

Ecole Polytechnique
Department of Applied Mathematics : **Enseignement d'approfondissement**

Supervised by : Emmanuel BACRY, Iacopo MASTROMATTEO

Optimal investment strategy in high-frequency trading

Bachir EL KHADIR, Mohamed Anas OUAALINE

Paris, le November 24, 2014

Time	Type	TradedPrice	BidPrice	BidQty	BidSplit	AskPrice	AskQty
28874.880073		-	141.57	51	1	141.58	81
28874.884174	AskMar	141.58	141.57	51	1	141.58	81
28874.884174		-	141.57	51	1	141.58	80
28874.887593	BidMar	141.57	141.57	51	1	141.58	80
28874.887593		-	141.57	50	1	141.58	80
28874.893296	AskMar	141.58	141.57	50	1	141.58	80
28874.893296		-	141.57	50	1	141.58	79
28874.89368	AskLim	-	141.57	50	1	141.58	84
28874.893949	AskLim	-	141.57	50	1	141.58	87

Contents

1	Presentation of the problem	2
1.1	Problematic	2
1.2	The model and its dynamic	2
1.2.1	The model	2
1.2.2	The dynamic	3
1.3	The optimization problem	3
1.4	The parameter c	4
2	Solving the minimization problem	5
2.1	First time of reach for a sum of a Markov chain	5
2.2	Expectation of stopped $(S_k)_k$	6
3	Numerical Results	7
3.1	Model results	7
3.2	Results from market data and validation	9
4	Model limitations and possible improvements	10
A	Technical proofs	11
A.1	Reccurence equation of $(g_{-1}^s)_s$	11
A.2	Taylor expansion of g_{-1}^s	11
B	The law of a sum of dependent random variables	12
C	Source Code	13

1 Presentation of the problem

1.1 Problematic

Given the state of an order book, we want to place purchase orders. Trought this document, we try to define an efficient and optimal strategy in order to reduce the cost of such orders. Since it is quite complex to define this kind of strategy, we focus on the case where we have only one purchase order to make and a time horizon before which the trade must be excuted. We have the choice between executing the order directly (market order), or placing a limit order. Since we are never sure that a limit order will be executed, we have to remove it after a certain time

k and place a market order at the best offer available. Notice that the case where $k = 0$ is when a market order is executed immediately. The cost related to these kind of strategies can be of two types:

- Cost related to the spread when executing a market order.
- Cost related to the uncertainty of future price movements that we denote c .

While the first is easily measurable, the second is dependant on many factors: volatility of the market, the risk aversion of the trader... We will discuss this issue further in the following sections. The goal of this document becomes:

Given the state of an order book and a single purchase order we want to make, What it is the best choice of k to minimize the cost of the trade ?

1.2 The model and its dynamic

Let us present the basic model which our study is based on.

1.2.1 The model

The order book is represented by its best ask and bid queues. The state of the limit order book is represented by :

- the bid price s_t^a and the ask price s_t^b
- the size of the bid queue q_t^a representing the outstanding limit buy orders at the bid, and
- the size of the ask queue q_t^b representing the outstanding limit sell orders at the ask

The bid and ask prices are multiples of the tick size $\delta = s_t^a - s_t^b$ which we consider constant. The state of the limit order book is thus described by the triplet $X_t = (s_t^b, q_t^b, q_t^a)$ which takes values in the discrete state space $\delta\mathbb{Z} \times \mathbb{N}^2$. With no loss of generality, we can assume that :

$$\delta = 1 \tag{1}$$

We assume that orders arriving at the best bid/ask follow a poissonion process and are all of the same size 1:

- A market/cancelation order arriving at the q_t^a (resp q_t^b) reduces its size by 1 if it is not empty
- A limit order arriving at the q_t^a (resp q_t^b) adds 1 to its size

Therefore q_t^a and q_t^b follow a birth/death process. We denote :

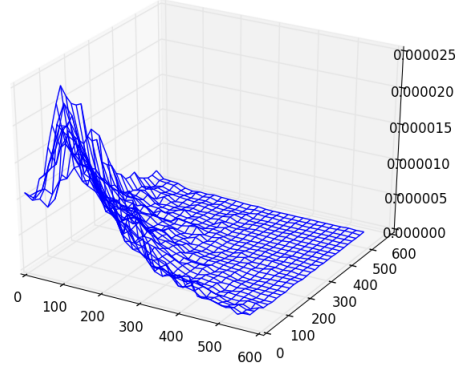
- $(\tau_i)_{i \geq 1}$ the sequence of random variables which stands for times between successive price changes. ¹
- $(X_i)_{i \geq 1}$ the successive moves in the price. Note that according to our hypothesis ($\delta = 1$) : $X_i \in \{-1, 1\}$.
- $(S_k)_{k \geq 0}$ the cumulative moves in the price after the n^{th} change, ie $S_n = \sum_{1 \leq i \leq n} X_i$. ³

1.2.2 The dynamic

The state of the order book is modified by order book events. When the bid (resp. ask) queue is empty, the price moves down (resp. up), and (q_t^a, q_t^b) (resp. (q_t^b, q_t^a)) is generated from a distribution $f(x, y)dx dy$ (resp. $\tilde{f}(x, y)dx dy$) independently from the rest of the historical variables. f will be estimated from historical data:

¹Cont and Larrard's article presents an expression for the cumulative distribution function of each τ_i ².

³see Annexe B for the generating function of S_n



This property give the market a markovian property: the new state of the market after a price move does only depend on the direction of this move. Therefore, $(X_i)_{i \geq 0}$ is markov process such as :

- $\forall n \geq 0 \quad X_n \in \{-1, 1\}$
- $\forall i \geq 0 \quad \mathbb{P}(X_{i+1} = 1 | X_i = 1) = \mathbb{P}(X_{i+1} = -1 | X_i = -1) = p = 1 - q$

Given the symmetry of the market, we may assume that $f(x, y) = \tilde{f}(y, x)$. This assumption makes the process $(\tau_i)_i$ a sequence of IID random variables independent from $(X_i)_i$.

1.3 The optimization problem

Recall how our strategy works:

- A limit order is placed at $t=0$ at the best bid.
- after some time $T = \sum_{i \leq k} \tau_i$, if the order is not executed, we cancel it and make a market order at the ask side.

k is the parameter of the strategy we seek to optimize. It stands for the maximum number of moves in price we are willing to wait for. Let N be the first time where the sequence $(S_n)_{n \geq 0}$ achieves -1 . Thereby :

$$N = \inf\{n \geq 0, S_n = -1\}$$

Finally, we denote c_{wait} the waiting cost per unit of time. The problem could be seen as a minimization problem :

$$\min_{k \in \mathbb{N}} f(k) := \min_{k \in \mathbb{N}} \mathbb{E} \left[\sum_{i=1}^{k \wedge N} \tau_i c_{wait} + (S_{k \wedge N} + 1) \delta \right] \quad (2)$$

We shall notice that $(\tau_i)_{(i > 0)}$ are assumed to be independent. The are also assumed to be independent of N . Thus we can rewrite f and the problem becomes :

$$\min_{k \in \mathbb{N}} f(k) = \min_{k \in \mathbb{N}} \mathbb{E} [(k \wedge N) \mathbb{E}(\tau) c_{wait} + (S_{k \wedge N} + 1)] \quad (3)$$

Since $\mathbb{E}(\tau)$ is a constant we can denote $c = \mathbb{E}(\tau) c_{wait}$ which is the waiting cost per price move. And the minimization problem finally becomes :

$$\min_{k \in \mathbb{N}} f(k) = \min_{k \in \mathbb{N}} c \mathbb{E} [k \wedge N] + \mathbb{E} [S_{k \wedge N}] + 1 \quad (4)$$

In order to solve this problem we need to do 2 different calculations :

- $\mathbb{E} [k \wedge N]$: we are going to explicitly calculate the Law N
- $\mathbb{E} [S_{k \wedge N}]$: since S_i and N are not independent, the law is quite hard to be explicitly calculated. Consequently, we are going to use a Monte-Carlo Algorithm.

1.4 The parameter c

In the previous section we have denoted :

$$c = \mathbb{E}(\tau)c_{wait}$$

It is an exogenous parameter which has to manually fixed. We notice that through its definition c depends on $\mathbb{E}(\tau)$ which is the average time between two moves in the price. However we can be more accurate and see that it also depends on the volatility of the market and the operator's risk aversion. We also notice that c needs to be lower than the thick δ otherwise it is not interesting to wait for a move in the price. Thus we came up with the following remarks.

- $c = \mathbb{E}(\tau)F(\sigma, v, r)$ where $\mathbb{E}(\tau)$ is the average time between two moves in the price, σ is the market volatility, v is the operator's risk aversion and r is some unknown parameters.
- $c < \delta$ where δ is the thick.
- c is increasing in σ : The more the market is volatile the more it risky is to wait.
- c is increasing in v : The more the operator is averse to risk the more likely he is to over-rate the waiting cost.

In our model we assume that $\mathbb{E}(\tau), \sigma, v$ and r do not depends on k .

2 Solving the minimization problem

2.1 First time of reach for a sum of a Markov chain

In this section we adress the issue of finding the law of the variable N . Let us define the probability-generating function of N such as :

$$g_s^x(z) = \mathbb{E}(z^N | S_0 = s, X_0 = x) \quad \forall x \in \{-1, 1\}, s \in \mathbb{N}, z \in [0, 1]$$

so that

$$\mathbb{E}(z^N | S_0 = s) = \alpha g_s^1(z) + (1 - \alpha) g_s^{-1}$$

where $\alpha := \mathbb{P}(X_0 = -1)$.

Given that (S_n, X_n) is a markov chain, we can infer that for all $x \in \{-1, 1\}$, $s \in \mathbb{N}$ and $z \in]0, 1[$ ⁴ :

$$\begin{cases} g_s^1 = z(pg_{s+1}^1 + qg_{s-1}^{-1}) \\ g_s^{-1} = z(qg_{s+1}^1 + pg_{s-1}^{-1}) \end{cases} \quad (5)$$

$$(6)$$

Denoting :

$$u = \frac{1}{2} \left(\frac{1}{zp} + z(2 - \frac{1}{p}) \right)^5 \quad (7)$$

and by imposing a Dirichlet boundary condition, we can prove⁶ that g_s^1 dynamic is given by the following second order system :

$$\begin{cases} g_{s+1}^{-1} - 2ug_s^{-1} + g_{s-1}^{-1} = 0 \\ g_{-1}^{-1} = 1 \\ \lim_{s \rightarrow +\infty} g_s^{-1} = 0 \end{cases} \quad (8)$$

$$(9)$$

$$(10)$$

Which is a classical linear recurrent sequence of order 2. The solution is given by :

$$g_s^{-1} = \left(u - \sqrt{u^2 - 1} \right)^{s+1} \quad (11)$$

Changing 1 to -1 in (5) and (6) we see that g_s^1 follow the same dynamic as g_s^{-1} if we change p to $1 - p$, therefore we can deduce all the properties of the first from the second.

Since g_s^x is the probability-generating function of N , the law of N conditional to $S_0 = 0$ and $X_0 = x$ is given by :

$$\begin{cases} g_s^{-1}(z) = \sum_{j=0}^{+\infty} \mathbb{P}_{X_0=-1}(N = j | S_0 = s) z^j \\ g_s^1(z) = \sum_{j=0}^{+\infty} \mathbb{P}_{X_0=1}(N = j | S_0 = s) z^j \end{cases} \quad (12)$$

$$(13)$$

$\mathbb{P}_{X_0=-1}(N = j | S_0 = s)$ can be calculated using taylor expansion, and $\mathbb{P}_{X_0=1}(N = j | S_0 = s)$ can be deduced as discussed above.

⁴For clarity, we omit to note the dependence of g in z .

⁵ $u \geq 1$

⁶see Annexe

2.2 Expectation of stopped $(S_k)_k$

$E[S_{k \wedge N}]$ is calculated using Monte Carlo algorithm.

```
#include <iostream>
#include <iomanip>
#include <string>
#include <map>
#include <random>

using namespace std;

int main()
{
    int k = 5;
    double p = 0.2;

    // Number of independant simulations
    int M = 1e6;

    // generator of  $X_i$ 
    random_device rd;
    mt19937 gen(rd());
    map<int, bernoulli_distribution> alea;
    // generator of  $X_i$  conditional to  $X_{i-1} = 1$ 
    alea[1] = std::bernoulli_distribution(p);
    // generator of  $X_i$  conditional to  $X_{i-1} = -1$ 
    alea[-1] = std::bernoulli_distribution(1 - p);

    double mean = 0;
    for (int m = 0; m < M; m++) {
        int S = 0;
        int X0 = 1;
        for (int i = 0; i < k && S >= 0; i++) {
            X0 = 2*alea[X0](gen)-1;
            S += X0;
        }
        mean += S;
    }

    mean /= M;
    cout << mean << endl;

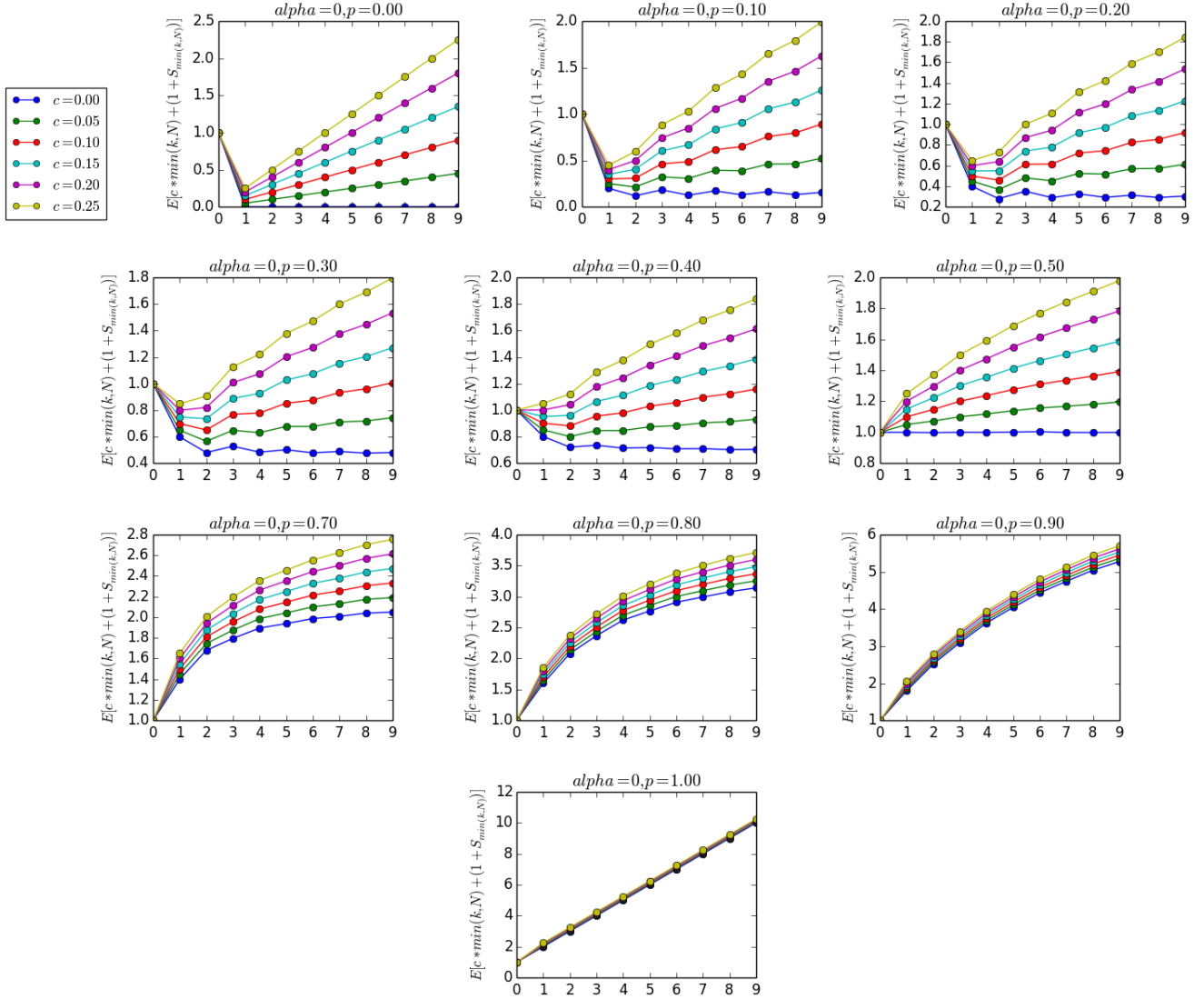
    return 0;
}
```

3 Numerical Results

3.1 Model results

Here are the results for different combination of the parameters (p, c, α) .

When $X_0 = 1$, ie $\alpha = 0$



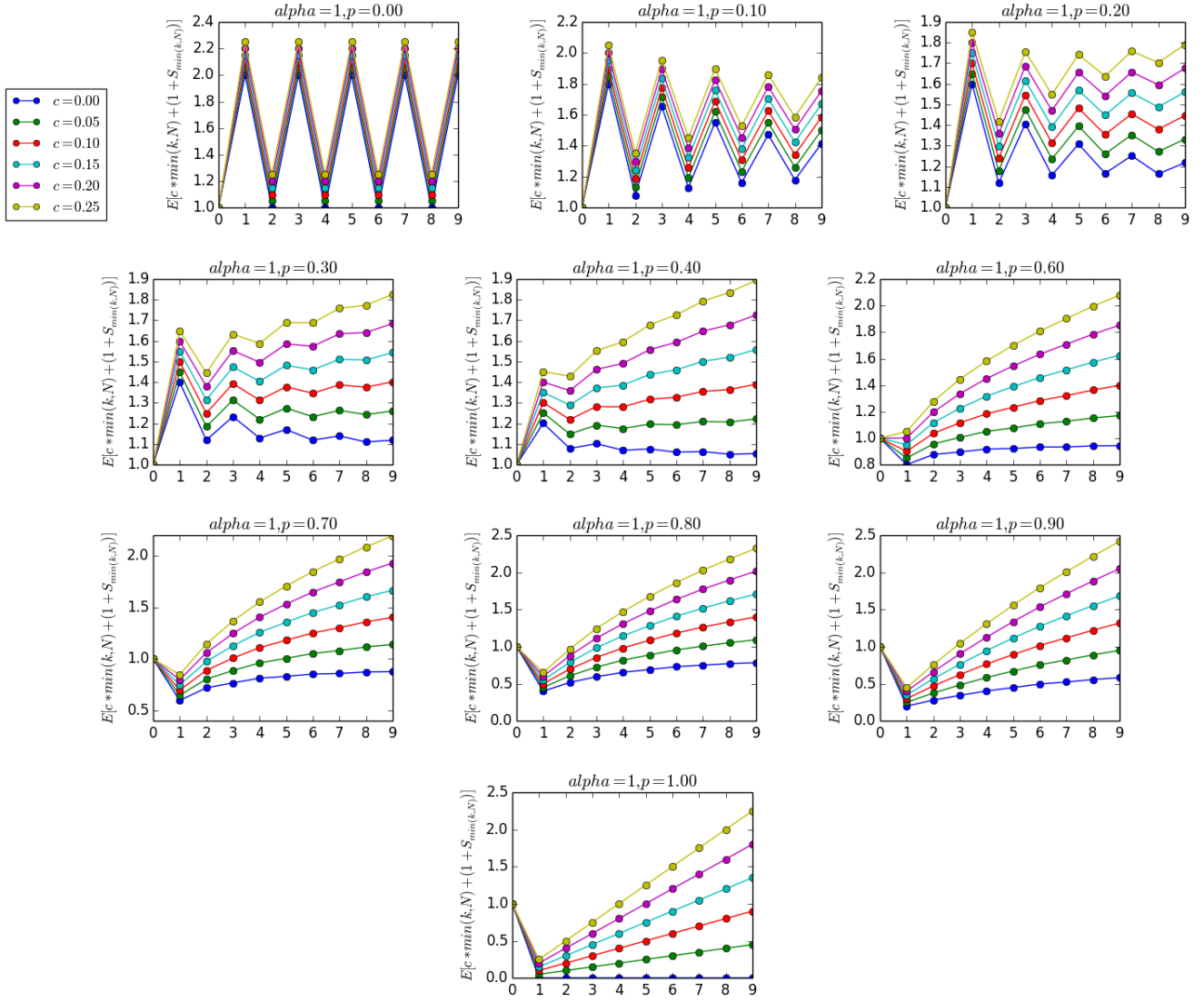
It is interesting to see that the behaviour changes when depending on whether $p < 1/2$ or not. Indeed, since $X_0 = 1$, for low values of p , there is a big chance that the trend changes:

$$p = \mathbb{P}(X_1 = 1 | X_0 = 1)$$

and therefore it is more profitable to buy directly at the opposite side.

The critical case when $p = 0.5$ is interesting since it is the case where the $(X_i)_i$ are independent. The random walk $(S_n)_n$ is a martingale whose expectation is constant. When $c = 0$, all strategies are equivalent, and as soon as $c > 0$, the additional cost of waiting makes it better to place a market order.

When $X_0 = -1$, ie $\alpha = 1$



3.2 Results from market data and validation

All experiments are done using the historical data of the Bund. Let's first try to estimate some parameters of the model.

$$p = \mathbb{P}(X_{i+1} = X_i) \approx \frac{\sum_{i \leq M} 1_{\{X_{i+1}=X_i\}}}{M}$$

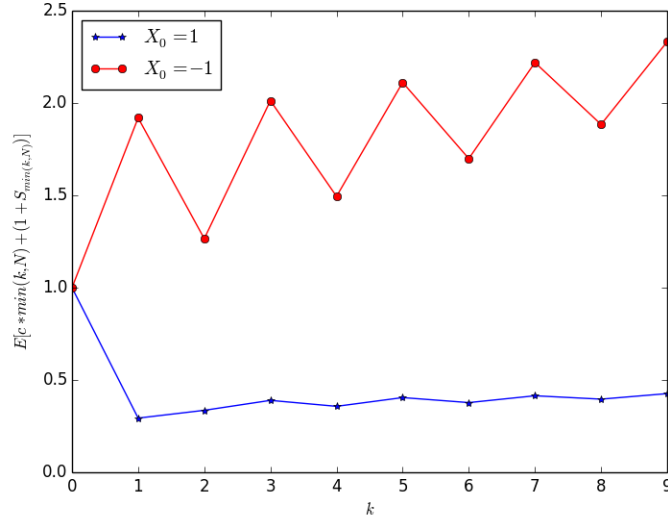
The following program gives $p \approx 0.1$

```
r=pickle.load(input)
X = r['BidPrice']
X = X[1:] - X[:-1]
X = X[ abs(X) > 1e-3] # Keep only the entries where the price changes
X = np.array(map(lambda x: 1 if x > 0 else -1, X))
p = np.count_nonzero(X[1:] == X[:-1]) / len(X)
```

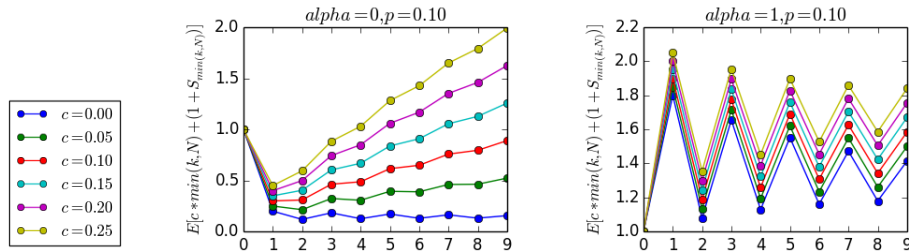
The same algorithm gives

$$\mathbb{E}[\tau] \approx 3.75ms$$

Here is the mean cost when applying Strategy(k) for $k \in \{0, \dots, 9\}$ periods using $c = 0.3$:



The shape of the graph is very similar to the ones we found in the previous section corresponding to $p = 0.1$



4 Model limitations and possible improvements

Our model presents three main limitations which can be summarized by the following :

- $\mathbb{E}(S_{k \wedge N})$ has been computed through a Monte Carlo Algorithm. While we are able to compute efficiently the law of S_i ⁷, we were not able to find a simple analytical expression. As a matter of fact, S_i and N are dependent and the calculus becomes quickly hard to do.
- Our model is markovian in the sense that we suppose that each move in price depends only on the previous move in price. This choice has been made to reduce the complexity of the calculus in particular for the law of N . This choice is also induced by a previous hypothesis that says that the state of the order book at the best bid and best ask level is randomly and independently drawn whenever the price moves. To improve the model at this point, we could consider other levels of the bid and ask queues. A more comprehensive, and thus more complex mode, would take in consideration the effect of previous trades on the intensity of arrival of orders. Hawkes processes can be considered for this.
- In our strategy, the maximum number of moves in price we are willing to wait for is bounded by 10. Whereas it could be interesting to wait more. We have made this arbitrary choice just to make a proof of concept. We could go much further with values of k above 100. The calculus remains quickly feasible.

⁷see Annexe B

A Technical proofs

A.1 Reccurence equation of $(g_{-1}^s)_s$

The purpose of this section is to find the equation verified by $(g_{-1}^s)_s$. Equations (5)-(6) implies that:

$$\begin{aligned} g_{s+1}^1 &= \frac{1}{zp}g_s^1 - \frac{q}{p}g_{s-1}^{-1} \\ &= \frac{1}{zq}g_s^{-1} - \frac{p}{q}g_{s-1}^{-1} \end{aligned}$$

Thereby : $g_s^{-1} = \frac{q}{p}g_s^1 + \frac{z(p^2-q^2)}{p}g_{s-1}^{-1}$.

Leading finally to: $g_{s+1}^{-1} = (\frac{1}{zp} + z(2 - \frac{1}{p}))g_s^{-1} - g_{s-1}^{-1}$.

A.2 Taylor expansion of g_{-1}^s

Recall that: $g_{s+1}^{-1} = (\frac{1}{zp} + z(2 - \frac{1}{p}))g_s^{-1} - g_{s-1}^{-1}$ Let us do the following denotations :

$$\begin{cases} u = \frac{1}{2}(\frac{1}{zp} + z(2 - \frac{1}{p})) = az + \frac{b}{z} & (14) \\ a = 1 - \frac{1}{2p}; b(p) = \frac{1}{2p} & (15) \end{cases}$$

$$(16)$$

Using a generalized Newton's formula we get :

$$\begin{aligned} (u^2 - 1)^{\frac{1}{2}} &= \sum_i \binom{\frac{1}{2}}{i} (-1)^i z^{2i-1} (b + az^2)^{1-2i} \\ &= \sum_i \binom{\frac{1}{2}}{i} (-1)^i z^{2i-1} \sum_j \binom{1-2i}{j} b^{1-2i-j} a^j z^{2j} \\ &= \sum_{i,j} \binom{\frac{1}{2}}{i} \binom{1-2i}{j} (-1)^i b^{1-2i-j} a^j z^{2(i+j)-1} \\ &= \sum_k \left(\sum_{2(i+j)-1=k} \binom{\frac{1}{2}}{i} \binom{1-2i}{j} (-1)^i b^{1-2i-j} a^j \right) z^k \\ &= \sum_r \left(\sum_{i+j=r} \binom{\frac{1}{2}}{i} \binom{1-2i}{j} (-1)^i b^{1-2i-j} a^j \right) z^{2r-1} \\ &= \sum_r \left(\sum_{i+j=r} \binom{\frac{1}{2}}{i} \binom{1-2i}{j} (-1)^i b^{1-2i-j} (1-b)^j \right) z^{2r-1} \\ &= \sum_r c_{2r-1} z^{2r-1} \end{aligned}$$

On the one hand, for even k :

$$c_k = 0$$

On the other hand, for odd $k \geq -1$ we get :

$$c_k = \sum_{2(i+j)=k+1} \binom{\frac{1}{2}}{i} \binom{1-2i}{j} (-1)^i b^{1-2i-j} (1-b)^j$$

$$c_{-1} = \binom{\frac{1}{2}}{0} \binom{1}{0} (-1)^0 b^1 (1-b)^0 = b$$

$$c_1 = \sum_{i+j=1} \binom{\frac{1}{2}}{i} \binom{1-2i}{j} (-1)^i b^{1-2i-j} (1-b)^j = a - p$$

Thereby

$$\begin{aligned} g_1^0(z) &= u - (u^2 - 1)^{\frac{1}{2}} \\ &= (a - c_1)z - (b + c_{-1})z^{-1} - \sum_{k>1} c_k z^k \\ &= pz - \sum_{k>1} c_k z^k \\ &= \sum_{k \geq 1} d_k z^k \end{aligned}$$

Where

$$d_{2r+1} = \begin{cases} p & \text{if } r = 0 \\ -\sum_{i \leq r+1} \binom{\frac{1}{2}}{i} \binom{1-2i}{1+r-i} (-1)^i (2p)^{2i-1} (2p-1)^{1+r-i} & \text{else} \end{cases}$$

B The law of a sum of dependent random variables

As previously, we define a Markov Chain $(X_n)_{n \geq 0}$ and its sum $(S_n)_{n \geq 0}$. In order to get the law of the r.v $(S_n)_{n \geq 0}$ we define its conditional probability-generating function as :

$$f_n^x(z) = \mathbb{E}(z^{S_n} | X_0 = x) \quad \forall x \in \{-1, 1\}, n \in \mathbb{N}, z \in [0, 1]$$

Let us consider $x \in \{-1, 1\}$, $n \in \mathbb{N}$, $z \in]0, 1]$ ⁸. Thanks to the fact that (S_n, X_n) is a markov chain, we can write :

$$f_{n+1}^x(z) = z \mathbb{E}(z^{S_n} | X_0 = 1) \mathbb{P}(X_1 = 1 | X_0 = x) + \frac{1}{z} \mathbb{E}(z^{S_n} | X_0 = -1) \mathbb{P}(X_1 = -1 | X_0 = x) \quad (17)$$

Leading thus to the following system with the appropriate initial conditions :

$$\begin{cases} f_{n+1}^1(z) = zp f_n^1(z) + \frac{1}{z} q f_n^{-1}(z) \end{cases} \quad (18)$$

$$\begin{cases} f_{n+1}^{-1}(z) = zq f_n^1(z) + \frac{1}{z} p f_n^{-1}(z) \end{cases} \quad (19)$$

$$\begin{cases} f_0^1(z) = f_0^{-1}(z) = 1 \end{cases} \quad (20)$$

$$\begin{cases} f_1^1(z) = zp + \frac{1}{z} q \end{cases} \quad (21)$$

$$\begin{cases} f_1^{-1}(z) = zq + \frac{1}{z} p \end{cases} \quad (22)$$

We then prove that f_n^1 and f_n^{-1} follow the same equation, ie for $x \in \{-1, 1\}$

$$f_{n+2}^x - p\left(\frac{1}{z} + z\right)f_{n+1}^x + (2p-1)f_n^x = 0 \quad (23)$$

Which is a classical linear recurrent sequence of order 2. The solution is given by :

$$f_n^x(z) = \frac{f_1^x - v + \sqrt{u^2 - (2p-1)}}{2\sqrt{v^2 - (2p-1)}} \left((v + \sqrt{v^2 - (2p-1)})^n - (v - \sqrt{v^2 - (2p-1)})^n \right) + (v - \sqrt{v^2 - (2p-1)})^n \quad (24)$$

⁸The case where $z=0$ can be obtained by continuity.

Where :

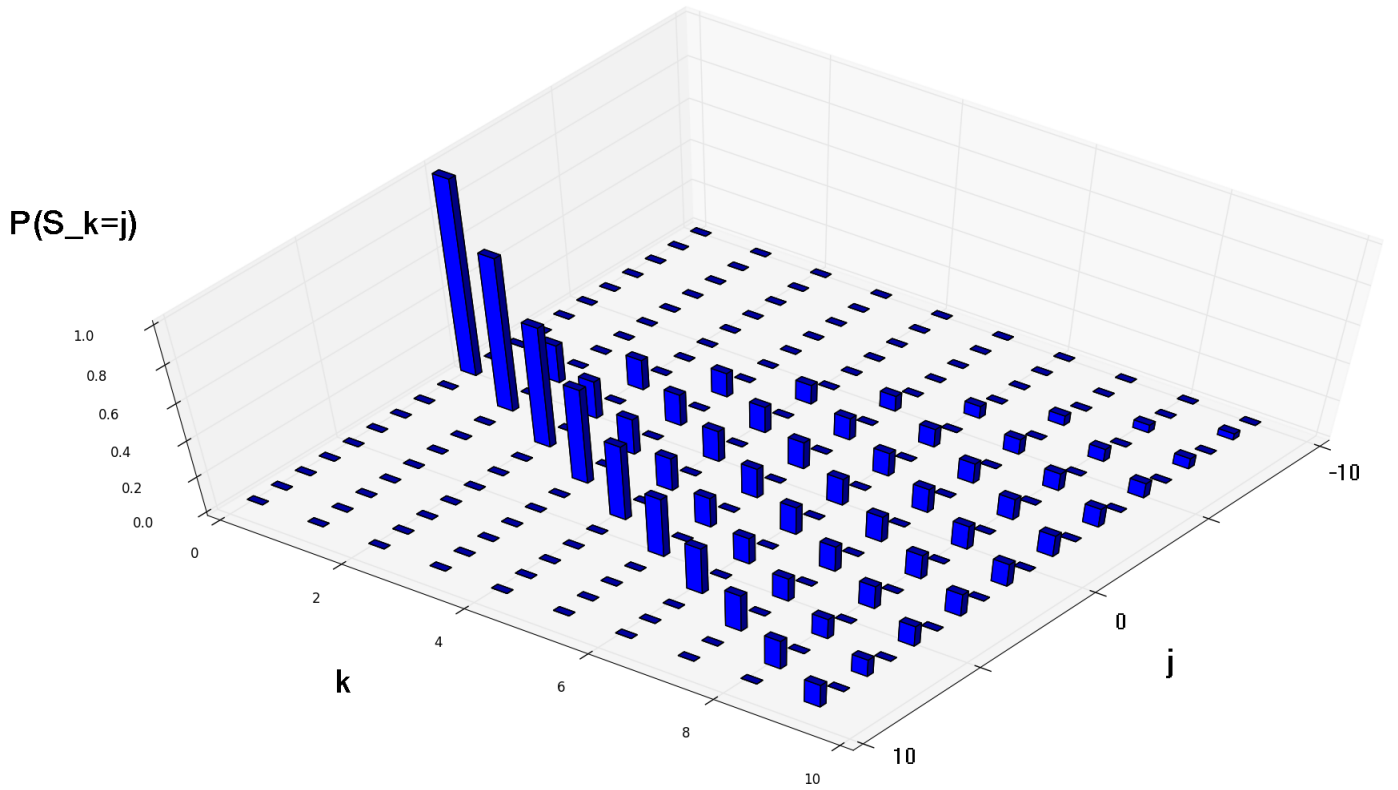
$$v = \frac{p}{2} \left(z + \frac{1}{z} \right)$$

f_n^x is the probability-generating function of $(S_n)_{n \geq 0}$. We shall notice that $\forall n \in \mathbb{N} - n \leq S_n \leq n$. We then get :

$$\left\{ \begin{array}{l} \mathbb{E}_{X_0=x}(S_n) = \frac{df_n^x}{dz}(z=1) \end{array} \right. \quad (25)$$

$$\left\{ \begin{array}{l} f_n^x(z) = \mathbb{E}(z^{S_n} | X_0 = x) = \sum_{-\infty}^{+\infty} \mathbb{P}(S_n = k) z^k = \sum_{-n}^{+n} \mathbb{P}(S_n = k) z^k \end{array} \right. \quad (26)$$

We use Mapple to compute the taylor expansion, we get for example for $X_0 = -1$ and $p = 0.2$:



We can see that $\mathbb{P}(S_0 = 0) = 1$ and $\mathbb{P}(S_1 = 1) = 1 - \mathbb{P}(S_1 = -1) = 0.2 = 1$

C Source Code