

Pricing Swing Options by Reinforcement Learning¹

Roberto Daluiso (Intesa Sanpaolo)

joint work with:

Emanuele Nastasi (Marketz S.p.A.)

Andrea Pallavicini (Intesa Sanpaolo)

Giulio Sartorelli (Intesa Sanpaolo)

Big Data and Machine Learning in Finance Conference
Politecnico di Milano, June 10, 2021

¹Daluiso et al. 2020 (<https://arxiv.org/abs/2001.08906>)

Contents

- 1 Introduction
- 2 Least-square Monte Carlo solution
- 3 Reinforcement learning solution
- 4 Conclusion

Contents

- 1 Introduction
- 2 Least-square Monte Carlo solution
- 3 Reinforcement learning solution
- 4 Conclusion

Motivation

Several pricing problems fit in the framework of stochastic control:

- Option contracts depend on decisions by the holder
- So their fair price is a supremum across all admissible choices

Challenges for traditional approaches

- 1 Complex dependence of the payoff on multiple decisions
- 2 High dimensional state space (e.g. basket options)

Reinforcement Learning (RL) may help:

- 1 Model-agnostic: learns the impact of actions by trial and error
- 2 Scalable: leverages machine learning to learn value functions

Swing Options

- A swing option guarantees a flexible supply of a commodity
- The underlying is the day-ahead future $F_{T_i} := F_{T_i}(T_{i+1}, 1d)$ for each fixing date T_1, \dots, T_{n_f} in the delivery period
- At each fixing date the holder chooses a quantity N_{T_i} and gets it at the strike price K . His choice is constrained by:

Local constraints: $N_{T_i} \in [N_m, N_M]$

Global constraint: $\sum_{i=1}^{n_f} N_{T_i} \in [C_m, C_M]$

Benchmarks

We purposely choose a well-studied problem, so that we can check the output of RL both quantitatively and qualitatively:

Numerical benchmark

- We implement a Least-Square Monte Carlo algorithm (LSMC), similarly to Barrera-Esteve et al. 2006
- Many alternative proposals exist in the literature

Theorem (Bardou et al. 2010)

*If the lower bound C_m and width $C_M - C_m$ of the global constraint are integer multiples of the width $N_M - N_m$ of the local constraint, then there exists an optimal control of **bang-bang** type, i.e. such that at each fixing date one chooses either the minimum or the maximum admissible consumption given local and global constraints.*

Contents

- 1 Introduction
- 2 Least-square Monte Carlo solution
- 3 Reinforcement learning solution
- 4 Conclusion

Dynamic Programming Principle

Suppose zero interest rates for simplicity, and define

$$C_{T_i} := \sum_{j=1}^i N_{T_j} \quad \text{total consumption up to time } T_i.$$

Then the value W_{T_i} of the option satisfies the backwards recursion

$$W_{T_i}(F_{T_i}, C_{T_{i-1}}) = \max_{N_{T_i}} \{ N_{T_i}(F_{T_i} - K) + \mathbb{E}_{T_i}[W_{T_{i+1}}(F_{T_{i+1}}, C_{T_i})] \}$$

Implicit assumptions

- 1 The day-ahead future F_t is Markov (because of the model)
- 2 W_{T_i} depends on past decisions only via $C_{T_{i-1}}$ (because of the type of constraints)

Hence the LSMC approach needs a redesign each time the model and/or contract details change.

Discretization

The backward recursion is implicit, because at step T_i the optimal C_{T_i} is yet unknown:

$$W_{T_i}(F_{T_i}, C_{T_{i-1}}) = \max_{N_{T_i}} \{ N_{T_i}(F_{T_i} - K) + \mathbb{E}_{T_i} [W_{T_{i+1}}(F_{T_{i+1}}, C_{T_i})] \}$$

A way out is to discretize the set of attainable total consumptions:

Discretized recursion

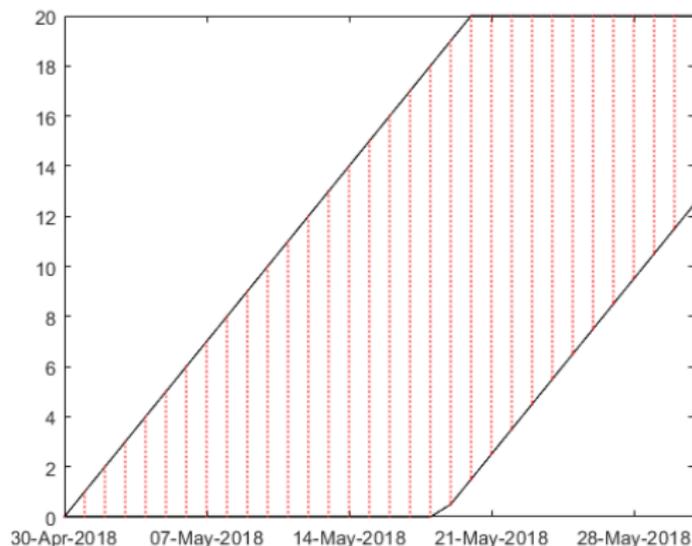
At step T_i and for each consumption C_ℓ^{i-1} of a grid C^{i-1} :

- 1 Define the set $Q^{T_i}(C_\ell^{i-1}) = \{ N_{T_i} \in [N_m, N_M] : N_{T_i} = C_j^i - C_\ell^{i-1} \}$
- 2 Approximate $W_{T_i}(F_{T_i}, C_\ell^{i-1})$ by a max on this finite set:

$$\max_{N \in Q^{T_i}(C_\ell^{i-1})} \{ N(F_{T_i} - K) + \mathbb{E}_{T_i} [W_{T_{i+1}}(F_{T_{i+1}}, C_\ell^{i-1} + N)] \}$$

- 3 Use regression to estimate the expected values.

Consumption Grid

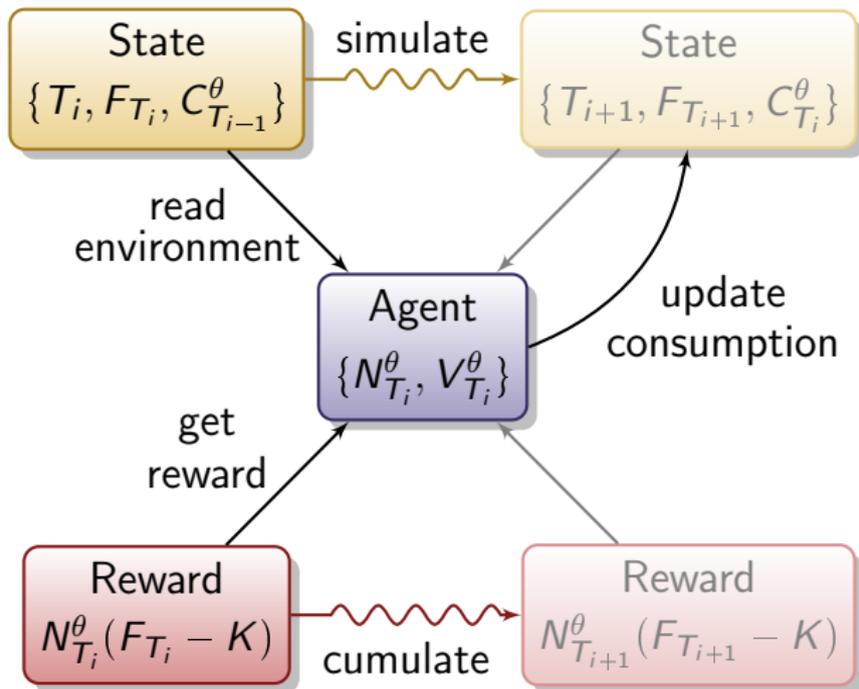


Total consumption grid for a swing delivering in the month of April 2019 with $N_M = 1$, $N_m = 0$, $C_M = 15.7$, $C_m = 5.2$.

Contents

- 1 Introduction
- 2 Least-square Monte Carlo solution
- 3 Reinforcement learning solution**
- 4 Conclusion

(Model-Free) Reinforcement Learning



A RL agent updates his policy while it interacts with the environment, so that the new N^θ improves his estimate V^θ of expected rewards.

RL Families

RL algorithms may be classified by the function which is learned (our choice in red):

Learned object

- *Value Iteration*: learn the expected cumulated reward $V(x)$ on the *optimal* policy as a function of the current state x
- *Q-Learning*: learn the expected cumulated reward $Q(x, a)$ of using an action a from state x and *then* acting optimally
- *Policy Optimization*: learn an action $\pi(x)$ maximizing an estimate of *the policy's* cumulated future rewards

Proximal Policy Optimization

We use the Proximal Policy Optimization (PPO) algorithm proposed in Schulman et al. 2017²:

- Well-suited for continuous control problems
- *Actor-critic* algorithm: both the policy π^θ and an estimate V^θ of its value function are maintained

Key points of PPO

- 1 Policies are randomized, hence represented by a pdf $\pi^\theta(a|x)$:
 - If \mathcal{A} is discrete, a network outputs the probability of each action
 - Otherwise, a network outputs the mean $N_{T_i}^\theta$ of a Gaussian
- 2 The algorithm updates θ by batch stochastic gradient descent
- 3 The objective function is clipped to discourage large steps, as in trust-region optimization

²We use the implementation of the algorithm found in OpenAI Baselines
<https://github.com/openai/baselines>

Setting of Numerical Examples

- Strike fixed at-the-money at inception
- Model in Daluiso et al. 2020 with mean reversion equal to 1
- Prices computed by MC on an out-of-sample set of 10^6 paths

Network inputs

Reasonably-normalized inputs help training considerably:

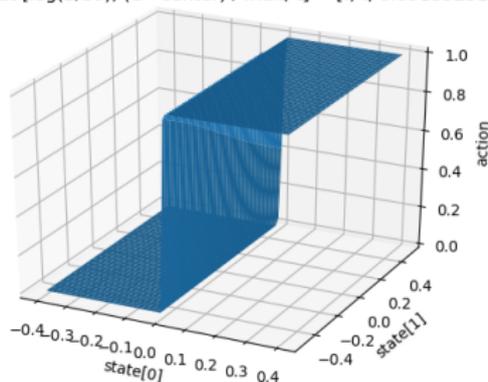
- 1 T_i is expressed as a year fraction
- 2 The total consumption to-date is remapped linearly at each time so that its domain is always $[-0.5, 0.5]$
- 3 The underlying future is remapped as $\log(F_{T_i}/F_{T_0})$

Optimized Policy

Main experiment

- 1 Delivery period: one month
- 2 Constraints in MWh: $[N_m, N_M] = [0, 1]$, $[C_m, C_M] = [12, 20]$, satisfying the hypotheses for existence of bang-bang optima
- 3 Fixed 5x4 architecture, tuned on a shorter-lived option

Actions at $[\log(S/S_0), (C - \text{center}) / \max, t] = [., ., 0.091552511]$



Normalized consumption as a function of normalized log-spot and consumption, on the fourth decision date. Although *a priori* the action was continuous-valued, *a posteriori* the PPO algorithm detected an optimal policy which is of bang-bang type.

Comparison to LSMC

Variants to the main experiment

- 1 C_m : 12 or 12.5 (to have non bang-bang optima)
- 2 Policies: continuous-valued or constrained to be bang-bang

	PPO	All Strategies	Only Bang-Bang
Out of Theorem Hyp.		7.92 ± 0.04	7.74 ± 0.04
Within Theorem Hyp.		8.40 ± 0.04	8.37 ± 0.04
	LSMC	All Strategies	Only Bang-Bang
Out of Theorem Hyp.		7.97 ± 0.04	7.76 ± 0.04
Within Theorem Hyp.		8.46 ± 0.04	8.46 ± 0.04

Contents

- 1 Introduction
- 2 Least-square Monte Carlo solution
- 3 Reinforcement learning solution
- 4 Conclusion

Summary and further developments

We priced swing options both by least-square Monte Carlo approach and by a recent RL algorithm (Proximal Policy Optimization).

Main findings

- 1 When bang-bang strategies are theoretically optimal, both algorithms can detect it; otherwise, a slight improvement is obtained if the bang-bang constraint on consumption is removed
- 2 RL can recover the prices obtained by LSMC, despite the fact that the latter was tailored to the swing option pricing task

Possible developments

Try RL in higher dimensional settings where traditional methods fail:

- 1 Swing options in more realistic models with jumps and/or multi-factor curve dynamics (Nastasi et al. [2018](#))
- 2 Pricing problems involving baskets/portfolios

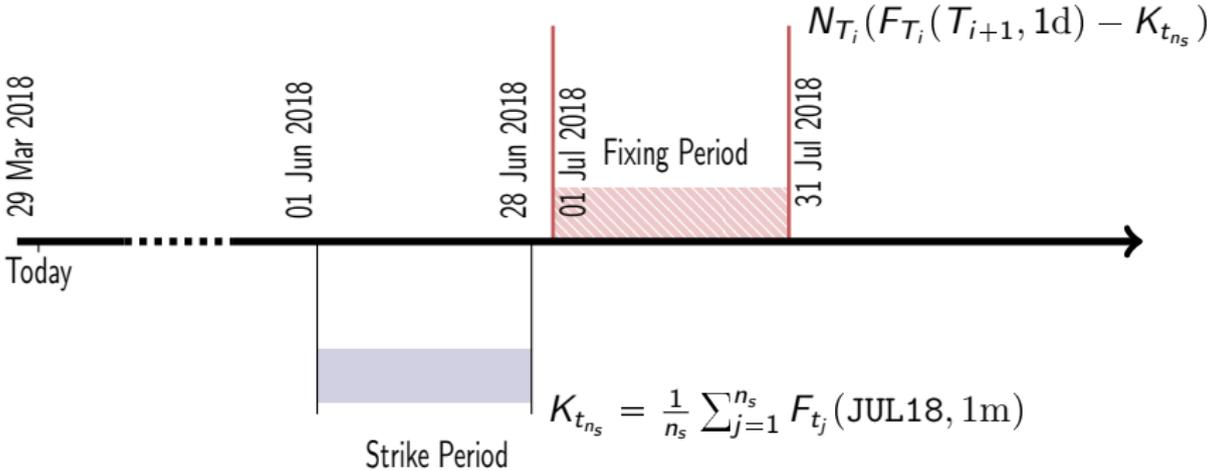
Selected references

- Bardou, O., S. Bouthemy, and G. Pagès (2010). “When are swing options bang-bang?” In: *International Journal of Theoretical and Applied Finance* 13.6, pp. 867–899.
- Barrera-Esteve, C. et al. (2006). “Numerical methods for the pricing of Swing options: a stochastic control approach”. In: *Methodology and Computing in Applied Probability* 8.4, pp. 517–540.
- Daluiso, R. et al. (2020). “Pricing Commodity Swing Options”. In: *Working paper*. URL: <https://arxiv.org/abs/2001.08906>.
- Nastasi, E., A. Pallavicini, and G. Sartorelli (2018). “Smile Modelling in Commodity Markets”. In: *Working paper*. URL: <https://arxiv.org/abs/1808.09685>.
- Schulman, J. et al. (2017). “Proximal Policy Optimization Algorithms”. In: *Working paper*. URL: <https://arxiv.org/abs/1707.06347>.

The opinions expressed in this work are solely those of the authors and do not represent in any way those of their current and past employers.

Backup Slides

Example Contract



Term-sheet data for a swing option with delivery in July 2018 in TTF natural gas market. The strike is fixed by averaging the JUL18 futures contract observed in the month of June.

Dynamic Programming by Regression

Now given a set of MC paths $F^{(k)}$ we approximate the expectation

$$\mathbb{E} \left[W_{T_{i+1}} (F_{T_{i+1}}, C_\ell^{i-1} + N) \mid F_{T_i} = F_{T_i}^{(k)} \right] \approx f_{T_i} (F_{T_i}^{(k)}; C_\ell^{i-1} + N)$$

supposing that it is quadratic in $F_{T_i}^{(k)}$:

$$f_{T_i} (F; C) := \alpha^i (C) + \beta^i (C) F + \gamma^i (C) F^2$$

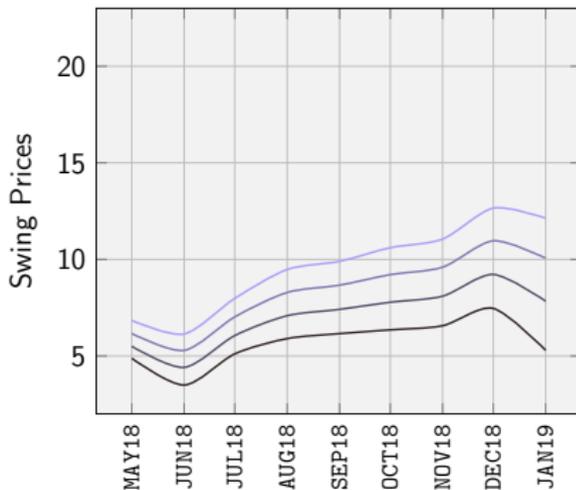
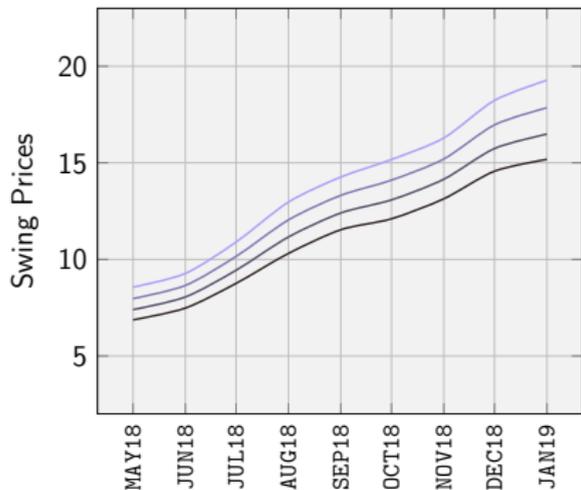
Remark

$C_\ell^{i-1} + N \in C^i$ by construction, hence we can estimate all needed coefficients by regressing for each j the realizations

$$y^{(k)} := W_{T_{i+1}} (F_{T_{i+1}}^{(k)}, C_j^i) \quad \text{against} \quad x^{(k)} := F_{T_i}^{(k)}.$$

The algorithm is initialized setting W_{T_f} to the contractual payoff.

Results



Swing option prices by varying the delivery starting date. We use the model proposed in Daluiso et al. 2020, with mean reversion parameter ranging from top to bottom from 1.5 to 0 with a step of 0.5. Left panel fixed-strike contracts. Right panel floating-strike contracts.

PPO Focus: Advantage

Notation

- $\tilde{N}_{T_i}^\theta$ randomized action, π^θ -distributed given the status at T_i
- $r_{T_i}(\tilde{N}_{T_i}^\theta) = \tilde{N}_{T_i}^\theta(F_{T_i} - K)$ reward at time T_i
- $\tilde{V}_{T_i}^\theta = \mathbb{E}_{T_i}[r_{T_i}(\tilde{N}_{T_i}^\theta) + \tilde{V}_{T_{i+1}}^\theta]$ value function of π^θ
- $\tilde{Q}_{T_i}^\theta(n) = \mathbb{E}_{T_i}[r_{T_i}(n) + \tilde{Q}_{T_{i+1}}^\theta(\tilde{N}_{T_{i+1}}^\theta)]$ action-value function of π^θ

- 1 In gradient-based policy optimization, one must calculate

$$\nabla_\theta \mathbb{E} \left[\sum_{i=1}^{n_f} r_{T_i}(\tilde{N}_{T_i}^\theta) \right] \stackrel{\text{Fact}}{=} \sum_i \mathbb{E} \left[A_{T_i}^\theta(\nabla_\theta \log \pi_{T_i}^\theta(n)) |_{n=\tilde{N}_{T_i}^\theta} \right]$$

where $A_{T_i}^\theta := \tilde{Q}_{T_i}^\theta(\tilde{N}_{T_i}^\theta) - \tilde{V}_{T_i}^\theta$ is called **advantage**

- 2 Since $\tilde{Q}_{T_i}^\theta$ is unknown, $A_{T_i}^\theta$ is then replaced by an estimator \hat{A}_i^θ
- 3 See paper for details

PPO Focus: Gradient Step

$$\theta_{k+1} = \theta_k + \rho \cdot \mathbb{E} \left[\nabla_{\theta} \sum_{i=1}^{n_f} (L_{T_i}^A(\theta) - \beta L_{T_i}^V(\theta)) \Big|_{\theta=\theta_k} \right]$$
$$L_{T_i}^A(\theta) := \min \left\{ \frac{\pi_{T_i}^{\theta}(\tilde{N}_{T_i}^{\theta_k})}{\pi_{T_i}^{\theta_k}(\tilde{N}_{T_i}^{\theta_k})} \hat{A}_i^{\theta_k}, \text{clip} \left(1 - \varepsilon, \frac{\pi_{T_i}^{\theta}(\tilde{N}_{T_i}^{\theta_k})}{\pi_{T_i}^{\theta_k}(\tilde{N}_{T_i}^{\theta_k})}, 1 + \varepsilon \right) \hat{A}_i^{\theta_k} \right\}$$
$$L_{T_i}^V(\theta) := \left(V_{T_i}^{\theta} - (\hat{A}_i^{\theta_k} + V_{T_i}^{\theta_k}) \right)^2$$

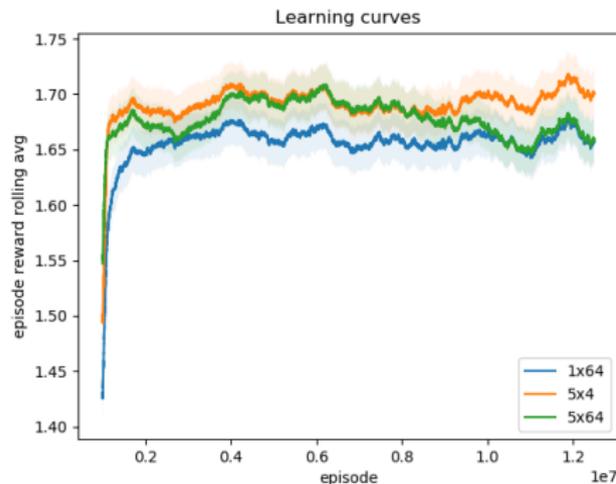
Intuition

- \mathbb{E} is estimated by a sample average on a batch
- ρ is the learning rate
- $\nabla_{\theta} L^A(\theta)$ is a clipping of the gradient derived in the previous slide
- $L^V(\theta)$ is smaller if V^{θ} represents better the expected rewards
- β is an hyper-parameter to balance the two effects

Neural Network Architecture Tuning

Preliminary experiment

- 1 Short delivery period (one week)
- 2 Constraints in MWh: $[N_m, N_M] = [0, 1]$, $[C_m, C_M] = [3, 5]$
- 3 Try wide and/or deep architectures



Learning curves. The solid lines are the moving average of the realized rewards on the last 10^6 episodes. The shadows represent the 98% confidence intervals.