

Statistical Arbitrage Trading Strategy

Backtesting Using KDB+/Q

Kim Kuen Tang

19. April 2012

Abstract

This talk is about the pairs trading strategy under the framework of stochastic control approach. We formulate the problem as the maximum of an expected terminal wealth under an utility function for a given risk parameter. This value function is solved using a separation ansatz which will help us to obtain a closed-form solution of the optimal positions. This strategy is then backtested completely in KDB+/Q. Numerical results are shown for the pnl, maximum drawdown and sharpe ratio.

Outline

What is a Trading Strategy?

Data-Driven Alpha Model

- Market Data

- Parameter Specification

Theory-Driven Alpha Model

- Introduction

- Model Specification

- Problem Formulation

- Closed Form Solution

Backtesting with KDB+/Q

- normal distributed random numbers

- Stochastic Differential Equations

- terminal pnl, maximum drawdown and sharpe ratio

What is a trading strategy?

Mathematician William F. Donoghue: “Thorp, my advice is to buy low sell high. “

Three steps for a successful alpha model:

step	skill
idea	visionary
development and backtesting	quantitative and programming
successful real world implementation	entrepreneurial

This talk is about the idea, development and backtesting.

Orderbook and Trade

- ▶ An order book is the list of orders a trading venue (in particular stock exchanges) uses to record the interest of buyers and sellers.
- ▶ The trade data contains the last sale, last size and last volume of the trade.
- ▶ Use metrics of the orderbooks and trades to predict the outcomes of the returns from midprices or trades.

We are interested in the following equation

$$h_k(r_{k+1}) = f_k(\bar{o}_k, \bar{t}_k)$$

Here r_{k+1} denotes the returns, (\bar{o}_k, \bar{t}_k) contains all the information of orderbooks and trades until k .

At least two ways to choose h_k and f_k

- ▶ Regression

$$r_{k+1} = \beta^T \cdot f(\bar{o}_k, \bar{t}_k) + \epsilon_{k+1}$$

- ▶ Support Vector Machine

$$\text{sign}(r_{k+1}) = \text{sign}(\beta^T \cdot f(\bar{o}_k, \bar{t}_k) + \beta_0)$$

This talk is not about these models. For completeness we just mentioned it here but we will not show any results.

Pairs Trading

- ▶ Identify a pair of stocks where the log prices is cointegrated
- ▶ Model the spread of the log prices using an Ornstein-Uhlenbeck process.
- ▶ If observations are larger (smaller) than the predicted value we take a long (short) position and unwind it when the spread reverts.

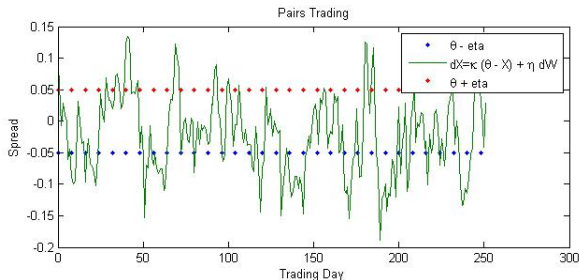


Figure: Example path for the Ornstein-Uhlenbeck process with $\theta = 0.0$

Model Specification

With $M = M_t$, $B = B_t$, $A = A_t$, $X = X_t$, $W = W_t$ and $Z = Z_t$ we describe the asset, spread and wealth dynamics by the following formulas:

- ▶ $dB/B = \mu dt + \sigma dZ$ (price dynamic)
- ▶ $X = \ln A - \ln B$ (spread definition)
- ▶ $dX = \kappa(\theta - X)dt + \eta dW$ (spread dynamic)
- ▶ $dZdW = \rho dt$ (correlation)
- ▶ $dM = rMdt$ (risk-free asset)

The investment behaviour is described by a process h which is progressively measurable. It models the portfolio weight for stock A and $-h$ for stock B . If V denotes the wealth dynamic of our self-financing portfolio then the dynamic is given by

$$dV/V = h dA/A - h dB/B + r dM/M$$

Observe that

$$A = B \exp(X) = f(B, X)$$

Using Ito's formula:

$$df = \dot{f} dt + (\nabla' f)(dB, dX) + 0.5(dB, dX)(\nabla^2 f)(dB, dX)$$

with

$$\dot{f} = 0; f_1 = A/B; f_2 = A; f_{11} = 0; f_{12} = A/B; f_{22} = A$$

we conclude that

$$\begin{aligned} dA/A - dB/B &= dX + dB/B \cdot dX + (dX)^2 \\ &= (\kappa(\theta - X) + \eta^2/2 + \rho\sigma\eta) dt + \eta dW \end{aligned}$$

Problem Formulation

Now we seek to find the optimal position h where it is able to solve the expression below

$$\sup_{h \in \alpha(t, T)} E^{t, v, x}[U(V(T))] \text{ for an utility function } U \quad (1)$$

s. t.

$$\begin{aligned} U(x) &= x^\gamma / \gamma \text{ for a } \gamma \leq -1 \\ dV/V &= h dA/A - h dB/B + r dM/M \end{aligned}$$

Here $\alpha(t, T)$ contains all progressively and admissible processes. If we denote with $G = G(t, v, x)$ the expression in (1) then it follows with the Bellman-Principle that

$$G(t, v, x) = \sup_{h \in \alpha(t, s)} E^{t, v, x}[G(s, V(s), X(s))]$$

The Ito-integral gives the following expression

$$0 = \sup_{h \in \alpha(t,s)} E^{t,v,x} \left[\int_t^s dG \right]$$

Divide the right hand side by $s - t$ and let s go to t to get

$$0 = \sup_{h \in \alpha(t,T)} \text{drift}(dG) = \sup_h h^2 \cdot \{*\} + h \cdot \{*\} + * \quad (2)$$

The first order condition implies

$$\begin{aligned} 0 &= 2 \cdot \hat{h} \cdot \{*\} + \{*\} \\ &= \hat{h} \eta^2 v \cdot G_{vv} + \eta^2 \cdot G_{vx} + (-\kappa(x - \theta) + 0.5 \cdot \eta^2 \rho \sigma \eta) \cdot G_v \end{aligned}$$

Assume that $G_{vv} < 0$, it yields

$$\hat{h} = -(\eta^2 \cdot G_{v,x} + (-\kappa(x - \theta) + 0.5 \cdot \eta^2 \rho \sigma \eta) \cdot G_v) / (\eta^2 v \cdot G_{v,v}) \quad (3)$$

Now plug back the above expression into (2) we get a differential equation which depends only on the value function G .

$$\begin{aligned} 0 = & \eta^2 G_t G_{v,v} - 0.5\eta^4 G_{v,x}^2 - 0.5b^2 G_v^2 - b\eta^2 G_v G_{v,x} \\ & + 0.5\eta^4 G_{v,v} G_{x,x} + r\eta^2 v G_v G_{v,v} - \kappa(x - \theta)\eta^2 G_x G_{v,v} \\ = & \mathfrak{L}G \end{aligned} \quad (4)$$

Here $b = -\kappa(x - \theta) + 0.5\eta^2 + \rho\sigma\eta$.

To obtain a closed form solution we will consider the following separation ansatz.

$$G(t, v, x) = f(t, x)v^\gamma$$

Here f denotes

$$f(t, x) = g(t) \exp(x\beta(t) + x^2\alpha(t))$$

Applying the Operator \mathfrak{L} on the expression $g(t) \exp(x\beta(t) + x^2\alpha(t))v^\gamma$ yields

$$0 = x^2\{*\} + x\{*\} + \{*\}$$

Setting the three coefficients to zero yields the following three differential equations

$$\dot{\alpha} = \{*\}\alpha^2 + \{*\}\alpha + \{*\}$$

$$0 = \mathfrak{B}(\alpha, \beta)$$

$$0 = \mathfrak{G}(\alpha, \beta, g)$$

The first equation is a Riccati equation and \mathfrak{B} and \mathfrak{G} are first order linear ordinary differential equations.

With $c = \sqrt{1 - \gamma}$ and $d = \exp(2\kappa(T - t)/c)$ we conclude that

$$\alpha(t) = \kappa(1 - c)/(2\gamma^2) \cdot \{1 + (2c)/(1 - c - (1 + c)d)\}$$

$$\beta(t) = \frac{1/(2\eta^2(1 - c - (1 + c)d))\gamma(1 - d)}{\{c(\eta^2 + 2\rho\sigma\eta)(1 - d) - (\eta^2 + 2\rho\sigma\eta + 2\kappa\theta)\}}$$

$$g(t) = \exp\left(-\int_t^T u(s)ds/(c^2\eta^2)\right)$$

where u satisfies a differential equation which is not of interest here.

We need to show that the expression $G = g \cdot \exp(x^2\alpha + x\beta)v^\gamma$ really satisfies the equation (3) and $G_{vv} < 0$. But this will not be shown here. Using (4) we get a closed form solution for the optimal position

$$\hat{h}(t, x) = (\beta(t) + 2x\alpha(t) - \kappa(x - \theta)/\eta^2 + \rho\sigma/\eta + 0.5)/(1 - \gamma) \quad (5)$$

In the next part we will show the backtesting implementation using KDB+/Q.

Introduction

We aim to backtest this strategy with 10000 paths. Each path will have $5 * 252$ points. Instead doing simulate-and-forget we will store all the paths in the memory for further path-dependent risk numbers calculation. For this approach we need a system that is able to store a lot of data in the memory and that is also very performant. We will use KDB+/Q.

What is KDB+/Q ?

- ▶ KDB+ is an in-memory, column-based database and Q is the query language developed by Arthur Whitney and commercialized by kx.
- ▶ Q is a terse variant of a programming language, short apl.
- ▶ The evaluation is from right to left.

We will use KDB+/Q to

- ▶ sample normal distributed random numbers.
- ▶ generate multivariate stochastic differential equations.
- ▶ calculate risk numbers like maximum drawdown, sharpe ratio and the terminal pnl.

Normal distributed random variables using Box Mueller Method

To sample normal distributed random numbers we will use the method from Box-Mueller.

Let u_1, u_2 be independent random variables uniform on $(0, 1)$ then the random variables

$$\sqrt{-2 \ln u_1} \cdot (\cos, \sin)(2\pi u_2)$$

will be normal distributed with mean zero and unit variance.

Code in Q.

```
tmp: 'int$num*step-1
N: sqrt[-2*log tmp?1 f]*:/:(cos; sin)@\:2*pi*tmp?1 f
'Z'W set'num cut/:L mmu N
```

How does it work?

```
tmp: 'int$num*step-1
```

```
num*step-1
```

```
'int$a
```

```
tmp:a
```

Putting everything together.

multiply num with step-1

convert a into an int

assign the value a to tmp

```
tmp : 'int$num * step - 1
```

convert right side to int

assign right side to tmp

How does it work?

```
N: sqrt[-2*log tmp?1 f]*/:(cos;sin)@\:2*pi*tmp?1 f
```

```
tmp?1f
```

sample tmp-times uniform distributed random numbers between zero and one

```
(cos;sin) @\: 1
```

apply the list l to each function on the left

```
l1 */: l2
```

multiply list l1 to each element on the right

Putting everything together.

```
sqrt[-2 * log tmp?1.0] */: (cos; sin) @\: 2 * pi * tmp?1.0
```

list of n floats

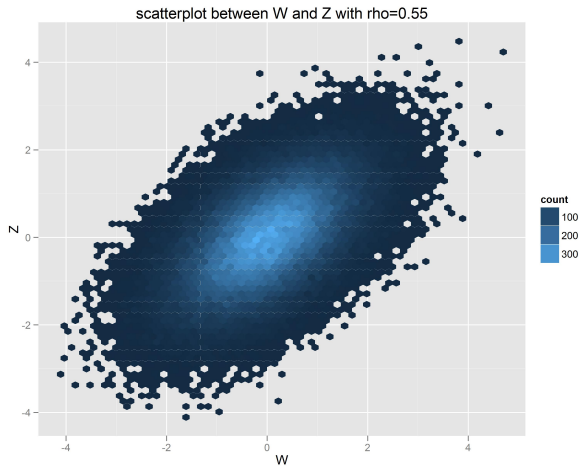
n random floats between 0 and 1

list of n floats

apply each element on left to right

element-wise multiply each element on the right to the left

Visualization of the normal random numbers



Price and Spread Process

The s.d.e. for the asset is

$$\begin{aligned}d(\log S_t) &= dS_t/S_t - 0.5(dS_t/S_t)^2 \\ &= (\mu - 0.5\sigma^2)dt + \sigma dW_t\end{aligned}$$

We solve the differential equation and take equidistant time intervals Δ for the simulation.

$$\begin{aligned}S_{n+1} &= S_n \exp((\mu - 0.5\sigma^2)\Delta + \sigma\sqrt{\Delta}\epsilon_{n+1}) \\ &= f(S_n, \epsilon_{n+1})\end{aligned}$$

Here ϵ_n are normal distributed random variables.

Code in Q:

```
sdelta:sqrt delta:matur%steps  
f:{x*exp(delta*mu-0.5*vola*vola)+vola*sdelta*y}  
A:flip f scan enlist[spot],Z
```

How does it work?

f:{}

enlist[spot]

11 , 12

flip m

f scan l

f is a dyadic function with x and y as its arguments

put the spot into a list

join the two lists to form a general list

transpose the matrix m

create a new list by applying the function f to successive items of the list argument l

Putting everything together.

```
flip          f scan  enlist[spot] ,Z
```

└──────────────────┘
put the spot into a list

└──────────────────┘
join the lists together

└──────────────────────────────────┘
apply the function successive to the element of the list

└──┘
transpose the matrix to get a vector of path

The calculation of the spread is similar.

```
ekh:exp neg kappa*delta
skh:sqrt eta*eta%2f*kappa*1f-ekh*ekh
g:{(x*ekh)+(theta*1f-ekh)+y*skh}
spr:neg A-B:A*exp lspr:flip g scan enlist [theta],W
```

Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

- └ Stochastic Differential Equations

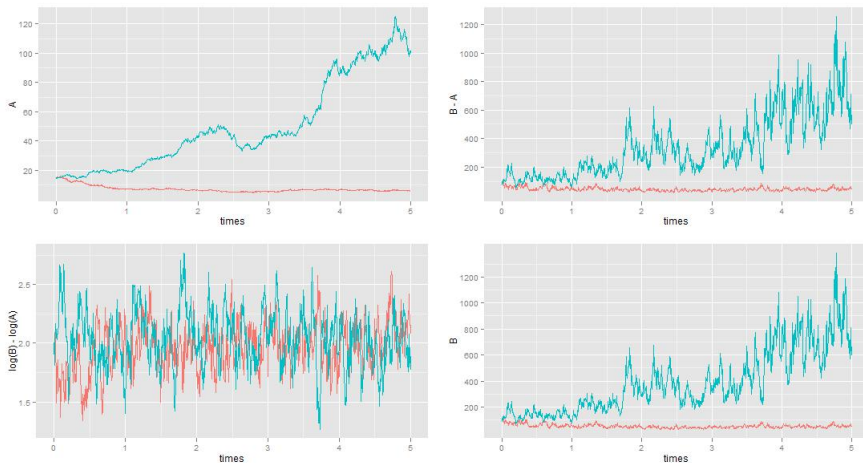


Figure: two example paths

The function h calculates the optimal position. Code in Q:

```
dpos: flip deltas flip pos:h[times;lspr]
scfs: flip sums flip cfs:neg spr*dpos
pnl:(tpnl:spr*pos)+scfs
ret:pnl%abs scfs
sh:sqrt[252]*(%) . stats:(avg;dev)@\: flip ret
mdd:{max neg x-maxs x}
mdds:mdd each pnl
```

How does it work?

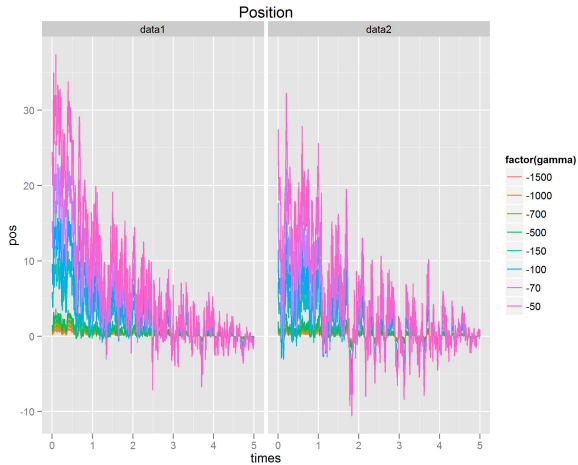
```
pos:h[times;lspr]
dpos:flip deltas flip pos
cfs:neg spr*dpos
scfs:flip sums flip cfs
pnl:(tpnl:spr*pos)+scfs
ret:pnl%abs scfs
stats:(avg;dev)@\:flip ret
sh:sqrt[252]*(%) . stats
{max neg x-maxs x} each pnl
```

optimal position
calculates the daily transaction
calculates the daily cashflows
calculates the total cashflows
calculation of the pnl
rate of return
mean and standard deviation
annualized sharpe ratio
maximum drawdown

Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

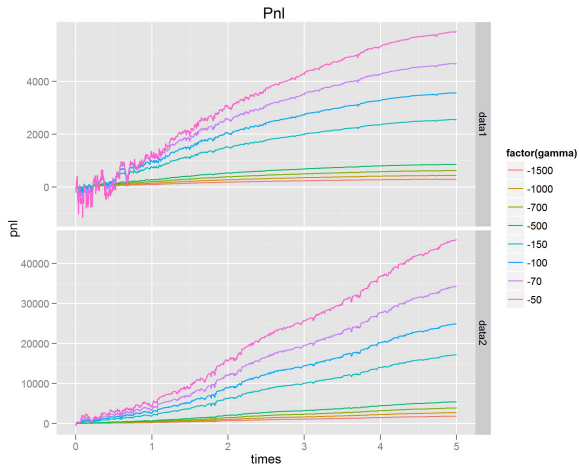
- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

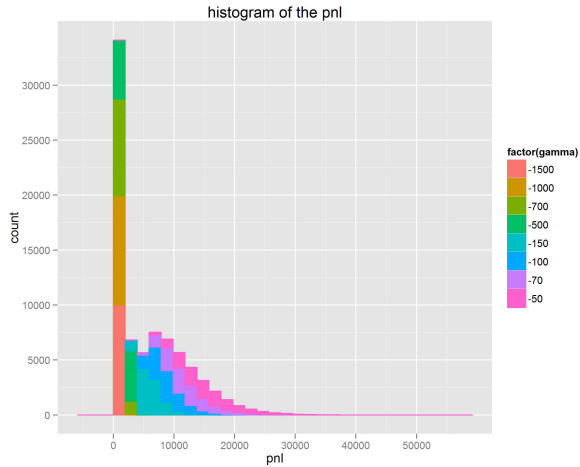
- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

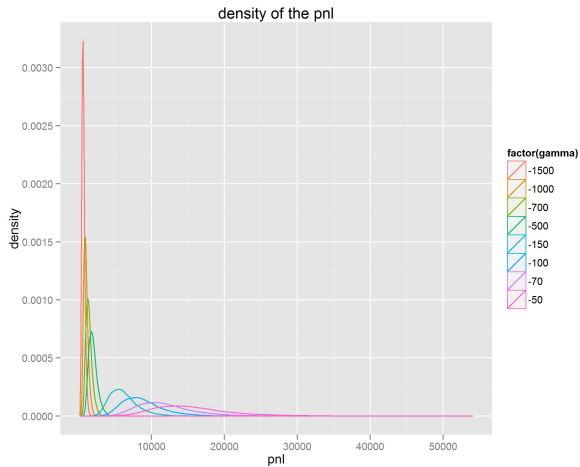
- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

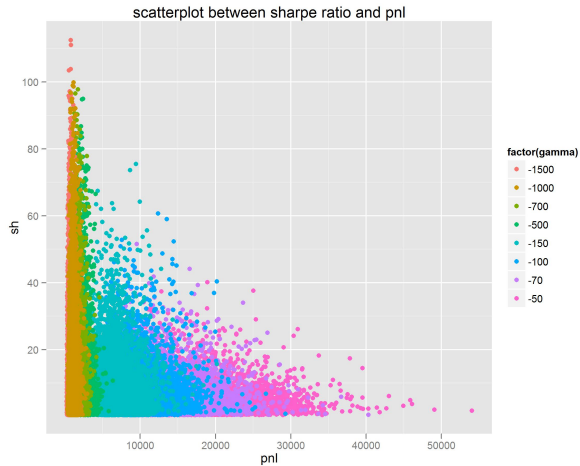
- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

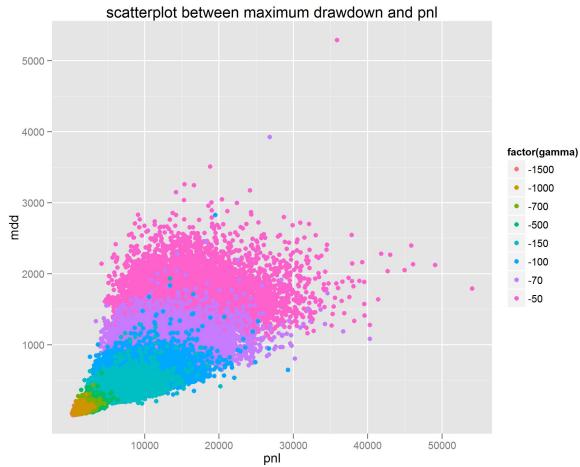
- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

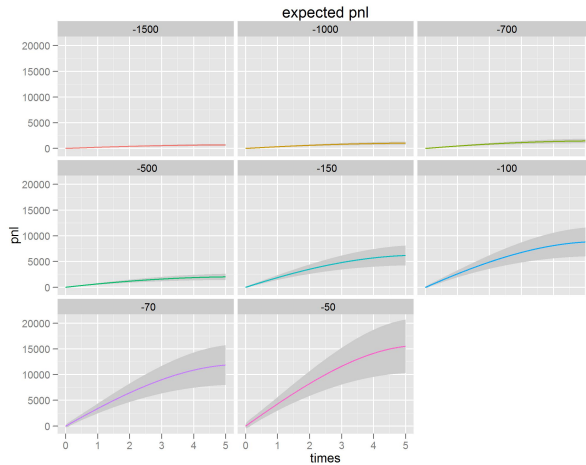
- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

- └ terminal pnl, maximum drawdown and sharpe ratio



Statistical Arbitrage Trading Strategy

- └ Backtesting with KDB+/Q

- └ terminal pnl, maximum drawdown and sharpe ratio

