

Taylor-based pseudo-metrics for random process fitting in dynamic programming.

Sylvain Gelly, Jérémie Mary, Olivier Teytaud

► **To cite this version:**

Sylvain Gelly, Jérémie Mary, Olivier Teytaud. Taylor-based pseudo-metrics for random process fitting in dynamic programming.. PDMIA, 2005, Lille, 16 p. inria-00000217

HAL Id: inria-00000217

<https://hal.inria.fr/inria-00000217>

Submitted on 13 Sep 2005

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Taylor-based pseudo-metrics for random process fitting in dynamic programming : expected loss minimization and risk management

Sylvain Gelly*, Jérémie Mary*, Olivier Teytaud* **

*TAO-Inria, LRI, UMR 8623(CNRS - Université Paris-Sud),

bat 490 Université Paris-Sud 91405 Orsay Cedex France

email: jeremie.mary@lri.fr

**Artelys, 12 rue du 4 Septembre, 75002 Paris, www.artelys.com

July 5, 2005

Abstract

Stochastic optimization is the research of x optimizing $E C(x, A)$, the expectation of $C(x, A)$, where A is a random variable. Typically $C(x, a)$ is the cost related to a strategy x which faces the realization a of the random process.

Many stochastic optimization problems deal with multiple time steps, leading to computationally difficult problems ; efficient solutions exist, for example through Bellman's optimality principle, but only provided that the random process is represented by a well structured process, typically an inhomogeneous Markovian process (hopefully with a finite number of states) or a scenario tree. The problem is that in the general case, A is far from being Markovian. So, we look for A' , "looking like A ", but belonging to a given family \mathcal{A}' which does not at all contain A . The problem is the numerical evaluation of " A' looks like A ".

A classical method is the use of the Kantorovitch-Rubinstein distance or other transportation metrics [Pflug, 2001], justified by straightforward bounds on the deviation $|EC(x, A) - EC(x, A')|$ through the use of the Kantorovitch-Rubinstein distance and uniform lipschitz conditions. These approaches might be better than the use of high-level statistics [Keefer, 1994]. We propose other (pseudo-)distances, based upon refined inequalities, guaranteeing a good choice of A' . Moreover, as in many cases, we indeed prefer the optimization with risk management, e.g. optimization of $EC(x, noise(A))$ where $noise(.)$ is a random noise modeling the lack of knowledge on the precise random variables, we propose distances which can deal with a user-defined noise. Tests on artificial data sets with realistic loss functions show the relevance of the method.

1 Introduction, notations

1.1 Informal overview

Consider a deterministic system in discrete time that depends on an exogenous random process and on your decision. The question is : what is the decision you should take ?

If the dynamics of the system are known, if the random process has a structured representation, if the dimension of the state space is not too large, a very classical solution is dynamic programming. Unfortunately, in real life, the random process can often only be represented in a Markovian manner by a huge random process, the random process made of the few relevant variables being far from Markovian. We therefore have to design a Markovian random process, without explosion of the state space, that is close to the real one, which is only known through a use Markov random process or through a non-Markovian random process.

A classical solution is the replacement of the non-Markovian random process by a Markovian one, minimizing the distance between the original random process and the created one. In this paper :

- we recall that the Kantorovitch-Rubinstein distance corresponds to the minimization of a 0-degree Taylor's expansion of the objective function ;
- we show how to use a degree 1 or degree 2 Taylor's expansion ;
- we show how to extend this to optimization robust to distribution shifts of the random process;
- we show experimental results illustrating the improved behavior, in particular with the degree 2.

1.2 The problem

Let S be a set of strategies. Let A be a random variable belonging to a given set \mathcal{A} of random variables, with domain X . Let $s, a \mapsto C(s, a)$ be a cost function (where $s \in S$ is a strategy and a is a realization of A). Our goal is to find a strategy s^* which ensure a low cost (i.e. the smallest $E(C(s^*, A))$ possible).

We suppose we already have an optimization method Opt which can find this optimum for a subset $\mathcal{A}' \subset \mathcal{A}$. So if our random variable A is in \mathcal{A}' , there is no problem. The trouble comes when we have a random variable $A \notin \mathcal{A}'$.

We look for $s \in S$ such that $EC(s, A)$, cost expectation, is minimal, and the trouble is that $A \notin \mathcal{A}'$.

In all the paper, s^* is an optimal strategy for A . s^* is unknown. We assume that we can use s_0 , a strategy close to s^* (which will be used for the evaluation of derivatives in Taylor's expansions). In usual cases, such a s_0 exist ; indeed in many practical cases, different versions coming from experts or from simplified models are used for comparison.

So we'll look for $A' \in \mathcal{A}'$, such that $E_A C(opt(A'), A)$ is as small as possible, i.e. we'll use as strategy for A the strategy which is only optimal for A' , expecting that

this strategy is nearly optimal for A as well. So the relevance of the distance used to choose A' close to A is very important. The Kantorovitch-Rubinstein one has been proposed ([Gröwe-Kuska et al, 2003]) but it does not take into account the form of the cost function of the problem. We propose a way to take advantage of our knowledge about the cost function to make a better choice for A' .

In particular, we'll show that in order to achieve this goal, A' can consistently be chosen through $A' \in \operatorname{argmin}_{A' \in \mathcal{A}'} \operatorname{distance}(A, A')$, for a well-chosen distance. This distance can be computed from usually available elements : the derivatives of the cost function $C(\cdot)$ at a reasonably good strategy s_0 , high-level statistics or simulations based upon $A' \in \mathcal{A}'$ and $A \in \mathcal{A}$. Figure 1 illustrates the method.

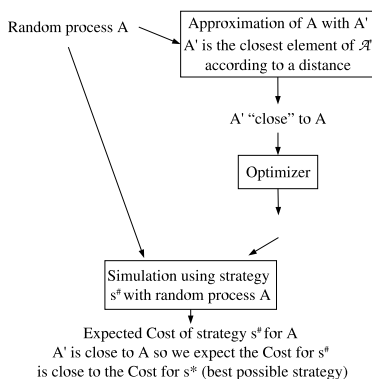


Figure 1: Typical paradigm for solving the problem. We study the upper-right box.

In all the paper, $d(A_1, A_2)$ where A_1 and A_2 are random variables, designs some distances between random variables, specified in the text. $d(s_1, s_2)$, where s_1 and s_2 are strategies, is a distance between the strategies such that the cost, for a given realization of the random process, is Lipschitzian in this distance (we use the norm of the gradient as the Lipschitz coefficient in the sequel, but a Lipschitz coefficient could be used if the gradient does not exist).

1.3 Motivations

Optimization in unknown environments and especially multi-stage stochastic optimization is a classical problem in discrete time optimal control. Many problems of this form, especially industrial problems, can be formulated in discrete time.

Such problems can be solved with different paradigms. Reinforcement learning, and its multiple variants ([Tesauro, 1989][Tesauro, 1999][Watkins, 1989][Neuneier, 1996][Neuneier, 1998][Neuneier et al, 1999]), is able to deal with very unstructured random processes, in particular without requiring that the model is Markovian ([Bertsekas et al, 1996][Sutton et al, 1998]). It allows to deal with risk through variance-penalized or more complex cost functions

([Aïd et al, 2003][Coraluppi et al, 1999][Geibel, 2001][Heger, 1994]), multi-actors environments ([Tesauro, 1999][Nagayuki et al, 2000]). It has been used in the domain of energy, in particular for its ability to deal with particular form of dynamic optimization such as multi-actor optimization or optimization with particular treatment of random processes ([Riedmiller et al, 2001][Aïd et al, 2003]).

These approaches, however, are notably slow in face of problems where the natural convexity (as in stock management) allows the use of efficient approaches with stable convergences : stochastic dynamic programming ([Bertsekas, 1995]). Whenever it is someone said that dynamic programming is a particular case of reinforcement learning, it is in fact a very particular case (it converges by Lagrangian decomposition and not by fix-point iterations), and the use of decomposition methods makes possible to treat real problems with very large size ; the convergence is proved and the precision is very good. However, these methods are very constrained from the point of view of the modelization of random processes ; only very structured random processes can be taken into account. The trouble is that in many cases, the random processes can only be known empirically through empirical time series and tendencies.

So the sum of the two constraints "reinforcement learning is too unstable and slow for the problem to be treated" and "stochastic dynamic programming can not deal with our random processes" leads to the following problem : how to choose a given family of random processes so as to fit a given random process ? Moreover, how to take risk into account ?

Assume that a criterion for choosing in a family of random process is provided. Then, the algorithmic resolution has been successfully performed in different papers ; in particular, [Hoyland et al, 2001] show that a complex criterion (therein based upon the respect of some high-level statistics) can be performed efficiently through non-linear programming, whenever the problem is non-convex ; it is easily differentiable and non-convexities can be handled thanks to multiple random initializations. Thus, a wide family of criterions can be concretely optimized.

Thus, the optimization of the criterion is tractable and the dynamic optimization thereafter is optimal within a very good precision after a long but controlled dynamic programming procedure. The choice of the criterion for choosing a model is then the main problem. In some cases the criterion was the preservation of some statistics, and in some cases the minimum Kantorovitch-Rubinstein distance. The latter is derived using bounds on the loss of precision depending upon the Kantorovitch-Rubinstein distance between the real random-process and its model. However, the bound is far from being optimal as it uses uniform bounds on the derivative. The purpose of this paper is the use of refined inequalities, using information about the derivatives, and to derive from these inequalities some improved distances for the choice of a model of the random process. Following [Gröwe-Kuska et al, 2003], we try to define the sense to be given to the term "look like", so that when a random process "looks like" the real random process, the optimum for this approximate random process is nearly optimal for the real random process. In particular, we search a distance more relevant than the Kantorovitch-Rubinstein distance. Moreover, we adapt our results to the case with risk.

The paper is organized as follows. Section 2 provides theoretical foundations (link with distances in 2.1, case of expectation 2.2, case of risk management 2.3). Section 3 provides experiments on artificial (reality-inspired) data sets. Section 4 concludes.

2 Theoretical results

Section 2.1 presents the general ideas. Section 2.2 presents and justifies a few distances for the optimization in expectation. Section 2.3 adapts these distances to the case of robust-to-noise optimization.

2.1 Stochastic optimization and distances

We want to optimize $s \mapsto EC(s, A)$. We look for $A' \in \mathcal{A}'$ such that $E_A C(Opt(A'), A)$ is as small as possible. We note that :

Lemma 1 :

$$EC(Opt(A'), A) \leq \inf_s EC(s, A) + 2 \sup_s |EC(s, A) - EC(s, A')|$$

Interpretation : This simply shows that if $EC(s, A)$ is uniformly close to $EC(s, A')$, then the cost in generalization, for the strategy optimized against A' instead of A , is surely close to the optimal strategy for A .

Proof :

Let's note $s' \in Opt(A')$, $s^* \in Opt(A)$, and $\epsilon = \sup_s |EC(s, A) - EC(s, A')|$.

$$\begin{aligned} EC(s', A) &\leq EC(s', A') + \epsilon \\ &\leq EC(s^*, A') + \epsilon \leq EC(s^*, A) + 2\epsilon \end{aligned}$$

■

2.2 Optimization of the expectation

We will note $A' = \pi(A)$. This (slight abuse of notation) notes that A' and $\pi(A)$ have the same law. As noticed in lemma 1, we are interested in $\sup_s |EC(s, A') - EC(s, A)|$; i.e. we look for a distance which is very related to this quantity (so that if we replace A by A' close to A for this distance, then by lemma 1 the optimal strategy for A' will be nearly optimal for A). We consider

$$\begin{aligned} \delta_s(A, A') &= |EC(s, A') - EC(s, A)| \\ \delta_s(A, A') &= \left| E_A \left(\nabla_A C(s, A)(\pi(A) - A) \right. \right. \\ &\quad \left. \left. + \epsilon \left\| \frac{\partial^2 C}{\partial A^2} \right\|_{S, [A, \pi(A)]}^\infty (A - \pi(A))^2 / 2 \right) \right| \end{aligned}$$

where $\epsilon \in [-1, 1]$. We will now assume that $d(s, s^*) \leq \eta$ and $d(s_0, s^*) \leq \eta$. Then

$$\begin{aligned} \delta_s(A, A') &= \left| E_A \left(\nabla_A C(s_0, A)(\pi(A) - A) \right. \right. \\ &\quad \left. \left. + \epsilon \left\| \frac{\partial^2 C}{\partial A^2} \right\|_{S, [A, \pi(A)]}^\infty (A - \pi(A))^2 / 2 \right) \right| \end{aligned}$$

$$+2\epsilon'\eta \left\| \frac{\partial C}{\partial S} \right\|_{\mathcal{S},[A,A']}^{\infty} \right)$$

where $\epsilon \in [-1, 1]$ and $\epsilon' \in [0, 1]$. So

$$\begin{aligned} \delta_s(A, A') \leq & \text{distance}(\pi, A, s_0) + \\ & \left| E_A \left(\epsilon \left\| \frac{\partial^2 C}{\partial A^2} \right\|_{\mathcal{S},[A,\pi(A)]}^{\infty} (A - \pi(A))^2 / 2 \right. \right. \\ & \left. \left. + 2\epsilon'\eta \left\| \frac{\partial C}{\partial S} \right\|_{\mathcal{S},[A,A']}^{\infty} \right) \right| \end{aligned}$$

with $\text{distance}(\pi, A, s_0) = |\text{deviation}(\pi, A, s_0)|$ and $\text{deviation}(\pi, A, s_0) = E_A \nabla_A C(s_0, A)(\pi(A) - A)$. This "distance" (which is not formally a distance), contains negative terms, and can be close to 0 whenever some terms are very large. We thus have another version which looks more safe : $\text{distance}'(\pi, A, s_0) = E_A |\nabla_A C(s_0, A)(\pi(A) - A)|$. However, we see below that this is not necessary.

We can go further with a Taylor expansion of order 2, the second derivative being indeed classically the criterion showing the transition to the risky domain :

$$\begin{aligned} \delta_s(A, A') \leq & \text{distance}^{(2)}(\pi, A, s_0) + \\ & \left| E_A \left(\epsilon \left\| \frac{\partial^3 C}{\partial A^3} \right\|_{\mathcal{S},[A,\pi(A)]}^{\infty} (A - \pi(A))^3 / 6 \right. \right. \\ & \left. \left. + 2\epsilon'\eta \left\| \frac{\partial C}{\partial S} \right\|_{\mathcal{S},[A,A']}^{\infty} \right) \right| \end{aligned}$$

with $\text{distance}^{(2)}(\pi, A, s_0) = |\text{deviation}^{(2)}(\pi, A, s_0)|$ with $\text{deviation}^{(2)}(\pi, A, s_0) = E_A \nabla_A C(s_0, A)(\pi(A) - A) + \frac{1}{2}(\pi(A) - A)^t H_A C(s_0, A)(\pi(A) - A)$ with $H_A C(s_0, A)$ the Hessian of $C(s_0, \cdot)$. As previously, one can prefer $\text{distance}'^{(2)}(\pi, A, s_0) = E_A (|\nabla_A C(s_0, A)(\pi(A) - A)| + |\frac{1}{2}(\pi(A) - A)^t H_A C(s_0, A)(\pi(A) - A)|)$.

We will consider as well a distance $\text{distance}^{(3)}$ in our experiments, Taylor expansion to order 3. The third derivative is not always very stable ; we prefer a term modelling the size of the error committed by the degree 2 approximation than a refined expansion. Thus, we replace the term of degree 3 by its absolute value.

Conclusion :

$$\begin{aligned} EC(\text{Opt}(A'), A) \leq & EC(s^*, A) + 2\text{distance}(\pi, A, s_0) \\ & + \text{error term} \end{aligned}$$

(error term of degree 1)

$$\begin{aligned} EC(\text{Opt}(A'), A) \leq & EC(s^*, A) + 2\text{distance}^{(2)}(\pi, A, s_0) \\ & + \text{error term}_2 \end{aligned}$$

(error term of degree 2) (where the error terms are upper bounded as explained above), i.e. we can optimize the choice of A' with the help of derivatives at s_0 .

2.3 Optimization with risk management

We will now consider an optimization in the sense of the expectation for risk ([Aïd et al, 2003][Coraluppi et al, 1999][Geibel, 2001][Heger, 1994]), here for a modified random process. We consider that the real random process is not $A \in \mathcal{A}$, but an unknown $\pi'(A)$, where π' is random. Instead of optimizing $EC(s, A)$, we thus want to optimize $EC(s, \pi'(A))$, where the expectation is with respect to both A and π' . Unfortunately, the straightforward application of the method above uses $\nabla_{\pi'(A)} C(s_0, \pi'(A))$ and we can only evaluate efficiently $\nabla_A C(s_0, \pi'(A))$. We need a distance between elements in \mathcal{A} and elements in \mathcal{A}' , taking into account high-level informations about the noise.

Precisely, we define $\pi'(a) = a + b$, where b is an additive random noise (π' is thus a random function). We will then optimize $EC(s, \pi'(A))$ instead of $EC(s, A)$.

We will so look for a distance between A' and A , that can be evaluated from available elements, such that $EC(Opt(A'), \pi'(A))$ is as small as possible. As previously, we will use lemma 1 and so try to optimize $sup_s |EC(s, A') - EC(s, \pi'(A))|$. We recall that we are interested in distances using only derivatives of C and distances between A and $\pi(A)$ for $A \in \mathcal{A}$ and π such that the law of $\pi(A)$ is in \mathcal{A}' .

So we have, with $\pi(A)$ of the same law as A' if A' represents A in \mathcal{A}' :

$$\begin{aligned} & EC(s, \pi'(A)) - EC(s, \pi(A)) \\ &= EC(s, \pi'(A)) - EC(s, A) + EC(s, A) - EC(s, \pi(A)) \\ &= deviation^{(2)}(\pi', A, s_0) - deviation(\pi, A, s_0) \\ &\quad + error\ term \end{aligned}$$

If b is invariant by rotation, the term of order 1 in $deviation^{(2)}$ disappears by symmetry and the remainder is

$$= \frac{1}{2} E Tr(H, A)K - deviation^{(2)}(\pi, A, s_0)$$

where K (chosen by the user) is $L_2(b)^2$ and $Tr(H, A)$ is the trace of the Hessian matrix $H_A C(s_0, A)$. So, the distance proposed for the case of risk is :

$$\begin{aligned} & distance_K^{(r)}(\pi, A, s_0) = \\ & |deviation^{(2)}(\pi, A, s_0) - \frac{1}{2} E Tr(H, A)K| \end{aligned}$$

It is not surprising that as K increases $\pi(A)$ gets further from A . Note that other forms of noise could be considered, for example with K depending upon A .

3 Experiments

We explain in 3.1 the overall paradigm which is tested here. The chosen problem is specified in 3.2. The pseudo-code of the experiments is provided in 3.3. The results are presented in section 3.4.

3.1 Proposed approach

Our goal is the comparison of different distances between random processes A and A' for the method illustrated by the figure 1. However, the distances we propose are based on a Taylor polynomial approximation and so we restrict A' to be "not too far" from A , so that the approximation is valid. The protocol is then the one illustrated in figure 1, but A' is chosen through optimization of $d(A, A')$ under the constraint that $dK(A, A') \leq dK_0$, where dK is the Kantorovitch-Rubinstein distance.

3.2 Problem specifications

Let $S = X = [0, 1]^d$. Consider a cost function $C(., .)$ and a random process A . Here, A is a sum of m Dirac masses, drawn uniformly on $[0, 1]^d$, with d the dimension of the considered problem (which is both the dimension of the random process and the dimension of the strategy). We have used typically \mathcal{A} the family of distributions equal to sum of m Dirac masses; A is randomly selected in \mathcal{A} by placing m Dirac masses uniformly in $[0, 1]^d$. \mathcal{A}' is the family of sums of p Dirac masses, where $p \in [[10, 50]]$ depending upon the number of Dirac masses used for A . The function π mapping A to A' is then defined in the following way :

- m/p Dirac masses of A are mapped to each Dirac mass of A' . We choose m/p to be integer to make the experiments easier.
- The way used to map the masses of A to the masses of A' is such as $A - \pi(A)$ is not large to ensure that the Taylor-approximations are relevant. In dimension 1, we simply sort the masses A and A' , and we map the m/p smallest values of A to the smallest value of A' and so on. In dimension > 1 , the method to map A to A' is more tricky, but the principle is the same.

The figure 2 gives, for a given distribution A the expectation of the cost in function of the strategies for a problem in dimension 1.

3.3 Algorithm

The experimental setup is as follow, with A a randomly drawn random process :

- Draw randomly a large number of A' (p Dirac masses, with $p < m$). We reject the A' which are not enough close to A , because the approximations we have made are only valid close to A . We construct $\pi(A) = A'$ as stated above.
- For each random process A' calculate the optimal strategy s ($s = Opt(A')$), and the expected associated cost (for the real random process A) $c(s) = EC(s, A)$.
- For each distance considered, sort the Dirac masses of the random process A' by their distance to A .
- Plot the expected cost of the optimal strategy ($c(s)$) for each distance : at abscissa n , plot the expected cost of the optimal strategy for the n^{th} closest A' according to this distance. If the distance is efficient, then the cost should get lower and lower as the distance decreases.

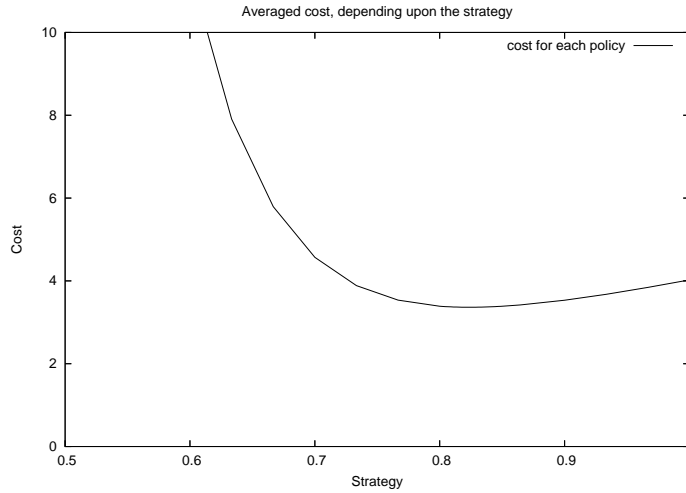


Figure 2: **Dimension 1** : expectation of the cost for a given distribution A , in function of the different strategies. We see that as in many supply chain problems, if the strategy is too short, then the cost is very large (strong derivative for small values of the strategy), whereas if the strategy is too long, the marginal cost is much smaller (small derivative). This means that being too short is much more detrimental than being too long (for the same deviation).

3.4 Results without noise

The figures 3 and 4 show the expected cost for random process A for the optimal strategy optimized for random process A' , in function of the rank of the distance between A and A' .

For a “good” distance we expect that small distances are associated to a small cost. We can see that the Kantorovitch-Rubinstein distance gives very bad results because the cost is not correlated to the distance, whereas the $distance^{(3)}$ is the best, and $distance^{(2)}$ gives also good results.

Hence, if one chooses a random process A' because of its small Kantorovitch-Rubinstein distance to the true random process A , and uses this A' as a model for the optimization he might get a very bad strategy, whenever the cost for A' is optimal.

We see on the zoomed part (below) of the figure 3 that $distance^{(2)}$ and $distance^{(3)}$ are very efficient. Indeed, the expected cost is very low in a wide range of chosen A' . The non zoomed part (above) of the figure 3 shows that the distances we have proposed measures efficiently the distance between A and A' in the sense that when the distance is large, the optimization of the cost for A' is not efficient as a strategy for A .

Figure 4 shows that the results remain in dimension 2.

3.5 Results with an additive noise

We keep the same protocol except we add a Gaussian noise of standard deviation 0.07 (quite important noise in regard to the data) as defined in section 2.3. As previously, we focus on the very first points of the curves (the smallest distances). As we have shown in section 3.4 that $distance^{(2)'}$ has no interest, we do not represent it any more. We always have $distance^{(3)}$ better than $distance^{(2)}$, itself better than the Kantorovitch-Rubinstein distance. The new distance " $distance^{(2)}$ with risk", tries to take in account the additive noise. Notice that we only assume we are able to evaluate the variance of this noise (a realistic hypothesis) but we verify it on a "nice" distribution shift, with the same law and the same parameters as in the derivation of the distance. This nevertheless shows the mathematical consistency of the derivation. The curve of the " $distance^{(2)}$ with risk" (see 5) is far better than the three other ones. This figure shows that our way to integrate noise into the distance can lead to significant improvements on the evaluations of the distance between random variables. Note that the level of noise is of course usually not available ; in our experiment, we used the *real* K value.

4 Conclusions

We have shown how to define a distance between A and A' taking into account derivatives of $C(s_0, A)$ in the research of a good representant A' of A . Roughly speaking, the interpretation of the criterion leads to two ways of selecting A' :

- put more points where the second derivative is strong ;
- counterbalance over-estimates $\pi(A)$ of A for some values of the random process A by under-estimates of similar amplitude.

Our approach quantifies numerically these elements. Figure 3 shows a very clear correlation between our distance and the loss function in A , which fully confirms our theoretical evaluations.

Another important feature of our method is that the risk can be taken into account in a very simple and efficient manner. Results of figure 5 show a very convincing improvement in the case of risk : adding the term of risk in our distance leads to a much better expectation on the noisy version of A . Our framework is natural and avoids the use of ad-hoc adaptations of the notion of risk such as conditional Value-At-Risk ; our notion of risk can take into account additive noise (possibly *dependent* on the realization of A through $K = K(A)$), thus modeling a variable uncertainty on A ; increasing K increasing the robustness to noise, in a simple, natural and convincing manner.

Our theoretical results are mathematically simple and sound (simply the use of Taylor's expansions). We point out the fact that they are just the natural extension of the derivation of Kantorovitch-Rubinstein bounds with more terms in the expansion. Our experiments confirm the theoretical results both for optimization in expectation and for the difficult problem of optimization with risk management.

We consider that these results are significant in the sense that :

- providing a framework that naturally extends to uncertainty management is important itself. In particular, the main argument for risk management is uncertainty ; dealing with risk management by the introduction of a noise on the random process is thus natural. This does not introduce any complexity in the model.
- stochastic dynamic optimization by dynamic programming methods is a very important real-world problem, which has not been much investigated by the AI community yet, whereas many future advances, according to us, will deal with the representation of random process and can come from statistics and paradigms of artificial intelligence ; dynamic programming remains the only reliable solution for many industrial problems ;
- only little attention has been devoted to the representation of random processes in stochastic dynamic optimization, whereas this part of the process appears in many cases as the most important part ; working on complex random processes in dynamic programming leads to very huge state spaces (the Bellman value has to be indexed on the full memory as soon as the random processes involved are not markovian) ; the nice representation of random processes is thus decisive with regard to the curse of dimensionality ;
- derivatives required here are often easily available (e.g., through Bellman values) ; the first derivative is a direct computation from the Bellman values and the transition function, the second derivative is more problem specific but follows naturally in many frameworks (e.g. parametric Bellman values parameterized by both the random process and the state space, or stock management with)
- the (relative) simplicity of our arguments is one more argument in favor of the method ;
- the optimization under constraint of our distance is easily performed (it is not convex, but it is polynomial and far from the state of the art of non-linear programming ; see e.g. [Hoyland et al, 2001]).

We now insist on the fact that our approach is directly operational and fully tractable. Implementing our approach requires 1) restricting the attention to processes such that $\pi(A) - A$ is small enough, for any realization, to ensure that the Taylor expansion is valid 2) compute the required distances 3) optimize one of our Taylor-based distances. 1) is performed by introducing constraints in the optimization. 2) is immediate for the first derivative and the diagonal of the Hessian if approximate Bellman values are available ; the general case of the second derivatives is more tricky and problem-specific but is tractable in many cases and in particular in stock management, which is an important area of dynamic programming, or whenever the bellman values are parametric and parameterized by both the state space and a random process value. 3) can be handled with non-linear programming with automatic differentiation for an optimization of the Taylor-based distance under constraint on $\pi(A) - A$; non-convexity can be handled thanks to multiple initial points as in [Hoyland et al, 2001].

In particular, direct applications are resource management for energy production, where modelization of short, middle and long-term meteorological random processes are very important and which is a main provider of real-world problems for stochastic dynamic programming, or portfolio management, where many different forms of unproperly defined random processes are concerned.

References

- [Bertsekas, 1995] D. Bertsekas. Dynamic programming and optimal control. Athena Scientific, 1995.
- [Aïd et al, 2003] R. Aïd, V. Grellier, A. Renaud, O. Teytaud. Application de l'apprentissage par renforcement à la gestion du risque. Proceedings of CAP 2003, in press.
- [Bertsekas et al, 1996] D.P. Bertsekas, J.N. Tsitsiklis. Neuro-dynamic programming, Athena Scientific.
- [Coraluppi et al, 1999] S. Corallupi, S. Marcus, Risk-sensitive and minimax control of discrete-time, finite state Markov Decision Processes. Automatica, 35, 301-309.
- [Geibel, 2001] P. Geibel, Reinforcement Learning With Bounded Risk, In: C. E. Brodley, and A. P. Danyluk, editors, "Machine Learning - Proceedings of the Eighteenth International Conference (ICML01)", pages 162-169. Morgan Kaufmann Publishers, San Francisco, CA.
- [Gröwe-Kuska et al, 2003] N. Gröwe-Kuska, H. Heitsch, W. Römisch, Scenario Reduction and Scenario Tree Construction for Power Management Problems, IEEE Bologna Power Tech Proceedings (A. Borghetti, C.A. Nucci, M. Paolone eds.), 2003.
- [Heger, 1994] M. Heger, Consideration of risk in reinforcement learning. Proceedings of ECML pp105-111, Morgan Kaufman, 1994.
- [Hoyland et al, 2001] K. Hoyland, S.W. Wallace, Generating Scenario Trees for Multistage Decision Problems, Informs 47, 2, 2001.
- [Keefer, 1994] D.L. Keefer, Certainty equivalents for three-point discrete-distribution approximations. Management Sciences 40, 760-773.
- [Nagayuki et al, 2000] Y. Nagayuki, S. Ishii, and K. Doya. Multi-agent reinforcement learning: An approach based on the other agent's internal model. In Proc. of the 4th Intl. Conf. on MultiAgent Systems.
- [Neuneier, 1996] R. Neuneir, Optimal asset allocation using adaptive dynamic programming, in Advances in Neural Information Processings Systems, D.S. Touretzky, M.C. Mozer, M.E. Hasjselmo, eds, vol. 8, MIT Press.

- [Neuneier, 1998] R. Neuneier, Enhancing Q-learning for optimal asset allocation, in Advances in Neural Information Processing Systems, M.I. Jordan, M.J. Kearns, S.A. Solla, eds, vol.10, MIT Press.
- [Neuneier et al, 1999] R. Neuneier, O. Mihatsch, Risk-sensitive reinforcement learning, MIT Press, NIPS'99.
- [Pflug, 2001] G.C. Pflug, Scenario Tree Generation for Multi-Period Financial Optimization by Optimal Discretization, Mathematical Programming, Series B 89:251-271, 2001.
- [Riedmiller et al, 2001] M. Riedmiller, A. Moore, J. Schneider, Reinforcement Learning for Cooperating and Communicating Reactive Agents in Electrical Power Grids, Balancing Reactivity and Social Deliberation in Multi-Agent Systems.
- [Sutton et al, 1998] R.S. Sutton, A.G. Barto, Reinforcement learning.
- [Tesauro, 1989] G. Tesauro. Neurogammon wins Computer Olympiad. Neural Computation 1, 321-323.
- [Tesauro, 1999] G. Tesauro. Pricing in agent economies using neural networks and multi-agent Q-learning. Proceedings of Workshop ABS-3: Learning About, From and With other Agents (held in conjunction with IJCAI '99).
- [Watkins, 1989] C. Watkins, Learning from Delayed Reward, Ph.D thesis, Kings College, University of Cambridge.

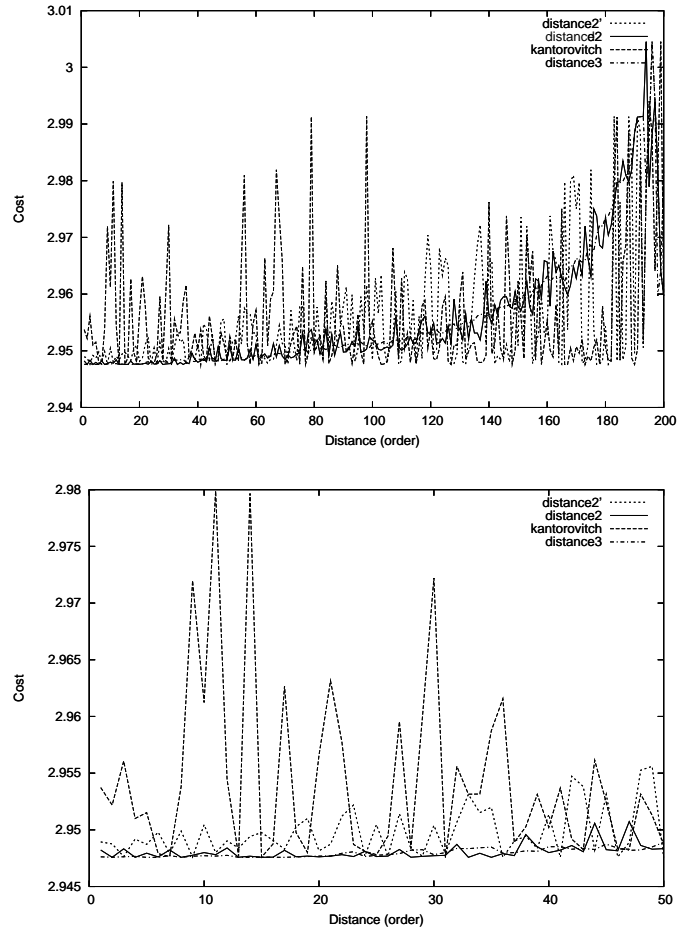


Figure 3: Results in **dimension 1** without noise. X-coordinate is the rank of A' according to the considered distance. For a “good” distance we expect that small distances are associated to a small cost. The second graphic is a zoom on the smallest distances between A and $\pi(A) = A'$. We can see that the Kantorovich-Rubinstein distance gives very bad results because the cost is not correlated to the distance, whereas the $distance^{(3)}$ is the best, and $distance^{(2)}$ gives also good results. Interpretation : a small rank for distance $distance^{(2)}(A, A')$ or $distance^{(3)}(A, A')$ leads to a good behavior when $Opt(A')$ is used on the random process A and these distances are highly correlated to the cost in generalization.

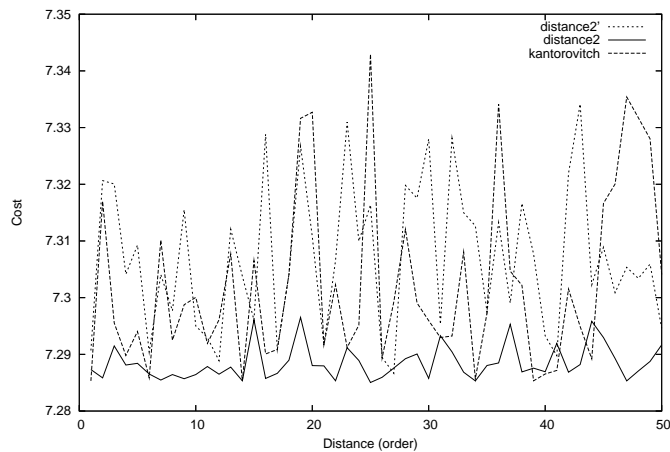


Figure 4: Results in **dimension 2**. This is the same graphic (zoom on the 50 smallest distances) using a two dimensions function instead of dimension one. distance3 is not plotted because we haven't implemented it in dimension greater than 1. We see that we get the same results in dimension 2, that is to say that the Kantorovitch-Rubinstein distance is not correlated to the cost, whereas the $distance^{(2)}$ leads to good results. Hence, choosing a random process A' "close" to A in the Kantorovitch-Rubinstein distance leads to have a very bad estimation of the best cost, and leads to get an apparently-optimal solution (when looking at the cost for A') very far from the real optimum.

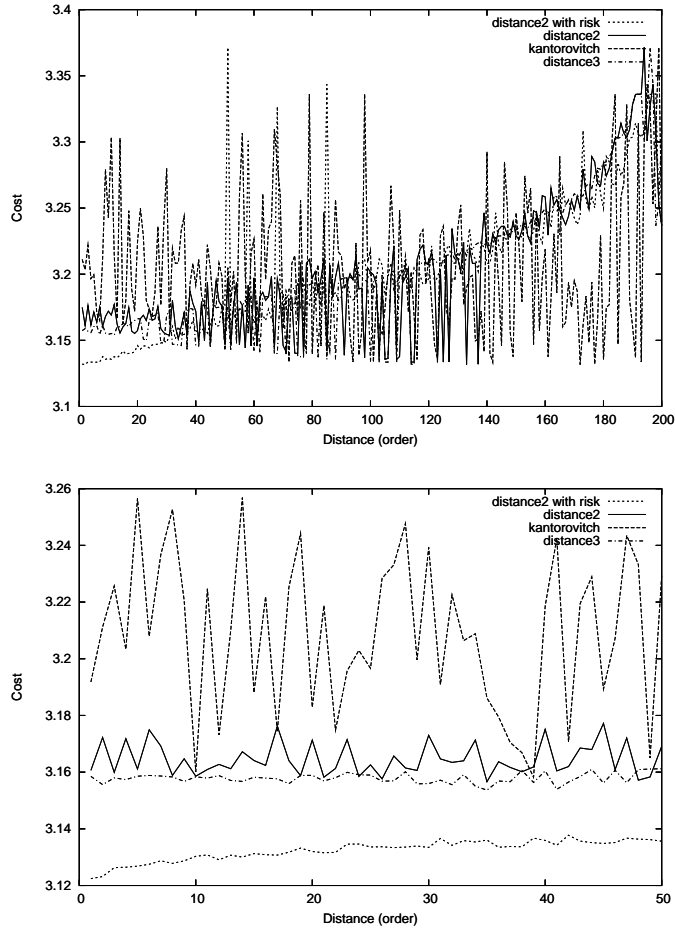


Figure 5: Results in **dimension 1 with noise**. X-coordinate is the rank of A' according to the considered distance. For a "good" distance we expect that small distances are associated to a small cost. The second graphic is a zoom on the smallest distances between A and $\pi(A) = A'$ and with more experiments (the graphic becomes unreadable with too much points on it). " $distance^{(2)}$ with risk" corresponds to the distance introduced in the section 2.3. We can see that the Kantorovitch-Rubinstein distance gives again bad results since the cost is not correlated to the distance, whereas the " $distance^{(2)}$ ", " $distance^{(3)}$ " and " $distance^{(2)}$ with risk" are better. The really interesting point here is that the " $distance^{(2)}$ with risk" gives much better results than even " $distance^{(3)}$ " which worked very well without noise. Hence, the added term in " $distance^{(2)}$ with risk" ($\frac{1}{2} E Tr(H, A)K$) is here validated.