

Cross-Entropy method: convergence issues for extended implementation

Frédéric Dambreville

Délégation Générale pour l'Armement, DGA/CTA/DT/GIP

16 Bis, Avenue Prieur de la Côte d'Or

Arcueil, F 94114, France

Web: <http://www.FredericDambreville.com>

Email: <http://email.FredericDambreville.com>

February 2, 2008

Abstract

The cross-entropy method (CE) developed by R. Rubinstein is an elegant practical principle for simulating rare events. The method approximates the probability of the rare event by means of a family of probabilistic models. The method has been extended to optimization, by considering an optimal event as a rare event. CE works rather good when dealing with deterministic function optimization. Now, it appears that two conditions are needed for a good convergence of the method. First, it is necessary to have a family of models sufficiently flexible for discriminating the optimal events. Indirectly, it appears also that the function to be optimized should be deterministic. The purpose of this paper is to consider the case of partially discriminating model family, and of stochastic functions. It will be shown on simple examples that the CE could fail when relaxing these hypotheses. Alternative improvements of the CE method are investigated and compared on random examples in order to handle this issue.

1 Introduction

The Cross-Entropy method has been developed by R. Rubinstein for the simulation of rare events[1]. The algorithm iteratively builds a near-optimal importance sampling of the rare event, based on a family of parameterized sampling laws. The construction of the importance sampling is obtained by iteratively:

- tossing samples,
- selecting the samples which are approximating the rare events,
- relearning the parameters of the sampling law by minimizing its Kulback-Leiber distance (cross-entropy) with the selection,
- computing the importance weightings.

By considering the optimal events related to an objective as rare events, the method has been extended to optimization problems.

The cross-entropy method has been implemented successfully on many combinatorial problems. However, attempted proofs of the method make some assumptions as preliminary requests[1, 4]. First, the proof has been made in a deterministic context. Secondly, the closure of the simulation law family should contain the dirac on the optimum (or laws with support on the optimum).

The first condition cannot be fulfilled properly, in case of stochastic problem. The second condition is an obvious requirement. But there are some cases, where it is not possible to handle all the solutions precisely by the law family. Indeed, the solutions may not be countable practically; this is typically the case for some dynamic problems (for example, the strategy tree against a deterministic computer chess player). Both difficulties are encountered in optimal planning with partial observation. The purpose of this paper is to point out on simple examples, that these hypotheses are necessary for the convergence of the classical CE method. The questions are:

- *Does the classical CE algorithm solve stochastic problems properly?* It appears that the quantile selection within the CE may not work properly, without a rather good estimation of the objective functional expectation. Nevertheless, smoother selection criteria seem to be a possible answer to these difficulties.
- Assume that the law family closure does not contain all the deterministic solutions. The CE algorithm will converge to a stochastic approximation of the optimal solution. *Is this approximation the best possible within the law family?* Our answer to this question is not absolutely negative. But it appears that some extensions of the CE, quite usually implemented, will fail on this question.

This paper presents some counterexamples to these questions. In the case of stochastic optimization, tests are done on simple random examples in order to compare the convergence of various CE methods with the global optimum.

Next section introduces shortly the principle of the CE method. Section 3 will consider the case, where the optimal solution is not caught properly by the sampling family. A counterexample is proposed and studied. In section 4, stochastic problems are considered. Two simple counterexamples are investigated, thus enlightening some typical convergence difficulties. Different evolutions of the cross-entropy are then compared to the basical method, by generating several random examples. In particular, a method with smooth sample selection is proposed as a possible alternative for the stochastic problems. Section 5 concludes.

2 Basis of the cross-entropy method

The reader interested in CE methods should refer to the tutorial [2] and the book [1] on the CE method. CE algorithms were first dedicated to estimating the probability of rare events. A slight change of the basic algorithm made it also good for optimization. We will not focus on the cross-entropy method for simulation, although this primary aspect of the method is quite interesting. Rather, the CE method for optimization is now presented and discussed. While there are different evolutions of the primary method related to the choice of the selective rate or to a smooth update, this presentation is restricted to the basical CE method. By the way, it is not difficult to attest that the counterexamples proposed in sections 3 and 4 still work with these evolutions.

2.1 General CE algorithm for the optimization

The Cross Entropy algorithm repeats until convergence the three successive phases in order to maximize a given reward criterion:

1. Generate samples of random data according to a parameterized random mechanism,
2. Select the best samples according to the reward criterion,
3. Update the parameters of the random mechanism, on the basis of the selected samples.

In the particular case of CE, the update in phase 3 is obtained by minimizing the Kullback-Leibler distance, or cross entropy, between the updated random mechanism and the selected samples. The next paragraphs describe on a theoretical example how such method can be used in an optimization problem.

Formalism. Let be given a function $x \mapsto f(x)$; this function is easily computable. The value $f(x)$ has to be maximized, by optimizing the choice of $x \in X$. The function f will be the reward criterion.

Now let be given a family of probabilistic laws, $P_\sigma |_{\sigma \in \Sigma}$, applying on the variable x . The family P is the parameterized random mechanism.

Let $\rho \in]0, 1[$ be a selective rate. The CE algorithm for (x, f, P) follows the synopsis :

1. Initialize $\sigma \in \Sigma$,
2. Generate N samples x_n according to P_σ ,
3. Select the ρN best samples according to the reward criterion f ,
4. Update σ as a minimizer of the cross-entropy with the selected samples:

$$\sigma \in \arg \max_{\sigma \in \Sigma} \sum_{n \text{ selected}} \ln P_\sigma(x_n),$$

5. Repeat from step 2 until convergence.

This algorithm requires f to be easily computable and the sampling of P_σ to be fast.

Interpretation. The CE algorithm tightens the law P_σ around the maximizer of f . Then, when the probabilistic family P is well suited to the maximization of f , it becomes equivalent to find a maximizer for f or to optimize the parameter σ by means of the CE algorithm. The problem is to find a good family, and convergence parameters.

Extensions.

Smooth update. The method has been extended by implementing a smooth update of the law. More precisely, assume the set $\{P_\sigma | \sigma \in \Sigma\}$ to be convex, and let $\alpha \in [0, 1[$ be a smoothing rate. The algorithm follows the synopsis :

1. Initialize $\sigma \in \Sigma$,
2. Generate N samples x_n according to P_σ ,
3. Select the ρN best samples according to the reward criterion f ,
4. Define σ_1 as a minimizer of the cross-entropy with the selected samples:

$$\sigma_1 \in \arg \max_{\sigma_1 \in \Sigma} \sum_{n \text{ selected}} \ln P_{\sigma_1}(x_n),$$

5. Define σ_2 such that $P_{\sigma_2} = \alpha P_\sigma + (1 - \alpha) P_{\sigma_1}$, and update σ by setting $\sigma := \sigma_2$,
6. Repeat from step 2 until convergence.

Adaptive parameters. The principle is to make the parameters α and ρ dependent of the iteration time of the algorithm or on other contextual informations. Adaptive parameters appears as a main ingredient in the different proofs of convergence of the method.

Sampling with rejection. In some examples (particularly the *salesman*) considered in the CE tutorial [2], the laws family $P_\sigma|_{\sigma \in \Sigma}$ does not match the set X of valid values for the variable x . More precisely, there is a set $Y \supseteq X$ such that $P_\sigma \in \mathcal{P}(Y)$, *i.e.* P_σ is defined as a probability over Y . The implementation of such a law family in the CE methods is possible by rejecting the invalid samples generated by P_λ . A slight change is implied in the step 2 of the CE algorithm:

2. Repeat the subsequent process for any $n \in \{1, \dots, N\}$:
 - (a) Generate a sample $x \in Y$ according to P_λ ,
 - (b) If $x \notin X$, then repeat from step (a),
 - (c) *At this step*, $x \in X$. Then, set $x_n = x$.

There is *no other change* implied to the algorithm. In particular, the update step is the same: *the update of P_λ is done from the selected values of the subset X .*

At first sight, *this update of the law is questionable in regards to the rejection.* Indeed, the law to be learned from the samples is $P_\lambda / \sum_{x \in X} P_\lambda(x)$ and not P_λ . This induces a different result while minimizing the cross-entropy with the selected samples.

However, the rejection could also be derived from a parameter adaptation: the idea is to interpret the invalid samples of $Y \setminus X$ as samples with very bad reward. Then, the classical CE scheme is recovered by adapting the number of samples N and the parameter of selection ρ in order to reject these invalid samples.

This last interpretation makes sense, when the process actually converges to a law with a support included in X . This is the case, for example, when the law converges to a dirac around the optimum. But otherwise, it will be shown in section 3 that the convergence may be biased.

Convergence. Different convergence results have been proposed for the method and its extensions [4, 1, 3]. The convergence needs a proper tuning of the parameters of the algorithm (selecting rate, smoothing, number of samples). Essentially, these results have been established for the optimization of deterministic functions. Another issue is the stability of the optimization process, when the family of law, $P_\sigma|_{\sigma \in \Sigma}$, does not necessarily match the optimal value properly. The questions investigated by this paper are:

- *Does the classical CE algorithm solve stochastic problems properly?* A negative answer is given subsequently. An evolution of the CE is proposed in order to deal with this problem.
- Assume that the law family closure does not contain the dirac, or dirac mixture, around the optimal solutions. *Does the CE process provide the best approximation possible within the law family?* A partial negative answer is provided in next section, by producing a counterexample based on a sampling law with reject. This counterexample does not work in the classical scheme of the CE. It is not clearly answered in this paper, what should be the conditions in the CE process for guaranteeing such stability of the convergence. But it is sure that one have to be more careful in the choice and manipulation of the family.

3 When the family of laws does not enclose the optimum

The subsequent example is inspired from a convergence flaw diagnosed within a practical trajectory planning experiment; an experiment achieved by Francis Celeste [5], which is working in our team.

Problem setting. It is assumed that an agent has two possible actions: the action *continue* or the action *end*. Each time the agent decides to *continue*, it receives the reward +1 and the process is continued. When the agent decides to *end*, it still receives the reward +1 but the process is terminated. Thus, the agent has to choose a sequence of action, which is a repetition of the action *continue* terminated by the action *end*:

continue; continue . . . continue; end .

The reward for a whole sequence of action is t , the length of the sequence. Now, a constraint of length is imposed to the actions. The sequence of action cannot contain more than T actions, so that $t \leq T$.

Optimal solution. The optimal solution is obvious. The agent will do as many action as possible. Its optimal sequence of action is thus:

$\underbrace{\text{continue}; \dots; \text{continue}}_{T \times}; \text{end} .$

The problem is actually a triviality. But we will see that for some laws family, the CE with rejection will fail in finding the optimal law.

Proposal of a laws family, and convergence issue. On such a simple example, the best choice is perhaps a law on the length of the process sequence. But in fact, this kind of problem could be easily generalized so as to involve more than two possible actions (not only *continue* or *end*). Then, a Markov chain is generally used for these problems. In the salesman problem, for example, the actions are the choice for a town; the salesman is solved in [1, 2] by means of a Markov chain with reject. A method with reject is investigated subsequently.

The purpose is to sample a sequence $(d_\theta | 1 \leq \theta \leq t)$ where $1 \leq t \leq T$, $d_\theta = \text{continue}$ for $\theta < t$, and $d_t = \text{end}$. This sampling will be done by means of a reject method:

- Generate a sample without size constraint: $(d_\theta | 1 \leq \theta \leq t)$ where $1 \leq t$, $d_\theta = \text{continue}$ for $\theta < t$, and $d_t = \text{end}$,
- Reject the sample when $t > T$.

Sampling a sequence without size constraint. The sampling will be generated uniformly and independently for each step, so that the sampling law of the sequence is characterized by the law p_λ for sampling a single action:

$$p_\lambda(d_\theta = \text{continue}) = \lambda \quad \text{and} \quad p_\lambda(d_\theta = \text{end}) = 1 - \lambda .$$

The whole process takes into account the ending state, so that the sample generation follows the following synopsis:

1. Set $t = 0$,
2. Set $t := t + 1$,
3. Generate d_t by means of the law p_λ ,
4. Repeat from step 2, until $d_t = \text{end}$.

As a consequence, the probability of a full sequence $d = (d_\theta | 1 \leq \theta \leq t)$ is given by:

$$P_\lambda(d) = \lambda^{t-1}(1 - \lambda) .$$

Optimal law. The optimal law is the one which yields the best gain expectation for the valid trajectories generated by P_λ . The gain expectation after rejection is given by:

$$E_{P_\lambda(\cdot | t \leq T)} t = \frac{\sum_{t=1}^T t \lambda^{t-1} (1 - \lambda)}{\sum_{t=1}^T \lambda^{t-1} (1 - \lambda)} = \frac{\sum_{t=1}^T t \lambda^{t-1}}{\sum_{t=1}^T \lambda^{t-1}} .$$

This expectation is maximized when $\lambda = 1$:

Within the family, the optimal law is P_1 .

Notice that this optimal distribution is not an optimum for the problem. The family $P_\lambda | 0 \leq \lambda \leq 1$ is not sufficiently rich to handle the optimum of the function.

It is sometimes not possible to provide a family able to handle the global optimum of the function. Then, it is often sufficient to find the optimal distribution among the family. *Is the CE able to provide such optimal distribution among the family?* Subsequently, it is shown on the example that the CE (with rejection) does not converge to the optimal law P_1 .

Updating the law. Assume $M = \rho N$ samples ($d^n | 1 \leq n \leq M$) being obtained after a sampling process (with reject) and a selection of the best samples. Denote t_n the ending time of sequence d^n (*beware*: it is not a power operation!).

The parameter λ for the upcoming loop of the CE algorithm is obtained by maximizing the distance with the selected samples:

$$\lambda \in \arg \max \frac{1}{M} \sum_{n=1}^M \ln P_\lambda(d^n).$$

Now:

$$\frac{1}{M} \sum_{n=1}^M \ln P_\lambda(d^n) = \frac{1}{M} \sum_{n=1}^M \ln(\lambda^{t_n-1}(1-\lambda)) = \left(\left(\frac{1}{M} \sum_{n=1}^M t_n \right) - 1 \right) \ln \lambda + \ln(1-\lambda).$$

The maximization then results to the relation:

$$\left(\left(\frac{1}{M} \sum_{n=1}^M t_n \right) - 1 \right) \frac{1}{\lambda} - \frac{1}{1-\lambda} = 0.$$

At last, the following update relation is derived:

$$\lambda = 1 - M \left/ \sum_{n=1}^M t_n \right. . \quad (1)$$

Convergence issue. Equation 1 and the rejection constraint imply that $\lambda \leq 1 - \frac{1}{T}$ after update. As a consequence, *the CE does not converge to P_1 , the optimal distribution among the family.* In fact, it is even proved by considering the CE process that $\lambda < 1 - \frac{1}{T}$. Let P_{λ^*} be the law obtained after convergence of the CE. Then:

$$E_{P_{\lambda^*}(\cdot | t \leq T)} t < \frac{\sum_{t=1}^T t(1-1/T)^{t-1}}{\sum_{t=1}^T (1-1/T)^{t-1}}.$$

Let us consider the simple case $T = 2$, and compare the expectations:

$$E_{P_1(\cdot | t \leq T)} t = (1+2)/(1+1) = \frac{3}{2} \quad \text{and} \quad E_{P_{\lambda^*}(\cdot | t \leq T)} t < (1+2 \times \frac{1}{2})/(1+\frac{1}{2}) = \frac{4}{3}.$$

The difference, at least 11%, is not negligible.

Convergence in the CE classical scheme. As it has been discussed in section 2, the update of λ within the classical scheme will be obtained by minimizing the cross-entropy of the *conditional* law:

$$P_\lambda^*(d | t \leq T) = \frac{\lambda^{t-1}(1-\lambda)}{\sum_{\theta=1}^T \lambda^{\theta-1}(1-\lambda)} = \frac{\lambda^{t-1}(1-\lambda)}{1-\lambda^T}.$$

with the selected samples. Thus, the update is expressed by:

$$\lambda \in \arg \max \frac{1}{M} \sum_{n=1}^M \ln \frac{\lambda^{t_n-1}(1-\lambda)}{1-\lambda^T}.$$

Defining $\bar{t} = \frac{1}{M} \sum_{n=1}^M t_n$, the optimization reduces to:

$$\lambda \in \arg \max \frac{\lambda^{\bar{t}-1}(1-\lambda)}{1-\lambda^T}.$$

The maximum of this function is not necessarily located at $\lambda = 1$. For example, when $\bar{t} = 1$, the optimum is obtained for $\lambda = 0$. Now, the function to be optimized could be rewritten:

$$\frac{\lambda^{\bar{t}-1}(1-\lambda)}{1-\lambda^T} = \frac{\lambda^{\bar{t}-T}}{\sum_{k=0}^{T-1} \lambda^{-k}}.$$

Then, it is deduced:

$$\bar{t} \geq T \implies 1 \in \arg \max_{0 \leq \lambda \leq 1} \frac{\lambda^{\bar{t}-1}(1-\lambda)}{1-\lambda^T}. \quad (2)$$

The equation (2) has a clear interpretation: when $\lambda > 0$ at initialization and the selective rate ρ is sufficiently small, then the CE algorithm (without rejection) converge to the optimal law P_1 . As a conclusion, our counterexample fails in the classical CE paradigm.

Discussion. The previous example has shown convergence issue of the CE with reject when the laws family cannot reach the optimum of the function. This counterexample does not work when using a classical CE scheme. In general, even when the family cannot handle the optimum exactly, the convergence still works rather well in the classical CE paradigm. Many questions arise however. In particular, how to evaluate and enhance the stability of the convergence in regards to the discrimination weakness of the laws family?

4 When the problem is stochastic

In this section, it is discussed about the convergence of the CE in case of stochastic optimization. Notice that it is still possible to bring such stochastic problems back to deterministic problems by computing the expectation of the objective function. But generally, this computation is obtained by simulation and is costly. A reduction of the cost could be obtained by means of the method described in section 4.2.1.

When the variable to be optimized and the stochastic variable of the system are dependent, the expectation will make necessary the use of a functional abstraction of the variable to be optimized (instead of conditional laws). This is again somewhat costly. Moreover, the cost reduction method described in section 4.2.1 is no more feasible (when the variable of the system depends on the variable to be optimized).

The purpose of this section is to consider the stochastic optimization by the CE without computing the expectation of the objective. It is shown on simple examples that there may be a true convergence difficulty of the CE method in such conditions.

In the first subsequent example, the value to be optimized is conditioned by another variable which is stochastic. In other word, the value to be optimized could be considered as a function of the stochastic variable. Such problems do not appear classically in the CE literacy, but explain clearly some typical difficulties in the convergence. The second example is unconditioned and more classical. These examples will be completed by a study of stochastic optimization problems (here, without conditioning), which will be generated randomly. Alternative solutions to the classical CE are proposed and compared then.

4.1 Examples

4.1.1 Example 1

Typically, there is an additional difficulty in evaluating the expectation of the objective function, when the variable to be optimized are conditioned by the variable of the system. For this reason, we will start by considering this kind of example.

Let us consider the following stochastic problem:

$$f_o \in \arg \max_{f: x \mapsto d} \sum_x p(x) V(f(x), x),$$

where $x \in \{0, 1\}$, $d \in \{0, 1\}$, $p(0) = p(1) = \frac{1}{2}$, $V(d, x) = 2x + d$,
and f is a mapping from x to d .

(3)

This problem could be seen from a probabilistic viewpoint:

$$h_o \in \arg \max_h \sum_{d,x} p(x) h(d|x) V(d, x),$$

where $x \in \{0, 1\}$, $d \in \{0, 1\}$, $p(0) = p(1) = \frac{1}{2}$, $V(d, x) = 2x + d$,
and $h(d|x)$ is a probability of d conditionally to x .

(4)

We will apply a cross-entropic method in order to solve the optimization (4). Notice that the method will differ slightly from usually, since we are handling a conditional laws family.

Direct solve. The obvious answer to this problem is $h(0|x) = 0$ and $h(1|x) = 1$; the optimal gain is 2.

Cross-entropic solve. A cross-entropic procedure is proposed here with quantile selection $\rho = 10\%$ (no smooth update, for simplicity) in order to solve (4):

- Initialize h by $h(0|0) = h(1|0) = h(0|1) = h(1|1) = \frac{1}{2}$,
- Make 100 samples and evaluate them by V ,
- [*] Select the 10% best samples, update $h(\cdot|x)$ from the selected samples, when it is possible.¹ Reiterate from previous step.

Since $V(d_1, 1) > V(d_2, 0)$ for any choice of d_i , it comes that samples $(d, 0)$ are (almost) never² selected. As a consequence, $h(\cdot|0)$ is (almost) never updated and stays unchanged. Thus, a practical convergence will stale to the solution $h(0|0) = h(1|0) = \frac{1}{2}$; $h(0|1) = 0$; $h(1|1) = 1$, which is sub-optimal. The expected gain is then $\frac{7}{4}$.

This example contains a specific difficulty: we are indeed optimizing the function $x \mapsto f(x)$ by mean of a conditional law. By the way, one may argue that [*] is not a good updating strategy, since the samples should be selected relatively to each condition x . But this is not possible, when there are many possible conditions x (this is often the case).

4.1.2 Example 2

It could be argued about the previous example that the use of a conditional family is not the classical scheme for applying the CE method. This forthcoming example will be related to a more classical scheme.

¹Leave $h(\cdot|x)$ unchanged when there are no selected samples conditioned by x .

²Probability is around 10^{-18}

Now, let us solve the following stochastic problem:

$$d_o \in \arg \max_d \sum_x p(x) V(d, x),$$

$$\text{where } x \in \{0, 1\}, d \in \{0, 1\}, p(0) = p(1) = \frac{1}{2},$$

$$\text{and } V(0, 0) = 2, V(0, 1) = -2, V(1, 0) = V(1, 1) = 1.$$
(5)

From a CE viewpoint, the problem becomes:

$$h_o \in \arg \max_h \sum_x p(x) h(d) V(d, x),$$

$$\text{where } x \in \{0, 1\}, d \in \{0, 1\}, p(0) = p(1) = \frac{1}{2},$$

$$V(0, 0) = 2, V(0, 1) = -2, V(1, 0) = V(1, 1) = 1.$$

$$\text{and } h(d) \text{ is a probability of } d.$$
(6)

Direct solve. The optimal solution of (6) is of course $h_o(0) = 0$ and $h_o(1) = 1$, resulting in the gain 1.

Cross-entropic solve. A cross-entropic procedure is proposed here with quantile selection $\rho = 10\%$ (no smooth update, for simplicity) in order to solve (6):

- Initialize h by $h(0) = h(1) = \frac{1}{2}$,
- Make 100 samples and evaluate them by V ,
- Select the 10% best samples, update h from the selected samples. Reiterate from previous step.

Since $V(0, 0) > V(d, x)$ for any $(d, x) \neq (0, 0)$, it comes that the samples $(d, x) \neq (0, 0)$ are (almost) never selected. As a consequence, the selected samples will be $(0, 0)$ from the beginning of the CE process. Consequently, the CE process will converge to the sub-optimal solution $h_*(0) = 1$ and $h_*(1) = 0$, thus resulting in the gain 0.

4.1.3 Discussion.

The two previous examples are enlightening. It appears clearly that the selection scheme of the CE (selection of a quantile) does not work properly, in regards to a stochastic objective. Indeed, some configurations of the problem, which are sampled by the law of the system but not by us, will be automatically discarded by the quantile selection process. By discarding these cases, a convergence bias is generated.

4.2 Alternative methods

4.2.1 Computing the expectation (reduced cost)

This method is not exactly an alternative: it is costly. But it will be provided as a reference for the test comparison. The idea is to replace the stochastic objective function $V(d, x)$ by an estimation of its expectation. This expectation is obtained by sampling over x according to the law p of the system. More samples are used, more accurate is the estimation. Here, we are using the same samples of x for computing the expected gain of the samples d_n . This will reduce greatly the complexity. But such method is not feasible, when the variables x and d are dependent. The whole algorithm is explained subsequently:

1. Initialize h ,

2. Generate N samples d_n according to h ,
3. Generate K samples x_k according to p ,
4. Evaluate each sample d_n by the estimated expectation $v_n = \sum_{k=1}^K V(d_n, x_k)$,
5. Select the ρN best samples d_n according to the expectation v_n ,
6. Update h as a minimizer of the cross-entropy with the selected samples:

$$h \in \arg \max \sum_{n \text{ selected}} \ln h(d_n),$$

7. Repeat from step 2 until convergence.

4.2.2 Using another selection scheme for the CE

The idea here is to change the selection scheme of the CE. The stochastic objective function $V(d, x)$ is directly used here. As in section 4.1.2, the stochastic pair (d, x) is sampled and evaluated at the same time.

Selection scheme. Assume N samples (d_n, x_n) being evaluated by $v_n = V(d_n, x_n)$. It is defined a non decreasing function R , which will characterize the importance $R(v_n)$ of each sample (d_n, x_n) . The update of h will be computed as a maximizer of the cross entropy with the discrete weighted distribution $\left(d_n, \frac{R(v_n)}{\sum_{n=1}^N R(v_n)}\right)$.

Algorithm. The whole algorithm is explained subsequently:

1. Initialize h ,
2. Generate N samples d_n according to h and N samples x_n according to p ,
3. Evaluate each sample pair (d_n, x_n) by $v_n = V(d_n, x_n)$,
4. Update h as a minimizer of the cross-entropy with the weighted samples:

$$h \in \arg \max \sum_{n=1}^N R(v_n) \ln h(d_n),$$

5. Repeat from step 2 until convergence.

This selection scheme is called *smooth selection scheme*. Notice that the quantile selection of Rubinstein is a particular case of the *smooth selection scheme*, where the function R is a heavyside function pointed on the quantile.

4.3 Method comparison by means of Randomly generated tests

The three methods, basic CE; CE with expectation computation; and smooth selection scheme, have been compared on random problems. The method for creating the problems is simple:

- There are 100 possibles states for d and for x , that is $d, x \in \{1, \dots, 100\}$,
- The parameters $V(d, x) \in]0, 1]$ are generated randomly, according to the uniform law, for any d and any x ,
- The probability p is generated randomly, according to the uniform law (that is the 99-dimensions vector characterizing p is generated uniformly),

Notice that it is quite easy to solve these problems, by enumerating the cases.

The test has been executed 1000 times. The parameters of the algorithm are:

- $K = N = 100$ and $\rho = 10\%$,
- The update is smoothed by $\alpha = 0.9$ (*i.e.* the innovation is 10%),
- The importance function R is defined by $R(v_n) = v_n$.

The following table gives the percentage of the optimum achieved by each method. These results are averaged over the 1000 executed tests, and the variance is given.

Optimal percentage	Mean	Variance
Basic CE	93.9%	3.7%
Expectation	99%	0.4%
Smooth scheme	99.1%	0.7%

The convergence speed of the expectation CE and the smooth selection CE was comparable. Since the expectation is computed with reduced cost, the methods run with similar computation cost.

5 Conclusion

This paper has investigated the convergence issues of the cross-entropy method when relaxing the constraints of use. A counterexample has been found for the CE with reject, when the laws family used for the CE is too weak and does not contain the optimum dirac. Counterexamples have been found when optimizing a stochastic objective function. Weakness of the family and stochastic objective are very important context of use of the CE algorithm. By the way, both difficulties are encountered when optimizing a control with partial observation[6]. An alternative evolution of the CE has been proposed for the stochastic optimization. It is based on a smooth scheme for the sample selection. The convergence of weak laws family is still an unsolved question. Next works will focus on this difficult problem. Moreover, the proof of convergence of the smooth selection scheme will be investigated; at this time, this method has been evaluated only by experimental means.

References

- [1] R. Rubinstein, D. P. Kroese, *The Cross-Entropy method. An unified approach to Combinatorial Optimization, Monte-Carlo Simulation, and Machine Learning*, Information Science & Statistics, Springer 2004.
- [2] De Boer and Kroese and Mannor and Rubinstein, *A Tutorial on the Cross-Entropy Method*,
<http://www.cs.utwente.nl/~ptdeboer/ce/>
- [3] R. Margolin, *On a Convergence of the Cross-Entropy Method*, Annals of Operations Research, Springer Netherlands, April 2005.
- [4] Homem-de-Mello, Rubinstein, *Rare Event Estimation for Static Models via Cross-Entropy and Importance Sampling*,
<http://users.iems.nwu.edu/~tito/list.htm>
- [5] F. Celeste, F. Dambreville and J.-P. Le Cadre, *Optimal path planning using cross-entropy method*, Conference Fusion 2006, Florence, Italy, July 2006.
- [6] F. Dambreville, *Cross-entropic learning of a machine for the decision in a partially observable universe*, Journal of Global Optimization, Springer Netherland, August 2006 (on line).