

# What is a good price model for high frequency data ?

Mathieu Rosenbaum

University Paris 6

February 2013

## High frequency data

- Order book.
- Different prices : last traded, mid-quote, best bid, best ask, Volume Weighted Average Price, . . .
- Durations.

# Example of order book



Carnet d'ordre ML					
<b>MICHELIN</b> 3			Heure	Prix 2	volume
FR0000121261 (ML)			13:15:05	83.65	50
<b>Dernier</b> 83,65			13:14:50	83.65	100
Var (%) : +4,56%			13:14:32	83.65	79
Var (pts) : 3,65			13:14:32	83.65	11
Ouvrir 80,15			13:14:30	83.65	49
Plus haut 85,1			13:14:30	83.65	3951
Plus bas 80,15			13:13:18	83.65	520
volume 3 023 872			13:13:18	83.65	507
			13:13:18	83.65	680
			13:13:18	83.65	232
<b>Demandes</b> (Acheteurs - bid) 1			<b>Offres</b> (Vendeurs - ask)		
Nb. ordres	Quantités	Cours	Cours	Quantités	Nb. ordres
1	4 771	83.65	83.70	365	1
1	225	83.60	83.75	4 135	3
3	2 810	83.50	83.80	1 105	12
2	1 946	83.45	83.85	1 275	6
4	5 955	83.40	83.90	339	5
<b>11</b>	<b>15 707</b>			<b>7 219</b>	<b>27</b>

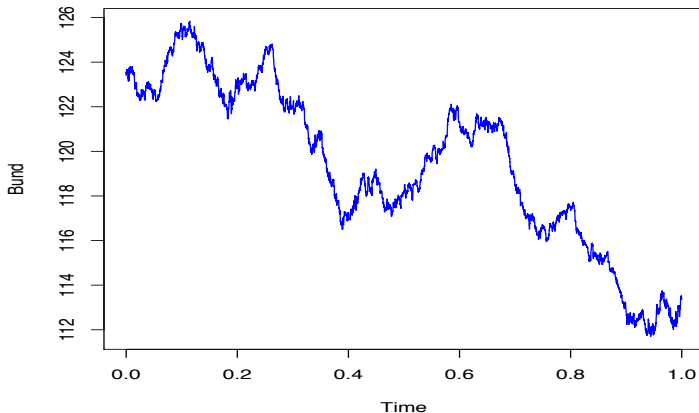
## Pricing-Hedging

- When pricing and hedging derivatives, the need for probabilistic models is clear.
- Their use is necessary in order to be able to manage the risks inherent to options trading.
- Assessing the quality of a model is quite an easy task.
- A model is suitable if it reproduces the behavior of important market quantities for the risk manager (smile, dynamic of the smile, etc.) and if it is computationally not too intricate so that it can be used in practice for pricing and hedging.
- In the derivatives context, the time scale of work is essentially macroscopic, in the sense that portfolio are usually rebalanced once or perhaps a few times within a day. Thus, the need for models at the macroscopic scale is obvious.

## Prices and Brownian type models

- At the macroscopic scale, for example if we consider one data per day over one year, prices trajectories look like sample paths of models classically used in mathematical finance such as Brownian motion, diffusions, stochastic volatility models,...

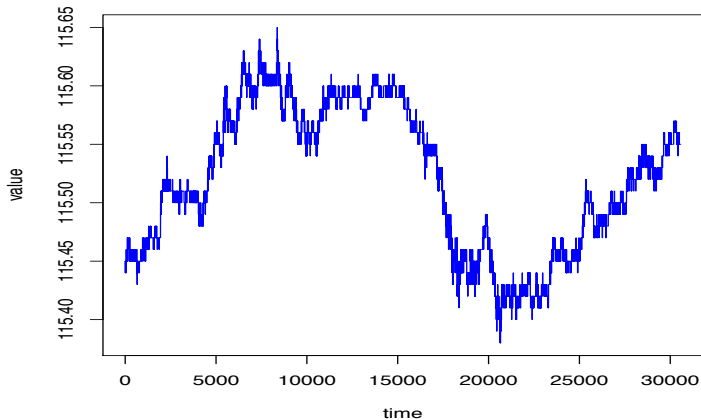
# Bund contract, 2005-09-01 to 2007-01-31, one data every hour



## Microstructure effects

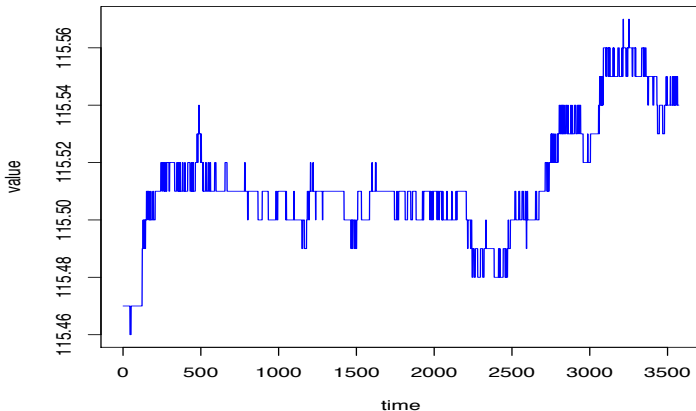
- At the microscopic scale, prices are very different from Brownian type sample paths.
- This phenomenon is called microstructure effect.
- Therefore, it is quite intricate to model such data.

# Bund, 2007-05-06, whole day, one data every second





Bund, 2007-05-06, from 10h am to 11h am, one data every second



## Relevance of high frequency models

- When considering high frequency data, the question of the relevance of the models is not trivial.
- In the economic approach to microstructure modeling, prices are explained through the behavior of market participants interacting together.
- Therefore, the scientific aim of the models is clear : trying to understand the behavior of market participants.
- Then, in this setting, a model is a good model if the assumptions on the market participants are reasonable, leading to acceptable market features and mostly if it enables to draw relevant conclusions or recommendations (for example to the regulator).

## Relevance of high frequency models

- In the case of probabilistic/statistical models, not deeply motivated by some agents behavior, what do one really want for a model ?
- Essentially, a model is good if it reasonably reproduces the stylized facts of the main quantities of interest and, mostly, if it is useful for market practitioners.

## Our approach

- We distinguish four levels of resolution for modeling “through scales”. A good probabilistic/statistical model should provide reasonable dynamics across these levels.

## Different scales

- L0** The ultimate level of the order book. One proposes a complex stochastic system in continuous time with discrete values in an appropriate state space which describes all the events of a limit order book.
- L1** The ultra high frequency level for the price. At this level, one wishes to model all transaction prices and durations between these transactions.
- L2** At an intermediate high frequency level. Here one does not focus on durations but essentially on the price, regularly sampled, for example every second or every minute.
- L3** The macroscopic level, where the price is viewed as a continuous semi-martingale. It is the dominant and historical approach.

## About these levels

- Of course the most reasonable approach is to obtain a valid model for all levels  $L_0$ ,  $L_1$ ,  $L_2$ ,  $L_3$ . This is also the most difficult approach because it requires to model all the complexity of the order book. We will essentially focus on price models here and so only consider levels  $L_1$ ,  $L_2$ ,  $L_3$ .
- Most models focus on a stability between levels  $L_2$  and  $L_3$  only. When it comes to practical use, the difficulty of such approaches is to “define” level  $L_2$ . Indeed, one has to choose an exogenous sampling scale. Thus, the question is : when is the model valid ? For 5 minutes sampling ? 1 minute sampling ? 1 second sampling ?
- Another issue is the question of the modeled price. Is it the last traded ? mid-quote ? best bid ? best ask ? Volume Weighted Average Price ?...

## Semi-martingales

- Since Black and Scholes, continuous time processes are standard modeling objects in mathematical finance. Among these processes, continuous semi-martingales play a major role.
- Indeed, it is well known that the no free lunch assumption is essentially only compatible with semi-martingale type dynamics for the price (Delbaen and Schachermayer).
- Since the middle of the nineties and the availability of a huge amount of very high frequency data, various empirical study have shown that over short time periods (order of magnitude of an hour or a day), it is not reasonable to see price observations as observations of a continuous semi-martingale.

## Signature plot

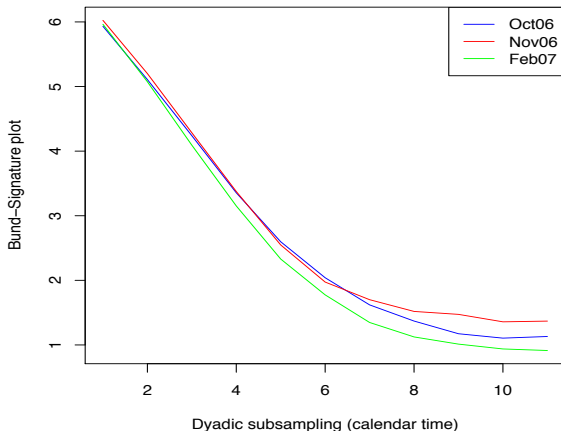
- In the analysis of empirical data, a usual tool is the “signature plot”. Assume we have price observations  $P_t$  at times  $t = i/n$ ,  $n \in \mathbb{N}$ ,  $i = 0, \dots, n$ , where  $t = 1$  represents for example one trading day. The signature plot is the function which to  $k = 1, \dots, n$  associates

$$RV_n(k) = \sum_{i=0}^{\lfloor n/k \rfloor - 1} (P_{k(i+1)/n} - P_{ki/n})^2.$$

- If  $P_t$  is a continuous semi-martingale, as soon as  $(n/k)$  is large enough,  $RV_n(k)$  is close to the quadratic variation of the semi-martingale. However, in practice,  $RV_n(k)$  is very often a decreasing functional, stabilizing for sampling periods larger than 10 minutes (depending on the asset of course).



# From the large scales to the fine scales



**FIGURE:** *Signature plot for Bund contract, one data every second, aggregation of all the trading days in each month.*

# From the large scales to the fine scales

## High frequency vs low frequency

- It is clear that high frequency prices time series are not of the same nature as low frequency series.
- Said differently, the scale invariance assumption associated to Brownian type dynamics is not reasonable above some sampling frequency.
- The goal of the coarse to fine approach to microstructure is the following : reconcile these different behaviors across scales starting from the coarse scale, that is the scale where a continuous semi-martingale type dynamics is a relevant model.

# From the large scales to the fine scales

## High frequency vs low frequency

- More precisely, one starts from a continuous semi-martingale, called efficient (latent) price, and apply a stochastic mechanism to it in order to derive the observed prices at higher frequencies.
- This mechanism must enable to obtain a suitable dynamics for the observed prices in the high frequencies. Also, its effect has to be negligible in the low frequencies, so that the observed price at the macroscopic scale remains close to the efficient price.

# From the large scales to the fine scales

## About the efficient price

- One specific element of the coarse to fine approach is the presence of the continuous semi-martingale.
- Of course, it is highly arguable to think that the “economic” notion of efficient price is meaningful in the very high frequencies.
- However :

# From the large scales to the fine scales

## About the efficient price

- It is clear that this efficient price, even though not well interpreted is meaningful. Indeed, one should not forget that a price is associated to an economic quantity (share, rate, currency, . . . ) and there is some kind of market agreement on the value of this quantity.
- One can be fully convinced that this efficient price exists looking at how prices can strongly react after an economic news.
- Also, it is important to understand that the use of a continuous semi-martingale at very high frequencies (typically the  $L1$  scale) does not mean that one believes that the notion of efficient price makes sense at these very fine scales. It can also only be seen as a probabilistic tool used in order to build the model.

# From the large scales to the fine scales

## Usefulness of the coarse to fine approach

- The advantage of the presence of this efficient price is that common quantities, meaningful for market practitioners, are perfectly well defined.
- For example, the notions of volatility or covariation are fully understood for continuous semi-martingales. Therefore, the coarse to fine approach enables to transfer these very usual notions at the high frequency level and to build statistical procedures for estimating them.

## Usefulness of the coarse to fine approach

- This is of first importance. Indeed, when the trading horizon is short, one still needs such concepts as volatility or covariation at the high frequency level, for example for optimal trading problems.
- Remark that the coarse to fine approach is also useful when one is only interested in these quantities at a low frequency level. Indeed, it enables to use high frequency data in order to estimate low frequency parameters of interest, which implies a much better accuracy.

## Additive microstructure noise models

- In the coarse to fine approach, one can consider the widely used notion of microstructure noise. This noise is defined for all  $t$  where the model price exists by

$$\varepsilon_t = \log P_t - \log X_t,$$

where  $P_t$  is the model price at time  $t$  and  $X_t$  is the semi-martingale used in order to build the model. In the so called additive microstructure noise approach, one focuses on modeling the process  $\varepsilon_t$ .



## Additive microstructure noise models

- More precisely, in such model one writes the log price  $P_t$  observed at time  $t = i/n$  as

$$\log P_{i/n} = \log X_{i/n} + \varepsilon_i^n,$$

where  $X$  is a continuous semi-martingale playing the role of efficient price and  $\varepsilon_i^n$  is the microstructure noise term. Such model is said to be an additive microstructure noise model if the noise is (essentially) centered, stationary and independent of  $X$ .

## Additive microstructure noise models

This type of approach is probably the simplest way to obtain

- data close to observations of a semi-martingale in the low frequencies :  $\log P_{i/n} - \log P_{j/n} \approx \log X_{i/n} - \log X_{j/n}$  if  $(i - j)/n$  is big enough,
- data very different from observations of a semi-martingale in the high frequencies, the noise term becoming predominant.

## Additive microstructure noise models

Let us make two important remarks :

- This approach is typically a  $L3 \rightarrow L2$  approach, which implies all the drawbacks and difficulties of use of a  $L3 \rightarrow L2$  model.
- In such models, one puts assumptions on the microstructure noise. Then, from this noise and the efficient price, the observed price is deduced.
- The drawback of such approach is that it models an unobservable quantity (the noise). A structural approach would be to directly model the observed price and possibly deduce the microstructure noise from the observed price and the efficient price.

## No $L2 \rightarrow L1$ stability in additive microstructure noise models

Some simple considerations enable to deduce some suitable properties of the observed price and the microstructure noise.

- (i) Observed prices are discrete. Indeed, market prices usually stay on a tick grid. This is not compatible with data from a continuous semi-martingale observed at exogenous times. Therefore, in this approach, this is a “reason” for the presence of microstructure noise.
- (ii) For many assets, there are quick oscillations of the transaction prices between two values (“bid-ask bounce”).

## No $L2 \rightarrow L1$ stability in additive microstructure noise models

Since  $(P_t)_{t \geq 0}$  is a jump process so that the number of jumps is finite over every bounded time interval, we deduce that

$$[\log P, \log X]_t = [\log X]_t + [\varepsilon, \log X]_t = 0$$

and

$$[\varepsilon]_t = [\log P]_t - [\log X]_t - 2[\varepsilon, \log X]_t = [\log P]_t + [\log X]_t.$$

Consequently,

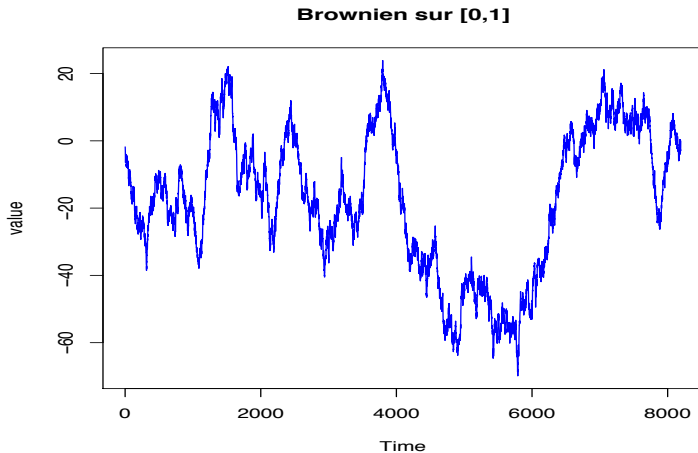
(iii)  $[\varepsilon]_t$  is almost surely finite.

### No $L2 \rightarrow L1$ stability in additive microstructure noise models

- The additive microstructure noise approach, which is the standard approach for modeling microstructure noise, considers a noise which is essentially iid, independent of the efficient price. As explained later, these models have been largely used in order to build volatility estimation procedures.
- These models have been extended in several directions : heteroskedastic noise correlated with the efficient price or contamination of the efficient price through a complex Markovian kernel, . . .

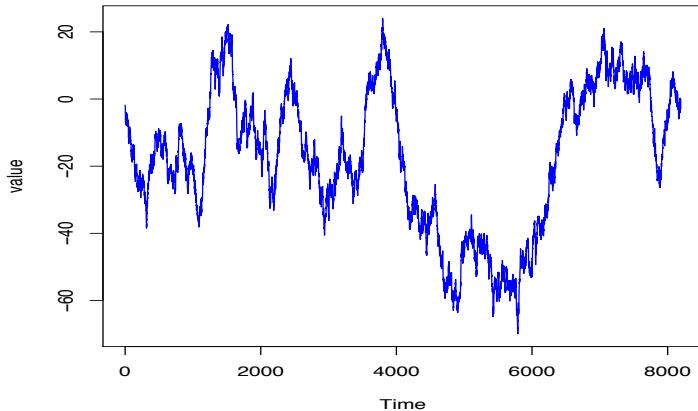
### No $L2 \rightarrow L1$ stability in additive microstructure noise models

- Additive microstructure noise models are very convenient to carry on computations, see later, and are often reasonable when considering sampling scales larger than about 5 minutes.
- However, they do not satisfy any of the properties (i), (ii), (iii) and the durations are not modeled. Consequently, they cannot be extended to level  $L1$ .
- Indeed, in the high frequencies, discretization and seasonality effects lead to very strong and intricate non linear dependence between the microstructure noise and the efficient price together, with a complex heteroskedasticity of the noise.

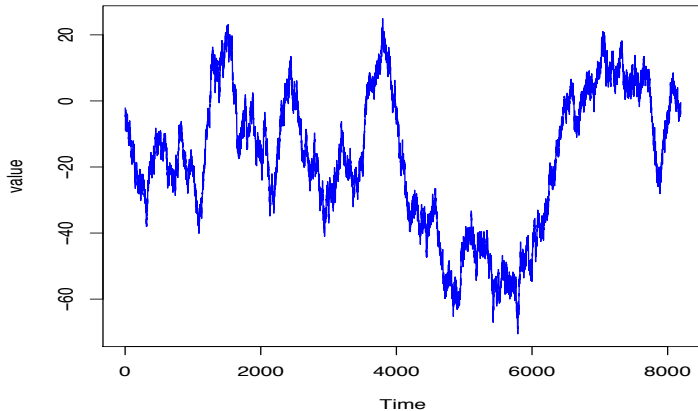




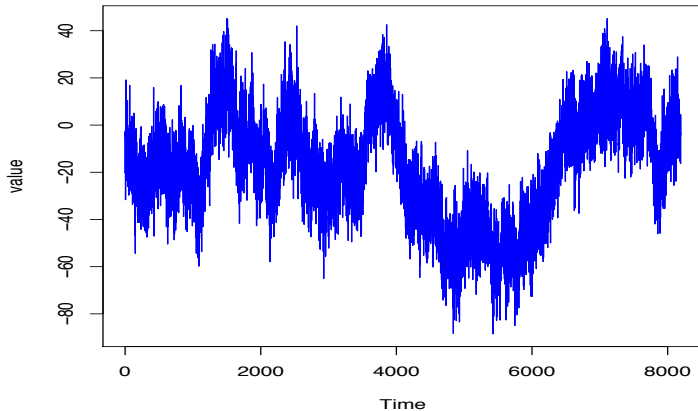
**Brownien+0.1\*gaussienne standard sur [0,1]**



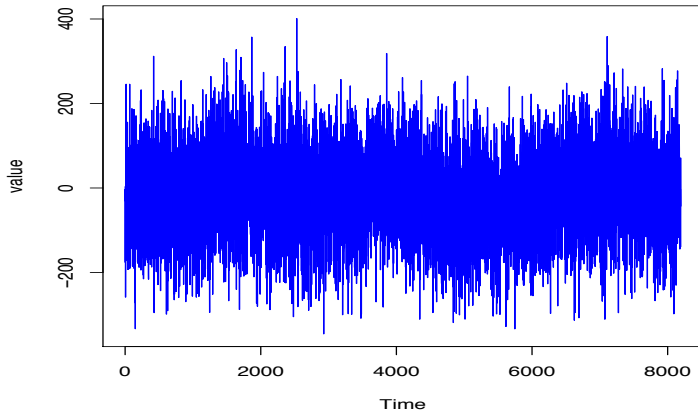
**Brownien+1\*gaussienne standard sur [0,1]**



**Brownien+10\*gaussienne standard sur [0,1]**



**Brownien+100\*gaussienne standard sur [0,1]**



## Rounding models

- Rounding models are a simple way to accommodate properties (i), (ii), (iii) and the assumption of an underlying semi-martingale efficient price.
- Here the idea is not to focus on the properties of the microstructure noise but directly on the properties of the observed price.
- Thus, it is very natural to consider the model of continuous semi-martingale observed with rounding error introduced by Delattre and Jacod.

## Rounding models

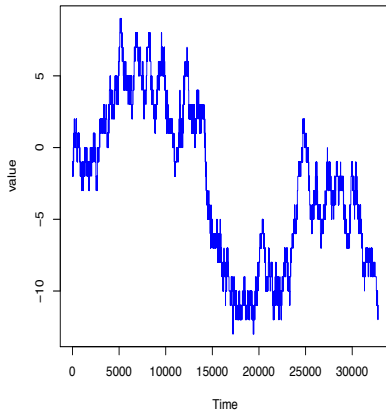
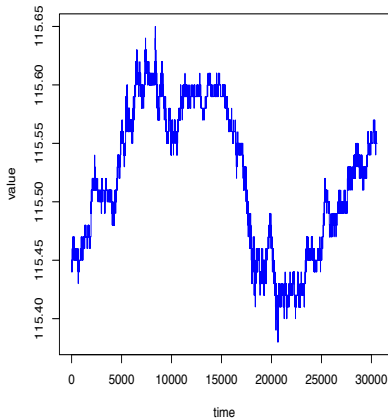
- In this model, the efficient price is modeled by a continuous semi-martingale  $X_t$  and the observed prices are given by the sample

$$(X_{i/n}^{(\alpha_n)}, i = 0, \dots, n),$$

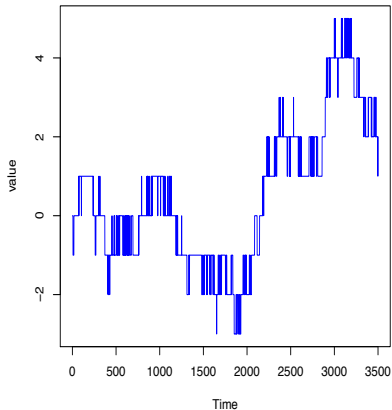
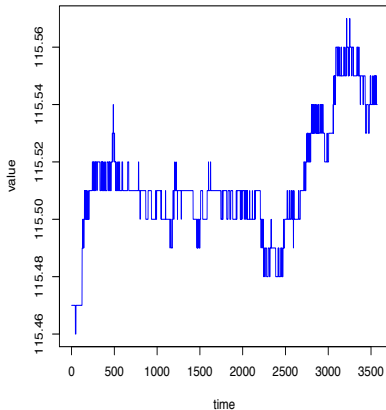
where  $X_{i/n}^{(\alpha_n)} = \alpha_n \lfloor X_{i/n} / \alpha_n \rfloor$ .

- Therefore,  $X_{i/n}^{(\alpha_n)}$  is the observation of the efficient price at time  $i/n$ , with rounding error  $\alpha_n$ , corresponding to the tick size.
- In this context, prices are discrete and behave as observations from a continuous semi-martingale in the low frequencies, the rounding effect becoming negligible. Furthermore, Properties (ii) and (iii) can be satisfied.

Bund, 2007-05-06, one data every second and  $\lfloor 20W_t \rfloor$ , on  $[0, 1]$ ,  $n = 2^{15}$



Bund, 2007-05-06, from 10h am to 11h am, one data every second and  $\lfloor 20W_t \rfloor$ , on  $[0, 0.1]$





## Rounding models

- Hence, compare to additive microstructure noise models, rounding models have several nice properties.
- However, the main drawback of these models is that, as additive microstructure noise models, they remain  $L3 \rightarrow L2$  models and cannot be extended to level  $L1$  : Assume that for any time  $t$  the observed price is given by the rounding value of a semi-martingale. This leads to an observed price with an infinite number of jumps on a finite interval, which is of course hardly acceptable. Therefore, the already mentioned drawbacks of  $L3 \rightarrow L2$  models apply in the rounding case.

## Suitable properties

It appears that several drawbacks are inherent to  $L3 \rightarrow L2$  models. Let us summarize the suitable properties for a  $L3 \rightarrow L1$  model :

- A continuous semi-martingale type behavior at large sampling scales.
- Model for prices and durations (in particular, no notion of sampling frequency is required).
- A clear definition of the price.
- Discrete prices.
- Bid-Ask bounce.

## Suitable properties

- Usual stylized facts of returns, durations and volatility (in a loose sense here). In particular, inverse relation between durations and volatility.
- An interpretation of the model.
- Finite quadratic variation for the microstructure noise.
- A testable model.
- A useful model, for example for building statistical procedures.

## A model

- We will describe later a model which satisfies the preceding properties.

What about using discontinuous semi-martingale in coarse to fine models ?

- Until now, as in most models of mathematical finance, we have always assumed that the efficient price was a continuous semi-martingale. This is of course possible to model it by a discontinuous semi-martingale.
- In this case, the observed price is for example, in the additive microstructure noise case, the sum of this discontinuous semi-martingale plus a noise term or, in the rounding case, the rounded value of this semi-martingale
- However, the difficulty of this approach is the identifiability question.

What about using discontinuous semi-martingale in coarse to fine models?

- The reason is the following : if the semi-martingale can be very general, it becomes not so clear why one should modify it in order to obtain the observed price.
- Indeed, if one considers for example the last traded price, this is a process defined in continuous time which is of finite variation. Therefore, it can be seen as the sample path of a semi-martingale and the decreasing signature plot can be considered as obtained on a (specific) semi-martingale.
- Thus, the need for noise in this case is not so clear, unless assuming some restrictive conditions on the jump part of the process. That is why in this paradigm, we consider the efficient price as a continuous semi-martingale.

## $L1 \rightarrow L3$ models

- As explained above, even though it reproduces some stylized facts, a model has no interest if it is not useful.
- Contrary to the coarse to fine models and their use for statistical estimation, there is no generic application of the fine to coarse models.
- Hence, the usefulness of the modeling depends on each model.
- One drawback of the majority of these models is that they are built without any reference to the some kind of efficient price. However, one must not forget that a price still symbolizes the value of an economical quantity (price discovery through the order book vs price formation).

## Is one more natural?

- Fine to coarse models seem more satisfying since they start from the small scales to arrive to the large scales. This is probably more natural.
- In particular an  $L0 \rightarrow L3$  approach is particularly satisfying since the large scale price is derived from the behavior of the order book.



## Is one more natural ?

- However, let us stress that  $L1 \rightarrow L3$  models are essentially phenomenological, in the sense that they use suitable stochastic processes in order to accurately reproduce the behavior of the data.
- Hence, it does not make sense to compare the fine to coarse approach  $L1 \rightarrow L3$  and the coarse to fine approach  $L3 \rightarrow L1$  in terms of “more natural point of view”.
- Indeed, an  $L3 \rightarrow L1$  model is in fact a model which is more easily defined thanks to an efficient price process and a stochastic mechanism applied to it than directly at level  $L1$ , whereas the contrary holds for an  $L1 \rightarrow L3$  model.