 68 the multi-state bandit problem, Section 3 describes a
 69 modification of the reinforcement pursuit algorithm
 70 that allows us to include the estimation of the hidden
 71 state. Section 4 justifies modifications to the EM al-
 72 gorithm, which permits the inclusion of the reward in-
 73 formation into the parameter estimation and to solve
 74 the problem of learning perceptual organization along
 75 with policy learning. The two parts of the estimation
 76 procedure are brought together in Section 5, which
 77 presents the complete reinforcement-driven EM algo-
 78 rithm. Experiments with this algorithm showing the
 79 results for different objectives are presented in Section
 80 6. The paper concludes with Section 7, which points
 81 out contributions and some of the problems with the
 82 algorithm.

83 1.1. Related work

84 The work of this paper is based on the EM algo-
 85 rithm [1,6,12] extended to situate it within the context
 86 of reinforcement learning and to take advantage of
 87 the additional information that is available as a result
 88 of interaction with the environment. Neal and Hinton
 89 [4] show a view of the EM algorithm that makes the
 90 extensions made in this paper possible. This view is
 91 concisely presented by Minka in [3].

92 At a certain parameter setting and binary reward,
 93 the algorithm shown here can be viewed as an on-line
 94 version of the λ EM, presented by Nigam et al. [5], for
 95 learning with partially labeled data (albeit for a Gaus-
 96 sian mixture) and transductive inference [10]. How-
 97 ever, the problem solved by the algorithm described
 98 here is in general more difficult than the problem
 99 of transductive clustering, which Nigam's algorithm
 100 solves as it does not have access to exact labels for
 101 the input data.

102 To learn the policy of the action selection the learn-
 103 ing algorithm developed here uses an N -armed bandit
 104 model (see, e.g., [8]). The multi-state policy is esti-
 105 mated with the aid of the reinforcement pursuit algo-
 106 rithm of Thathachar and Sastry [9], which is applied
 107 to a set of states simultaneously.

108 A problem similar to the one presented here was
 109 explored by Likas [2], who used a variant of the RE-
 110 INFORCE algorithm [11] to learn vector quantization
 111 on the batch data, aided by a reward signal.

112 The technique of probability matching for reinforce-
 113 ment learning used here is similar to that shown by
 114 Sabes and Jordan [7]. Using this technique, the algo-
 115 rithm presented here constructs a reward-dependent
 116 probability distribution to guide the algorithm towards
 117 the configuration resulting in higher value of the ex-
 118 pected reward.

119 2. Agent training

120 Let us for a moment return to the example of the
 121 agent training and introduce some notation. The train-
 122 ing consists of a series of steps, where superscript in-
 123 dicates the number of the step. At each step the trainer
 124 randomly selects a state, s , out of a finite set of states
 125 (see Fig. 2a). From that state the trainer produces an
 126 observation for the agent, x^n , according to the proba-
 127 bility distribution $p^*(x|s)$. The chosen state switches
 128 the trainer to one of its "reward modes", where any
 129 action, a^n , that the agent might perform in response
 130 to x^n will be rewarded, punished or ignored, by dis-
 131 pensing to it a reward, r , from a probability distribu-
 132 tion conditioned on both trainer's state and the agent's
 133 action, $p^*(r|s, a)$.

134 To clarify this, let us use an example. "A state"
 135 of the trainer corresponds to the trainer's desire to
 136 evoke a particular response from the agent. For ex-
 137 ample, the trainer might choose to have the agent
 138 go away. The external manifestation of that state
 139 ($s = \text{want-it-to-go-away}$) might be a verbal command
 140 "Go away!" ($x^n = \text{"Go away!"}$). At this point the
 141 agent might erroneously choose to come up closer
 142 ($a^n = \text{approach-trainer}$). The trainer lets the agent
 143 to get away with it this time by not punishing it
 144 (drawing a value of 0 from the probability distribu-
 145 tion $p^*(r|s = \text{want-it-to-go-away}, a = \text{approached-}$
 146 trainer).

Environment model

- (1) Nature can be in one of a discrete set of states $S = \{s_i\}_{i=1}^K$;
- (2) Nature selects a state, s_i , from a probability distribution $p^*(s)$;
- (3) From a conditional distribution $p^*(x|s = s_i)$ nature generates an observation x^n ;
- (4) Upon observing the agent taking an action a^n , nature produces a reward from the distribution $p^*(r|s = s_i, a = a^n)$.

The goal of the agent is to learn how to respond to the commands given to it by the trainer to maximize the intake of the reward. To solve this problem we propose an agent architecture that consists of two parts—a perceptual model and a policy. The perceptual model is a set of clusters in the input space, while the policy is a probability distribution, which for a given state gives a probability with which every agent's action would produce a reward.

The agent's task is then to estimate parameters of the model shown in Fig. 2b. In that figure $p(x|s)$ is the probability model of a category s , which is subject to parameter estimation; $p(s|x^n)$ is a belief state, calculated for the incoming observation from current model parameters; $p(a|s)$ is a policy, that is related to the past correlations between state, action and reward; and $p(a|x^n)$ is the probability distribution from which the action is selected upon observing x^n .

It is important to note that the perceptual organization of the agent cannot be formed by providing it with supervised examples of correct perceptual class assignment. The trainer can only judge the quality of the agent's perceptual organization based in the actions the agent selects in response to the external stimuli.

3. Estimation of the associative policy

The policy in the problem addressed here has a special and simple character—since all observations are assumed to be equally likely and independent—the trainer can produce any observation (e.g., utterance or gesture) at will—there is no need to keep the history as a part of the action context. In other words, the agent's behavior does not change the state of the world, which only changes at random. The task of the policy estimation in such a simplified setting is termed *associative search* and requires no history of prior observations to be stored. However, it is complicated by the necessity of perceptual learning, or state estimation,

which places it in the category of *associative search with hidden state*.

If the state space is modeled with a mixture distribution, then the problem can be described as follows: given an observation, estimate the state of the world from a finite set of states, $S = \{s_i\}$. Given the belief about the state membership, select an action (label), which will result in the maximum amount of expected payoff received once the action is performed. With that payoff, update the parameters of the policy of the action selection and of the input distribution. This section will deal with solving the problem of policy estimation for such a setup.

3.1. Multi-state bandit problem

Due to the assumption that the observations are independent of each other, this problem can be formulated as a multi-state N -armed bandit [8]. The N -armed bandit is a gambling device with a set of N arms (see Fig. 1). Each arm has a probability distribution associated with it, according to which the sample reward is drawn every time the arm is pulled. Most frequently the reward generating process is assumed to be stationary or, at most, slowly varying.

Now imagine that the bandit can be in any one of M states, each of which have different distributions of the reward. Before each trial the bandit switches to a new state and produces an observation, x^n , from

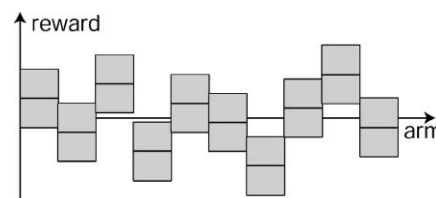


Fig. 1. The 10-armed bandit model. Each of the 10 arms produces a reward by drawing a sample from a corresponding distribution. Each box signifies the reward distribution with some mean (horizontal bar) and variance (height of the box).

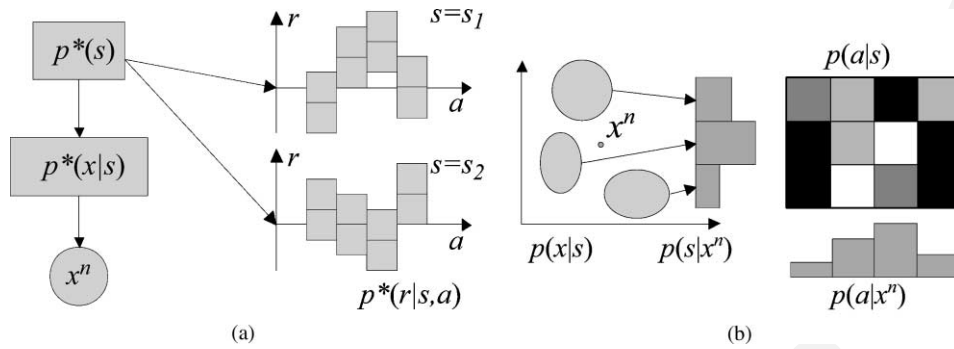


Fig. 2. (a) A generative model of the environment, shown as a 2-state 4-armed bandit. The bandit randomly switches between two states, according to a sample drawn from $p^*(s)$. After selecting the state, s , an observation, x^n is produced from a distribution $p^*(x|s)$. (b) The estimator consists of two parts—a perceptual model and a policy. Upon receiving the observation, x^n , the distribution $p(a|x^n)$ is constructed and an action is selected by drawing a sample from it. Upon delivery of the reward parameters of both the policy and the perceptual model are updated.

213 the distribution associated with this state (see Fig. 2a).
 214 The player's goal is to identify this state and perform
 215 action selection and model update for that state. When
 216 the state is perfectly known the problem reduces to
 217 M independent N -armed bandits. It is more difficult
 218 when the state is hidden and must be estimated.

219 *3.2. Solutions with known state*

220 When the state is exactly known, then the solu-
 221 tion for the multi-state bandit is achieved by independ-
 222 ently solving a set of single-state bandits. A variety
 223 of *action-value* methods, such as *sample average*,
 224 *reinforcement comparison* and *reinforcement pursuit*,
 225 have been proposed to solve the single-state bandit
 226 problem.¹ The general idea is to stochastically search
 227 the action space while updating the estimate of the re-
 228 ward function. A probability distribution over the ac-
 229 tion space (*action preference*) is built based on this es-
 230 timate and action selection is done via sampling from
 231 this distribution.

232 The simplest pursuit method, adapted for the
 233 multi-state agent, maintains an estimate of the pay-
 234 off structure of the bandit via action-value func-
 235 tion, $Q_t(a, s)$. This function is updated at each step

¹ A classical method for solving bandit problems, which includes balancing of exploration with exploitation involves computation of the so called Gittins indices. This method provides an optimal solution to a large class of problems, but assumes the knowledge of prior distribution of possible problems.

based on the reward received from the bandit after 236
 pulling the arm a by, e.g., an exponentially weighted 237
 sample-average method 238

$$Q_t(a, s) = Q_{t-1}(a, s) + \alpha(r - Q_{t-1}(a, s)). \quad (1) \quad 239$$

Based on the value of $Q_t(a, s)$, the pursuit method 240
 updates its action preference model, $\hat{p}_t(a|s)$, such that 241
 the action with the highest value of $Q_t(a, s)$ increases 242
 the probability of being selected by a small fraction, 243
 γ . Actions that are currently found to be suboptimal 244
 decrease their probability correspondingly. Let $a_{t+1}^* =$ 245
 $\arg \max_a (Q_t(a, s))$, then 246

$$\begin{aligned} \hat{p}_{t+1}(a^*|s) &= \hat{p}_t(a^*|s) + \gamma(1 - \hat{p}_t(a^*|s)), & 248 \\ \hat{p}_{t+1}(a|s) &= \hat{p}_t(a|s) + \gamma(0 - \hat{p}_t(a|s)), \quad \forall a \neq a^*. & 249 \end{aligned} \quad (2) \quad 250$$

The convergence of the pursuit method is dependent 251
 upon values of α and γ , which in all experiments of 252
 this paper are set to be $\alpha = 0.1$ and $\gamma = 0.01$. In addi- 253
 tion, it is readily combined with ϵ -greedy techniques 254
 to allow for non-stationary environments. 255

256 *3.3. Solutions with hidden state*

In the presence of the hidden state the problem of 257
 estimating the optimal action becomes more difficult. 258
 The uncertainty about the state can be dealt with by 259
 distributing the reward proportionally to the current 260

261 belief about the state membership of the observation
262 x^n .

263 Most of the bandit search algorithms allow for formulating a policy, or a probability distribution over
264 actions, given a state, $p(a|s)$. This is an arbitrary distribution which only expresses the current estimate of
265 “action preferences”. The action is selected by sampling the conditional probability distribution $p(a|x^n)$,
266 which can be calculated from the belief state and the policy, by marginalizing the joint, $p(a, s|x^n)$
267

$$272 \quad p(a|x^n) = \sum_s p(a, s|x^n) = \sum_s p(a|s, x^n)p(s|x^n) \\ 273 \quad = \sum_s p(a|s)p(s|x^n). \quad (3)$$

274 The action selection now takes into account the uncertainty about the state, encoded in the state posterior.
275 For the purpose of bandit updates, the reward is distributed among M bandits in proportion to their contribution to $p(a|x^n)$
276

$$280 \quad Q_t(a, s) = Q_{t-1}(a, s) + \alpha(rp(s|x^n) - Q_{t-1}(a, s)). \\ 281 \quad (4)$$

282 The action preference update equations, Eq. (2) are left unchanged.
283

284 4. Clustering under reward

285 Given the agent’s architecture in Fig. 2b it is clear that the agent is not provided with direct supervision
286 for the task of state estimation. This is because the feedback from the trainer rewards its action selection
287 and not the sample classification in its perceptual system. On the other hand, for the purposes of perceptual
288 clustering this situation is better than having no feedback at all, since it does provide some degree of
289 guidance, however small.
290

291 State estimation under reward can be performed with the aid of the EM algorithm, which is often used
292 for unsupervised clustering. This section introduces a technique for including the reward function into
293 the EM re-estimation procedure. The new objective function is simply implemented in the EM framework
294 while allowing the algorithm to “fall back” to the unsupervised mode if no reward is provided.
295

4.1. Weak transduction

302

The EM algorithm is a powerful tool for solving supervised and transductive problems. It is often used
303 as a clustering algorithm with the objective of maximizing the likelihood of the data. This is a good heuristic
304 to use for learning perceptual organization when no other evidence is available. However, by using an
305 *unsupervised* technique for learning the perceptual organization, one disregards its utility for the agent.
306

307 The utility of a perceptual configuration is measured by the reward that the agent collects while using it.
308 Therefore, an algorithm is sought, which while capable of learning from patterns in the input data alone,
309 can be “directed” with the reward to choose a different configuration providing higher payoff. That is,
310 the solution should be an EM-type algorithm, which would allow the inclusion of reward into its objective
311 for state estimation, while learning the policy of action selection.
312

313 The EM algorithm is frequently used for *unsupervised* clustering of data by spatial proximity in the
314 space of features. For a given number of clusters the algorithm proceeds iteratively first to calculate from
315 the current cluster statistics the probability of data to be generated by each cluster, a state posterior, $p(s|x)$;
316 and then to average the data, weighted by this posterior to update cluster statistics.
317

318 When the data comes with the known state attribution, $s^n = z$, then the posterior of each data point
319 turns into a deterministic function, having 1 at the slot of the corresponding state and 0 elsewhere:
320

$$321 \quad p(s^n = z|x^n) = 1, \quad p(s^n \neq z|x^n) = 0. \quad (5) \quad 322$$

323 Averaging with respect to this posterior, let us call it $p_{01}(s|x)$, results in the parameter estimation to decompose
324 into several independent *supervised* estimation problems.
325

326 When only part of the data has an exact label, $s^n = z$, then the solution to the clustering problem results
327 in a mixture of the supervised solution for the labeled data and the unsupervised solution for the unlabeled
328 set. This is an example of *transduction*, where the knowledge from the labeled data is transduced onto the
329 unlabeled. The setting of a transductive classification is illustrated in Fig. 3a. In the figure the supervised
330 solution disregards the additional indirect information provided by the density of the unlabeled data, while
331

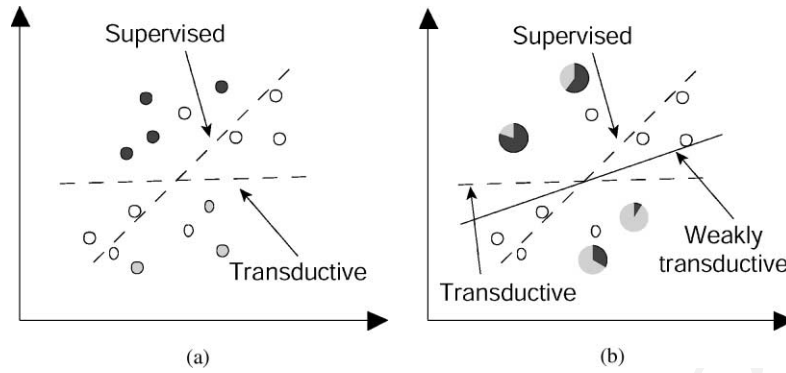


Fig. 3. (a) Illustration of transductive inference for classification. Empty circles indicate the data with unknown class membership while the colored circles indicate the labeled data, which belongs to one of two classes. Classification boundaries are different for estimators with and without unlabeled data. (b) Weak transduction has no direct indication of the class label, but a probability distribution over labels. Clear circles, again, show the unlabeled data, while large colored ones show the data that is weakly labeled. The area of the circle shaded with a particular color is proportional to the probability with which this data point belongs to this class.

348 transductive inference takes into account the complete
349 data set.

350 Supervised, unsupervised and transductive learn-
351 ing methods view the label information in a binary
352 fashion—it is either present or absent. In contrast, un-
353 der the circumstances of the problem where the knowl-
354 edge about the label is inexact and subjective, the situ-
355 ation is a bit worse than in the transductive setting,
356 but better than unsupervised. With the current setting
357 of the model parameters the posterior $p(s|x^n)$ is com-
358 puted as the state membership of the query point. If,
359 consequently, some value of reward results from this
360 assignment, it indicates the quality of the posterior
361 given the current parameter settings. This is the situa-
362 tion, which the algorithm being described encounters
363 in the task of estimating a state. That is, in line with
364 the above taxonomy, the data is labeled with a *prob-*
365 *ability distribution*, as illustrated in Fig. 3b. It is con-
366 venient to call data labeled in this fashion *weakly la-*
367 *beled*, and the problem—a *weak transduction*. These
368 terms properly place the problem among traditional
369 machine learning tasks and emphasizes its relation to
370 already existing techniques for learning with labeled,
371 unlabeled and partially labeled data [5].

372 4.2. Reward-driven variational bound

373 Typically, the EM algorithm for density estimation
374 is used for unsupervised maximization of the likeli-

hood function of a parametric density model when ob- 375
taining an analytical solution for the gradient in the 376
parameter space is difficult. This is the case when we 377
need to learn parameters of a mixture density. In the 378
algorithm of this paper the input space is represented 379
by a mixture density, $p(x; \theta) = \sum_i p(s_i) p(x|s_i; \theta_i)$, 380
parameters of which, θ , need to be estimated. The goal 381
of the algorithm, however, is to not simply maximize 382
the likelihood of the data, but also take into account 383
the external reward signal if such is present. To do so, 384
in this section a new cost function is justified, which 385
allows for inclusion of the reward in the traditional 386
EM framework. 387

388 4.3. EM as a variational bound optimization 389

$$f(\theta) = \int p(x, s, \theta) \frac{q(s)}{q(s)} ds \geq \prod_s \left(\frac{p(x, s, \theta)}{q(s)} \right)^{q(s)} = g(x, \theta). \quad (6)$$

392 Here, $g(x, \theta)$ is a lower bound of the likelihood, $f(\theta)$,
393 and $q(s)$ is some positive function of s , integrating to
394 1. Typically, for the purposes of optimization of $f(\theta)$,
395 the logarithm of $g(x, \theta)$ is optimized

$$G(x, \theta) = \int q(s) \log p(x, s, \theta) - q(s) \log q(s) ds. \quad (7)$$

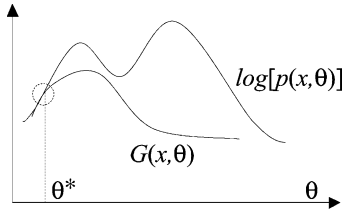


Fig. 4. The step in the direction of the gradient of the lower bound of Eq. (8) is the step in the direction of the gradient of likelihood.

399 It follows from Eq. (6) that for any density $q(s)$,
 400 $G(x, \theta)$ is a lower bound on $\log f(\theta)$. Now the den-
 401 sity $q(s)$ needs to be found, which touches $\log f(\theta)$
 402 at θ . The cost function in Eq. (7) can be re-written as
 403 follows [4]:

$$404 \quad G(x, \theta) = -D(q(s) \| p(s|x, \theta)) + \log f(\theta), \quad (8)$$

405 where $D(p \| q)$ is a Kullback–Leibler divergence be-
 406 tween distributions p and q . From here it is eas-
 407 ily shown that $G(x, \theta) = \log f(\theta)$ when $q(s) =$
 408 $p(s|x, \theta)$, i.e., the bound will be touching the likeli-
 409 hood function at the current θ , as shown in Fig. 4.

410 4.4. Augmented reward bound

411 In order to let EM include the expected reward into
 412 the optimization, the EM bound shown above needs to
 413 be augmented with a reward-dependent term. It is easy
 414 to do using the probability matching technique [7].

415 To learn preferred cluster configurations, one
 416 can consider observation-state pairs and construct a
 417 reward-dependent probability distribution, $p^*(s|x; r)$.
 418 The task of the learning algorithm is to select from a
 419 set of conditional distributions $p(S|\mathcal{X}, \theta)$, aided by
 420 rewards that are provided by the environment for some
 421 of the data points. These rewards can be thought of
 422 as inverse energies—pairs (s, x) receiving higher re-
 423 wards correspond to lower energy states. Energies can
 424 be converted to probabilities via the Boltzmann dis-
 425 tribution, which represents the ideal observation-state
 426 mapping— (s, x) pairs receiving higher rewards being
 427 more likely than pairs receiving low reward. If the
 428 parameters of $p(s|x, \theta)$ are adjusted so that it is close
 429 to $p^*(s|x; r)$, then the output of the algorithm will
 430 typically result in higher rewards.

431 Following this line of reasoning $p^*(s|x; r)$ is made
 432 proportional to the Boltzmann distribution as shown

later in the text. Starting with Eq. (8), an additional 433
 term penalizes the estimator for being different from 434
 this distribution in the posterior: 436

$$437 \quad F(x, \theta) = -D(q(s) \| p(s|x, \theta)) \\ 438 \quad + E_{q(s)} \left[\log \frac{p^*(s|x; r)}{p(s|x, \theta)} \right] + \log f(\theta). \quad (9)$$

When $q(s)$ is set to the posterior distribution, 439
 $p(s|x, \theta)$, the expectation term turns into negative 440
 divergence between the posterior and, $p^*(s|x; r)$: 442

$$443 \quad E_{q(s)} \left[\log \frac{p^*(s|x; r)}{p(s|x, \theta)} \right] \Bigg|_{q(s)=p(s|x, \theta)} \\ 444 \quad = -D(p(s|x, \theta) \| p^*(s|x; r)). \quad (10)$$

In fact this term induces a different but very intuitive 445
 bound for the likelihood maximization (see Appendix 446
 A and Theorem 1 for proof) 447

$$448 \quad F(x, \theta) = -D(q(s) \| p^*(s|x; r)) + \log f(\theta). \quad (11)$$

This function has the same form as Eq. (8), which 449
 implies that for practical purposes one may simply 450
 substitute the EM-induced posterior with the fictitious 451
 probability distribution, $p^*(s|x; r)$. It provides the tra- 452
 ditional bound for the likelihood function in the ab- 453
 sence of the reward. With the reward present, the al- 454
 gorithm performs only a *partial E-step*. However, the 455
 step in the direction of the gradient of this bound leads 456
 uphill in the future expected reward. 457

Now $p^*(s|x; r)$ needs to be constructed in a con- 458
 venient form. The main constraint that should be im- 459
 posed is that the additional term in Eq. (9) vanishes 460
 when after producing a label s for an observation x , 461
 the reward r received from the environment is 0. That 462
 is, 463

$$464 \quad E_{q(s)} \left[\log \frac{p_{r=0}^*(s|x; r)}{p(s|x, \theta)} \right] = 0 \quad (12)$$

which implies that $p_{r=0}^*(s|x; r) = p(s|x, \theta)$. The dis- 465
 tribution $p_{r=0}^*(s|x; r)$ can be set to be proportional to 466
 the Boltzmann distribution: 467

$$468 \quad p^*(s|x; r) = \frac{p(s|x, \theta) \exp(\beta r p(s|x, \theta))}{Z_\beta(x, \theta)}. \quad (13)$$

This form of $p^*(s|x; r)$ is used throughout the rest of 469
 this paper. 470

The resulting bound is illustrated in Fig. 5. The 471
 augmented bound behaves just like the traditional EM 472

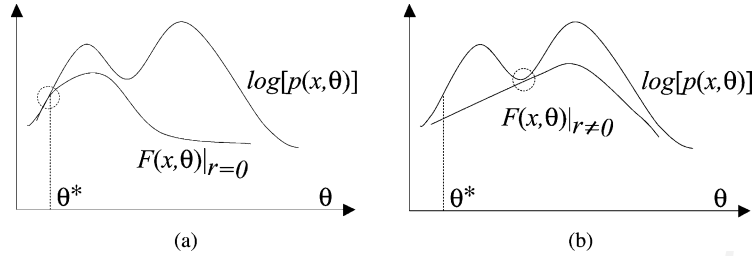


Fig. 5. (a) The augmented bound behaves just like the traditional EM bound when no reward is present. (b) With the reward present the bound is no longer in contact with the likelihood at the current parameter setting, leading uphill in the expected reward.

473 bound when no reward is present. With the reward
 474 present, the bound is no longer in contact with the like-
 475 lihood at the current parameter setting, leading uphill
 476 in the expected reward. The point of contact with the
 477 bound is the value of parameter at which the posterior
 478 $p(s|x^n)$ equals $p^*(s|x; r)$.

479 **5. Reward-driven EM**

480 Now the two parts of the estimation procedure can
 481 be joined to get the complete solution to perceptual
 482 learning under reward. The algorithm is shown below
 483 and is illustrated in Fig. 6.

484 The algorithm folds the action selection policy esti-
 485 mation into the *expectation* step of the EM algorithm
 486 while using the immediate reward signal to control the
 487 entropy of the posterior for the *maximization* step. The
 488 algorithm is iterative and incremental, performing one
 489 iteration per data point, keeping only the sufficient
 490 statistics about the density function of the input space.
 491 The goal of the algorithm is to estimate the structure
 492 shown in Fig. 2. It proceeds as follows:

REM algorithm

- (1) *Initialize*: set parameters to starting values;
for each new data point:
- (2) *E-step*:
 - (a) calculate $p(s|x^n)$ using the Bayes rule and the
current parameters of the
observation model, $p(x)$;
 - (b) *Forward pass*:
 - (i) compute $p(a|x^n)$ (Eq. (3));
 - (ii) select an arm by sampling $p(a|x^n)$;
 - (c) *Backward pass*:
 - (i) collect reward and distribute it among the
states in fractions of $p(s|x^n)$;
 - (ii) calculate $p^*(s|x^n, r^n)$ (Eq. (13));
- (3) *M-step*: maximize the resulting bound, Eq. (A.1)

In the forward pass of the algorithm the processing
 breaks out of the EMs expectation step to select an
 action and update the bandit model as shown in Fig. 6.
 The yielded payoff serves as a control parameter for
 the EM.

6. Experiments

The experimental analysis of the algorithm pre-
 sented in this paper is performed on a series of tasks
 of increased difficulty. The first experiment does not
 include policy learning and is designed to simply test
 estimation of the perceptual model alone for a fixed
 optimal policy. Next, two experiments are performed,
 which involve the policy estimation. In the first exper-
 iment the reward is delivered by a bandit with only
 one arm per state producing a unit of reward. In the
 second experiment the binary restriction on the bandit

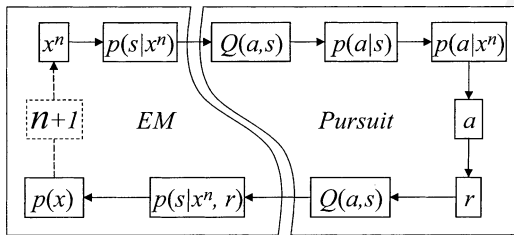


Fig. 6. The reward-driven perceptual learning algorithm breaks out of the expectation step of EM to compute the improved posterior. Then the parameter estimation is performed with respect to $p(s|x^n, r)$.

510 is removed allowing each arm to produce some value
 511 of reward, positive or negative. Finally, an experiment
 512 is performed with a variation on the reward structure
 513 such that the reward reinforces arbitrary objective, not
 514 related to the likelihood of the data.

515 6.1. EM for state estimation

516 The first experiment confirms the conclusions of the
 517 previous section, showing that it is in fact possible to
 518 use the EM framework for partially supervised tasks.
 519 It has to be shown that, given the context of the clas-
 520 sification task, the algorithm will result in choosing
 521 the clustering configuration that provides a higher ex-
 522 pected reward.

523 In the experiments of this section, the performance
 524 of the algorithm is compared with the traditional EM.
 525 However, it should be understood that this comparison
 526 is for reference only, as the EM is not designed to
 527 perform the task that REM is targeting and can only
 528 provide the “worst case” performance.

529 As a source of the data a Gaussian mixture, $q(x) =$
 530 $\sum q(s)q(x|s)$ is used. The algorithm estimates the
 531 density $p(x) = \sum p(s)p(x|s)$ by adjusting its param-
 532 eters in an on-line fashion, upon seeing every data
 533 point, x^n . The reward is delivered after an attempt to
 534 classify x^n to be generated by a particular component
 535 of $p(x|s_i)$. The experiment proceeds as follows:

-
- (1) Initialize the generator mixture, $q(x)$: for each state, s_i , randomly select a Gaussian observation model— $\mu_i \sim N(\mathbf{0}, 2I)$ and $\sigma_i = I$;
 - (2) Iterate:
 - (a) randomly choose a generator state, s_k ;
 - (b) generate an observation, x^n , distributed with μ_k and σ_k ;
 - (c) using current parameters of the model, $p(x)$, select a label l^n ;
 - (d) if $l^n = s_k$, deliver a reward of 1, otherwise -1 ;
 - (e) update parameters of the model
 - (i) compute $p^*(s|x^n; \hat{\nu})$ via Eq. (13);
 - (ii) perform the E-step of the EM algorithm using $p^*(s|x^n; \hat{\nu})$ in place of $p(s|x^n)$.
-

537 The results of the incremental reinforced binary
 538 classification experiments are shown in Fig. 7. The top
 539 plot shows the attained likelihood of the data after a
 540 number of randomly generated samples. The horizon-

tal axis shows the number of iterations (data points
 541 seen so far) with the likelihood plotted along the verti-
 542 cal axis. It is curious to see that the unguided EM (with
 543 $\beta = 0$) attains the lowest likelihood. This is partially
 544 due to the fact that the EM is more likely to get stuck
 545 in the local maxima, while the reward signal delivers
 546 some extra energy for the algorithm to get out of it.

547 The intuition behind choosing the parameter β is
 548 that as it increases, the entropy of the probability dis-
 549 tribution from which a label is selected drops. Char-
 550 acteristic behavior of the algorithm can be observed
 551 at extreme values of β with $\beta = -\infty$, the distribution
 552 over labels is uniform and the label selection is per-
 553 formed purely by chance, with no regard to neither the
 554 reward nor mixture parameters. At $\beta = 0$ the distribu-
 555 tion over labels exactly equals to the mixture posterior,
 556 that is the algorithm disregards the reward completely,
 557 performing the unsupervised parameter estimation as
 558 mixture parameters dictate. Setting β to $+\infty$ results
 559 in a “winner-take-all” label assignment.
 560

561 The second plot in Fig. 7 complements the likeli-
 562 hood plot by showing the classification accuracy of
 563 the algorithm at different values of the parameter β . It
 564 is expected that the accuracy of the EM used for clas-
 565 sification should not be better than chance, since even
 566 when EM converges to the correct set of classes it does
 567 not care which source cluster corresponds to which
 568 estimated component. Positive values of the param-
 569 eter β drive the extended EM towards correct labeling,
 570 while negative β drives the algorithm away from it,
 571 as can be seen in the accuracy plot. It is interesting
 572 that none of the settings of the parameter β result in
 573 the optimal accuracy of 1. There are two reasons for
 574 this. First, any fixed value of β less than ∞ will cause
 575 a sub-optimal label to be selected, albeit with small
 576 probability. The second reason is related to the fact
 577 that even optimal Bayes classifier will not achieve the
 578 perfect classification rate as randomly placed source
 579 Gaussian components may significantly overlap.

580 The influence of β is further illustrated in Fig. 8.
 581 The figure shows the resulting clustering attained with
 582 different values of β . It can be seen that the clusters
 583 for positive and negative values of β have opposite
 584 labeling while zero-valued β is labeled by chance. In
 585 this run the source distribution has the component 1
 586 (green) at the position (5, 5) and component 2 (red) at
 587 (0, 0), which is correctly identified by the algorithm
 588 with large positive value of β .

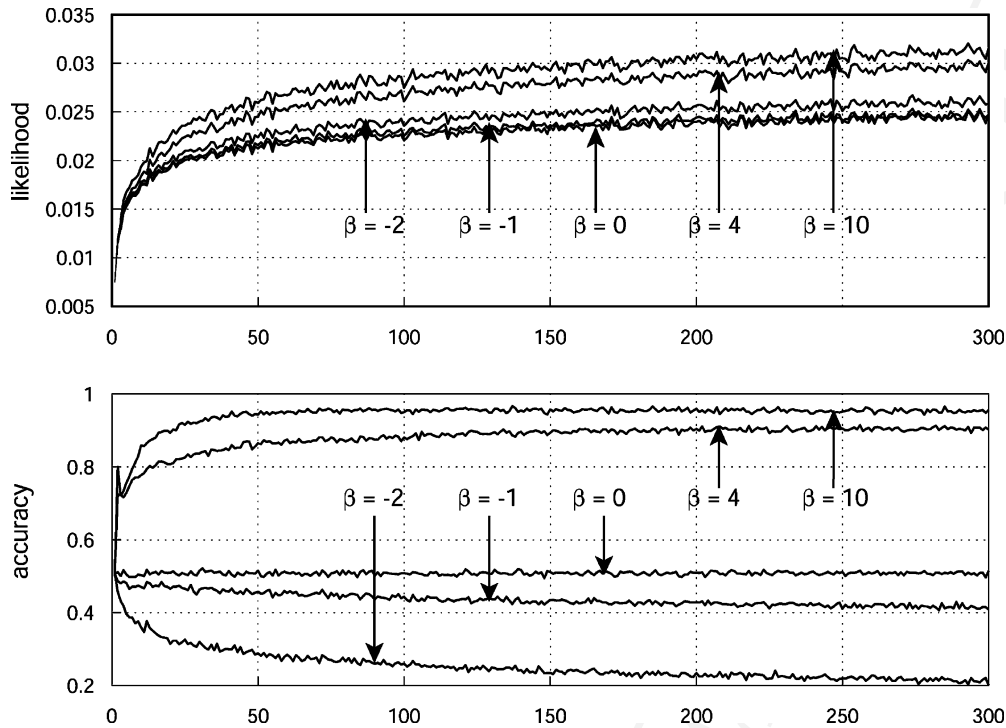


Fig. 7. Performance of the REM averaged over 1000 runs for different values of the parameter β as compared with EM. Curiously, even negative values of β result in higher likelihood than that attained by EM.

589 6.2. Multi-state bandit with hidden state

590 In contrast to the previous experiment a policy esti-
 591 mation is now introduced. The estimation of the per-
 592 ceptual state has to be performed on the basis of *in-*
 593 *direct* reward attribution, i.e., the state now becomes
 594 hidden.

595 6.2.1. Maximization of the likelihood—binary bandit

596 This section shows the results on problems in which
 597 the reward function is well aligned with the likeli-

hood, i.e., the problems where maximization of the 598
 reward results in maximization of the likelihood. Re- 599
 sults for this task are shown in Fig. 9. Unlike in 600
 the experiments of the previous section, the cluster 601
 identity is not important, as long as they correctly 602
 partition the input space. The multi-state bandit es- 603
 sentially implements the mapping from clusters to 604
 labels. 605

It is particularly interesting to see if the reward-based 606
 estimator of the input density results in a better fit 607
 of the resulting observation density to the one that 608

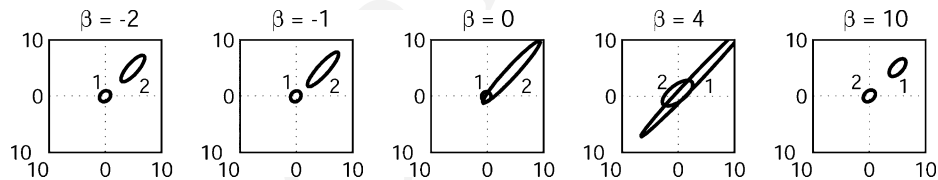


Fig. 8. Results of a run of the algorithm for different values of β starting from the same initial conditions. For coefficients with opposite signs the labeling is reversed, while the EM produces the labeling by chance. In this run the source distribution has the component 1 at the position (5, 5) and component 2 at (0, 0).

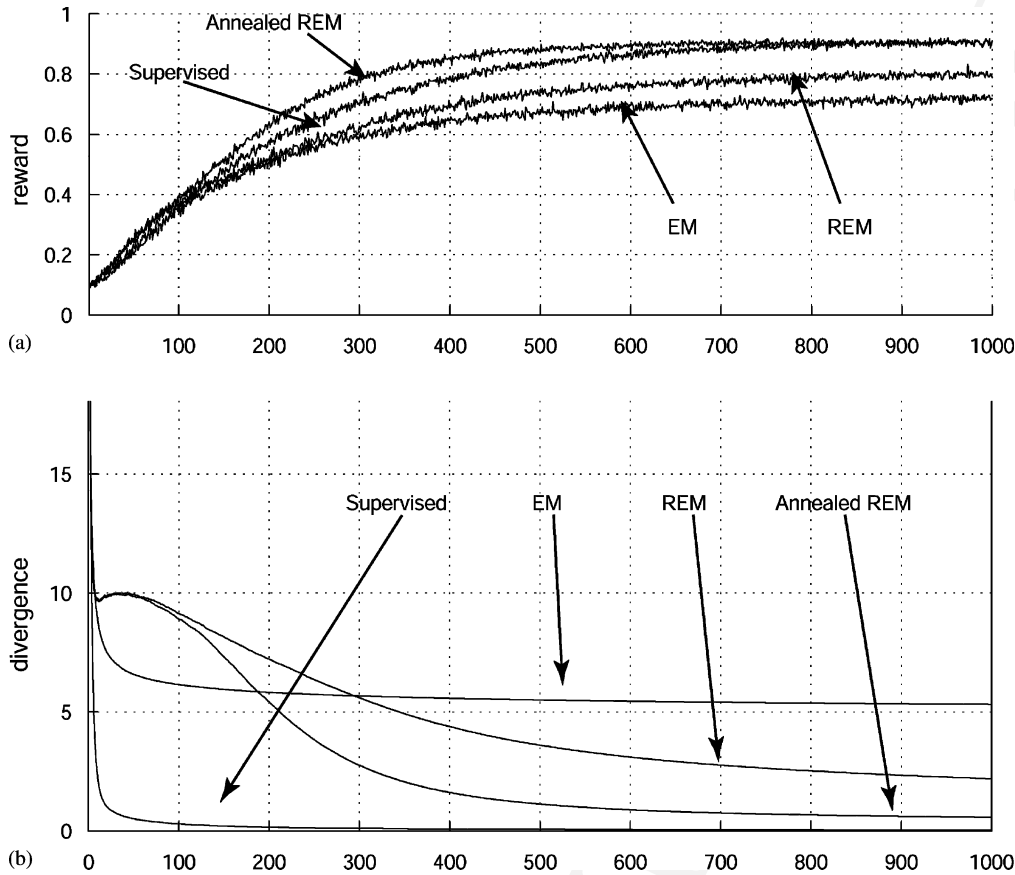


Fig. 9. (a) Performance on the 2-state 10-armed binary bandit. (b) Divergence between estimated and true source distributions.

609 gets reinforced than the regular EM. In the case
 610 of a Gaussian mixture density with a known num-
 611 ber of components (known number of states), the
 612 fit can be measured with the symmetric KL diver-
 613 gence

$$615 \quad S(p||q) = \frac{1}{4}[(\mu_p - \mu_q)^T(\Sigma_p^{-1} + \Sigma_q^{-1})(\mu_p - \mu_q) \\ 616 \quad + \text{tr}(\Sigma_p^{-1}\Sigma_q + \Sigma_q^{-1}\Sigma_p - 2\mathbf{I})]. \quad (14)$$

617 For a lack of a better analytical method, this quan-
 618 tity is computed for every combination of source
 619 and estimator components and the minimum value is
 620 selected.

621 The experiment with a 2-state 10-arm bandit is per-
 622 formed as follows:

-
- (1) *Initialize*: for each state, randomly select a Gaussian observation model: $\mu_i \sim N(\mathbf{0}, 2I)$, $\sigma_i = I$;
 - (2) *Iterate*:
 - (a) randomly choose a generator state, s_k ;
 - (b) generate an observation, x^n , from $\mathcal{N}(\mu_k, \sigma_k)$;
 - (c) using current parameters select an action a^n ;
 - (d) if a^n is the same as the optimal arm deliver a reward of 1, otherwise -1 ;
 - (e) update parameters of the model;
-

623
 624 One variation on the algorithm described in this
 625 paper is the REM with the parameter β changing
 626 over time. For example, slowly increasing β , start-
 627 ing with the value of 0 will cause the algorithm to

628 not pay any attention to the reward initially, while
629 slowly shifting towards the “winner-take-all” mode
630 after some period of time. Let us call it annealed
631 REM.

632 Fig. 9a shows the average amount of reward collected
633 by bandits trained with the EM, REM and
634 annealed REM algorithms compared to the case
635 where the input space is estimated via a supervised
636 estimator. As the goal is an accurate reproduction of
637 the source mixture, these plots need to be considered
638 along with the divergence plots Eq. (14), given in
639 Fig. 9b. The annealed REM algorithm, which slowly
640 increases the value of the parameter β , performs
641 very well, converging even faster than the supervised
642 case. It is somewhat puzzling, but easily explained
643 by the fact that the annealing amounts to simultane-
644 ous exploration of all states of the bandit in the
645 initial stages. This gives a good set of initial condi-

646 tions for subsequent search in each bandit when β
647 increases.

6.2.2. Maximization of the likelihood—full bandit

648 The algorithm works with the full bandit, where
649 each action taken by the algorithm results in some
650 value of the reward—positive or negative, with no
651 modifications. The results are shown in Fig. 10a. As in
652 the case with the binary bandit, the initial convergence
653 of both REM and annealed REM is faster than the
654 supervised case. The advantage, compared to EM, how-
655 ever, seems less spectacular than in the binary case.
656 The divergence plots (Fig. 10b), as before, show better
657 fit of REM and annealed REM to the source distribu-
658 tion.

659 This experiment shows the worst case scenario for
660 the algorithm. The reward structure here has many local
661 maxima and is “distracting” for the on-line search.
662

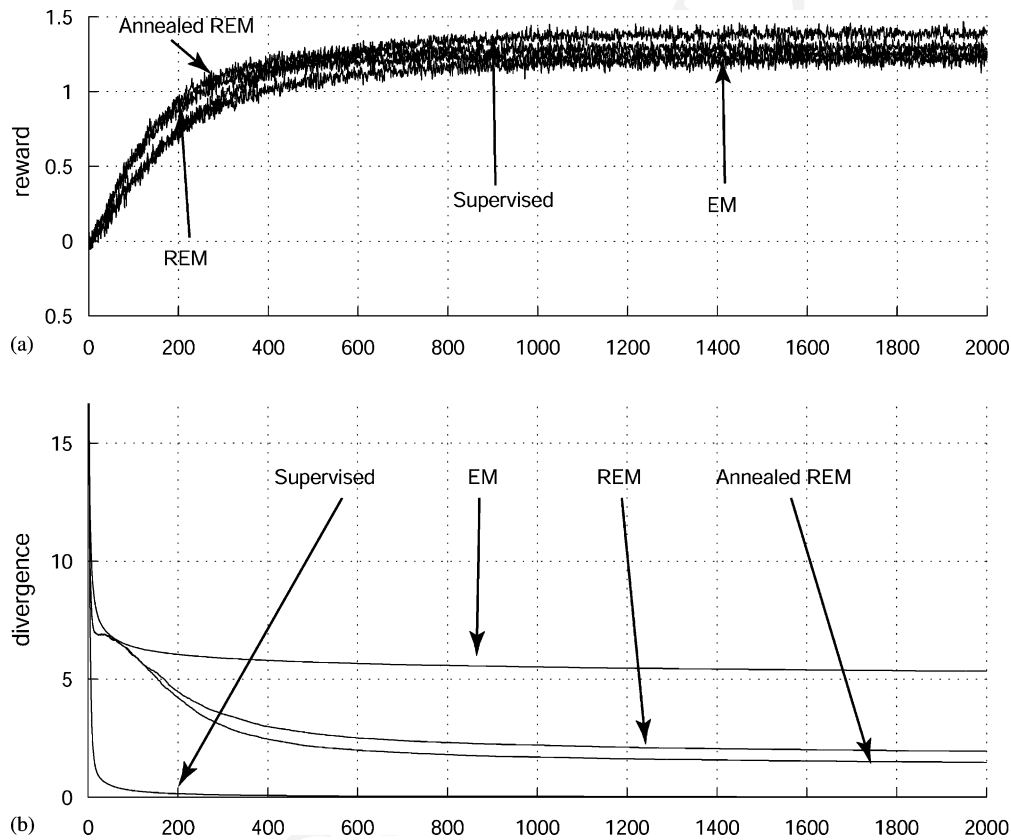


Fig. 10. (a) Performance on the full 2-state 10-armed bandit. (b) Divergence between estimated and true source distributions.

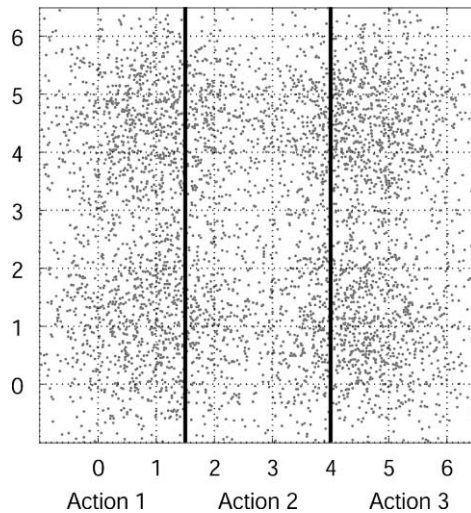


Fig. 11. Source and the reward structure for the reward bound maximization task. The data forms four strong clusters, while the reward is delivered for selecting action 1 if the data comes from the area marked “action 1”, etc.

663 The search becomes more difficult and the limitations
 664 of the search algorithm become the deciding factor
 665 in the achieved performance. However, despite the in-
 666 consistencies in the reward, the perceptual system cap-
 667 tures the input distribution better when aided by the
 668 reward than when no feedback is given.

669 6.2.3. Maximization of the reward

670 It is interesting to see how this model performs
 671 on a problem in which the reward function is not

aligned with the likelihood. The problem in this sec- 672
 tion is as follows—the input data is generated from 673
 four 2-dimensional Gaussians. However the reward is 674
 delivered in such a way that action a_1 is rewarded 675
 when $x_1^n < 1.5$, a_2 when $1.5 \leq x_1 < 4$ and a_3 when 676
 $x_1 > 4$, as shown in Fig. 11. 677

The performance of the model on this task is shown 678
 in Fig. 12. After 2000 iterations the EM estimator 679
 yields an average reward of 0.58, annealed REM— 680
 0.82 and supervised estimator—0.96 with the maxi- 681
 mum possible reward of 1. 682

Fig. 13 shows results of a single run of the algo- 683
 rithm. The left column of the figure shows the 684
 resulting positions and outlines of the mixture compo- 685
 nents. The middle column shows the classifica- 686
 tion decision regions corresponding to the cluster- 687
 ing shown on the left. The right column shows the 688
 “cluster assignment”—matrices that map states to 689
 actions, $p(a|s)$. A value in k th position of l th row 690
 of the matrix indicates the probability of selecting 691
 an action k once the point x^n is classified as be- 692
 longing to the cluster l . Figure (a)–(c) demonstrates 693
 the performance of the annealed REM algorithm, 694
 (d)–(f)—that of the supervised model, and the bot- 695
 tom row (g)–(i)—the performance of the unguided 696
 EM. The supervised case gives the best possible 697
 partitioning of the input while using three Gaus- 698
 sians (component 4 is never used and therefore has 699
 a mixing coefficient 0). The REM uses all four 700
 components and aligns them with the reward parti- 701
 tioning. Note that both clusters 2 and 4 select action 702
 a_1 . 703

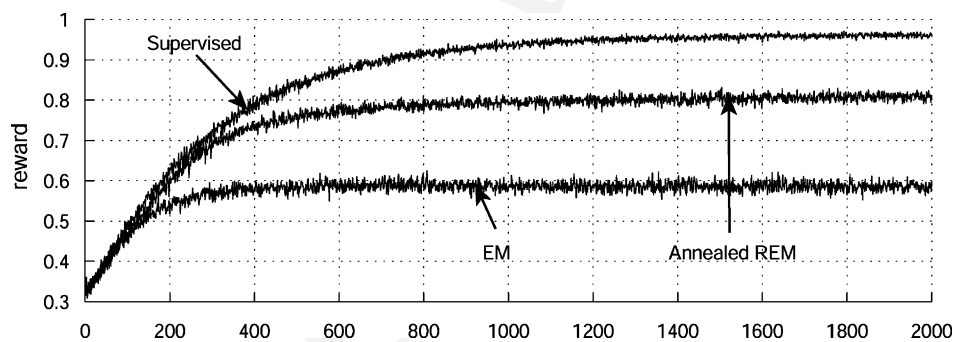


Fig. 12. Performance of EM, REM and a fully supervised estimator on the problem where reward structure does not coincide with the likelihood (averaged over 2000 runs).

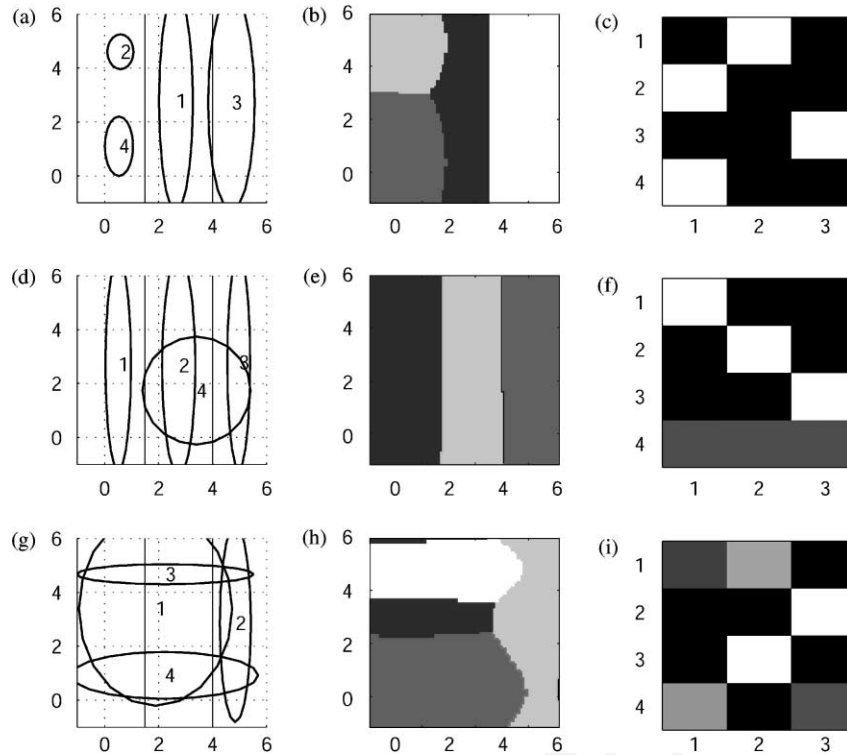


Fig. 13. Final cluster positions (left column), decision regions (middle column) and cluster assignment matrices (right column) for REM (top row), supervised (middle row) and EM (bottom row) estimators after a single run.

704 7. Conclusions

705 This paper presented an extension to the EM al-
 706 gorithm that allows for solving a range of learning
 707 tasks—from fully unsupervised, to fully supervised,
 708 including the partially and weakly labeled data. The
 709 justification for entropic variations of the posterior
 710 to achieve arbitrary component assignment goals is
 711 provided in the text. The algorithm allows for smooth
 712 blending between likelihood- and reward-based
 713 costs.

714 The paper shows that inclusion of the reward signal
 715 into the process of state estimation is important if we
 716 want to design agents without explicit programming
 717 of their perceptual states. The feedback is not only
 718 important for computing the policy of action selec-
 719 tion, but also as a guiding mechanism for developing
 720 a robust grounded perception. In contrast to unsuper-
 721 vised techniques, where the final cluster configuration
 722 is aligned with the likelihood of the data, in the algo-

rithm shown in this paper the grounding is achieved in
 a procedure allowing us to develop the configuration
 with a high utility to the agent.

723
 724
 725 One of the problems of the algorithm is the appro-
 726 priate choice of the parameter β . In some cases it is
 727 convenient to have an asymmetric schedule for posi-
 728 tive and negative rewards, which adds another param-
 729 eter to the set.

730
 731 In other cases special care must be taken about the
 732 fact that both reward signal for the clustering algorithm
 733 and the state assignment for the action selection are
 734 non-stationary.

Appendix A

735
 736 Theorem 1 shows that the a reward objective func-
 737 tion $F(x, \theta)$ (Eq. (9)) is a lower bound on a log like-
 738 lihood, $\log f(\theta)$ and can be used for EM-type estima-
 739 tion.

740 **Theorem 1.** $F(x, \theta)$ is a lower bound on $\log f(\theta)$.

743 **Proof.** Starting from (9), one can write

$$\begin{aligned}
 744 \quad F(x, \theta) &= -D(q(s) \| p(s|x, \theta)) \\
 &+ E_{q(s)} \left[\log \frac{p^*(s|x; r)}{p(s|x, \theta)} \right] + \log f(\theta) \\
 745 \quad &= \int q(s) \log \frac{p(s|x, \theta)}{q(s)} ds \\
 &+ \int q(s) \log \frac{p^*(s|x; r)}{p(s|x, \theta)} ds + \log f(\theta) \\
 746 \quad &= \int q(s) \left[\log \frac{p(s|x, \theta)}{q(s)} + \log \frac{p^*(s|x; r)}{p(s|x, \theta)} \right] ds \\
 &+ \log f(\theta) = \int q(s) \log \frac{p^*(s|x; r)}{q(s)} ds \\
 747 \quad &+ \log f(\theta) = -D(q(s) \| p^*(s|x; r)) \\
 &+ \log f(\theta). \tag{A.1}
 \end{aligned}$$

752 In the last line of Eq. (A.1) the divergence,
 753 $D(q(s) \| p^*(s|x; r)) \geq 0$ from which it follows that

$$754 \quad F(x, \theta) \leq \log f(\theta), \quad \forall q(s), \theta, \text{ s.t. } \sum (q(s)) = 1 \tag{A.2}$$

757 with equality holding iff $q(s) = p^*(s|x; r)$. This con-
 758 cludes the proof. \square

759

References

- | | |
|--|--------------------------|
| [1] N.M. Dempster, A.P. Laird, D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, <i>Journal of Royal Statistical Society B</i> 39 (1977) 185–197. | 761
762
763 |
| [2] A. Likas, Reinforcement learning approach to online clustering, <i>Neural Computation</i> 11 (1999) 1915–1932. | 764
765 |
| [3] T.P. Minka, Tutorial: expectation-maximization as lower bound maximization, Tutorial Note, 1999. | 766
767 |
| [4] R.M. Neal, G.E. Hinton, A view of the EM algorithm that justifies incremental, sparse and other variants, in: M.I. Jordan (Ed.), <i>Learning in Graphical Models</i> , MIT Press, Cambridge, MA 1998, pp. 355–368. | 768
769
770
771 |
| [5] K. Nigam, A. McCallum, S. Thrun, T. Mitchell, Text classification from labeled and unlabeled documents using EM, <i>Machine Learning</i> 39 (2/3) (2000) 103–134. | 772
773
774 |
| [6] R. Redner, H. Walker, Mixture densities, maximum likelihood, and the EM algorithm, <i>SIAM Review</i> 26 (2) (1984) 195–239. | 775
776 |
| [7] P.N. Sabes, M.I. Jordan, Reinforcement Learning by Probability Matching, Vol. 8, NIPS, 1996. | 777
778 |
| [8] R.S. Sutton, A.G. Barto, <i>Reinforcement Learning: An Introduction</i> , MIT Press, Cambridge, MA, 1998. | 779
780 |
| [9] M.A.L. Thathachar, P.S. Sastry, A new approach to the design of reinforcement schemes for learning automata, <i>IEEE Transactions on Systems, Man and Cybernetics</i> 15 (1985) 168–175. | 781
782
783
784 |
| [10] V.N. Vapnik, <i>Statistical Learning Theory</i> , Wiley, New York, 1998. | 785
786 |
| [11] R.J. Williams, A class of gradient-estimating algorithms for reinforcement learning in neural networks, in: <i>Proceedings of the IEEE First International Conference on Neural Networks</i> , San Diego, CA, 1987. | 787
788
789
790 |
| [12] L. Xu, M.I. Jordan, On convergence properties of the EM algorithm for Gaussian mixtures, <i>Neural Computation</i> 8 (1) (1996) 129–151. | 791
792
793 |