# Clamping Down on Arbitrage

Peter Jäckel\*

First version:20th September 2013This version:24th April 2014

#### Abstract

We present a practical method for the interpolation of implied Black volatility designed to avoid spurious arbitrage even for extreme or marginal inputs.

#### 1 Introduction

Forty years ago, Fischer Black and Myron Scholes published their article on what is now known as the *Black-Scholes-Merton* formula for the valuation of plain vanilla options [BS73] on equities. Four years later, Fischer Black showed that essentially the same formula, but based on the *forward* instead of the *spot* value of the underlying, can be applied to vanilla options on pretty much any underlying financial observable class<sup>1</sup>. Since then, option traders worldwide have been expressing and comparing option prices in terms of their implied Black volatilities. One of the many reasons of convenience for doing so is that, when expressed as implied volatility, option prices can be compared easily *across* different strikes without having to do some kind of non-linear adjustment for the moneyness of the contract. The implied volatility for any given strike is a direct measure for the relative uncertainty associated with that contract. After all, *uncertainty*, in a manner of speaking, is what the word *volatility* actually means in plain English. Ironically, various authors in mathematical finance have taken this custom of expressing option prices as implied volatilities by trading practitioners to indicate that traders *believe* that some kind of mathematical process volatility is a certain constant for the life of that contract struck at its own specific strike, whilst it is a different constant for other strikes, and have used this as a criticism of the Black-Scholes-Merton formulation. In reality, though, traders

<sup>\*</sup>Deputy head of Quantitative Research, VTB Capital

Key words and phrases. implied volatility interpolation, arbitrage-free.

<sup>&</sup>lt;sup>1</sup>Fischer Black's article [Bla76] was actually written specifically on commodities but the *Black* formula was very soon adopted across all asset classes.

express option prices as implied volatilities *precisely* because they *know very well* that there are no constants, in order to compare different levels of riskiness across strikes on a like-for-like scale.

Nowadays, the use of implied volatility extends to the marking of *implied volatility* surfaces in aid of being able to value large portfolios of option positions across a diverse range of strikes and expiries. For this purpose, in practice, implied volatility tends to be interpolated both in the strike direction, and along the expiry axis. This is a process fraught with danger since the conditions on avoiding arbitrage in the form of implicit negative calendar spread or butterfly prices are, when expressed as constraints on the functional form of implied volatility, both non-linear and non-trivial. Even if we focus only on interpolation in the strike direction, somewhat surprisingly, there is very little literature and research on the subject of implied volatility interpolation. One viable approach is to use a parametric form for the density resulting from the concept of maximum entropy subject to the constraint of matching the given data points for implied volatility [Ave98, HB04, BK96]. The parametric form for the density is then piecewise exponential, and calibration to the input data requires a multi-variate numerical fitting procedure. Whilst intellectually very appealing, this approach is somewhat limited in the shapes it can attain due to the piecewise exponential nature of the density function. An alternative was proposed in [Kah04]. That approach is more akin to a conventional interpolation method. Being based on a functional form that is written in terms of the option price formula, it is essentially the addition of an affine function of the strike (not the log-strike!) to the Black option price formula, and as a consequence, little control is retained about the shape of interpolation in between nodes even if the input data would be perfectly amenable to other, smoother, interpolation forms. In addition to that, the approach in [Kah04] heavily relies on multivariate non-linear root finding (i.e., parametric calibration) and invariably incurs all the associated potential issues such as costly evaluation, residual calibration inaccuracy, and so on.

In this article, we present a procedure that retains, as much as possible, the similarity to an originally chosen smooth interpolation method, but avoids arbitrage when that original interpolation method would give rise to it. The new method still avoids arbitrage even under marginal conditions such as when option prices are co-linear, thus implying regions of zero density, whilst, in a manner of speaking, keeping the overall implied volatility profile as smooth as possible. Importantly, the presented method is fully analytic and requires no numerical calibration<sup>2</sup> at any stage which we find desirable for what is intended to be used as a parametric interpolation method.

 $<sup>^{2}</sup>$  The only arguable exception being the involved calculations of Black volatilities implied from prices, though, this is also required in [Kah04], inside an outer calibration of their interpolation method. We say *arguable* because the implied volatility calculation method in [Jäc13] attains maximum accuracy with precisely two iterations and can for all intents and purpose be considered to be at least semi-analytical, or possibly even as analytical as the cumulative normal function and its inverse.

# 2 No-arbitrage conditions for implied volatility interpolation across strikes

The (undiscounted) Black option formula can be written as

$$B(F, K, \sigma, \theta) = \theta \cdot \left[ F \cdot \Phi \left( -\theta \cdot \left( \frac{z}{\sigma} - \frac{\sigma}{2} \right) \right) - K \cdot \Phi \left( -\theta \cdot \left( \frac{z}{\sigma} + \frac{\sigma}{2} \right) \right) \right]$$
(2.1)

where we have set the time to expiry to 1 without loss of generality, and defined

$$\theta := \pm 1$$
 for calls/puts (2.2)

and

$$z := \ln(K/F). \tag{2.3}$$

When volatility is a function of strike, i.e.,  $\sigma = \sigma(K)$ , the first and second derivative of the vanilla option price

$$v(K) := B(F, K, \sigma(K), \theta)$$
(2.4)

with respect to K are:-

$$v'(K) = K \cdot \varphi \left(\frac{z}{\sigma} + \frac{\sigma}{2}\right) \cdot \sigma' - \theta \cdot \Phi \left(-\theta \cdot \left(\frac{z}{\sigma} + \frac{\sigma}{2}\right)\right)$$
(2.5)

$$v''(K) = \frac{\varphi\left(\frac{z}{\sigma} + \frac{\sigma}{2}\right)}{4K\sigma^3} \cdot \left[4\sigma^2 + 4\left(\sigma^2 - 2z\right)K\sigma\sigma' + 4K^2\sigma^3\sigma'' + \left(4z^2 - \sigma^4\right)K^2\sigma'^2\right] \quad (2.6)$$

For there to be no arbitrage to be implied by the volatility function  $\sigma(K)$ , it is necessary for the second derivative, which is the risk-neutral *Bronzin-Breeden-Litzenberger* density ([Bro08, page 51, equation (17.a)], and [BL78])

$$\psi_{\text{BBL}}(K) = v''(K) \tag{2.7}$$

to be non-negative:

$$\psi_{\rm BBL} \ge 0. \tag{2.8}$$

Recasting equations (2.5) and (2.6) by the aid of the transformation

$$f(z) := \sigma(K)^2 \tag{2.9}$$

leads to

$$v'(K) = \frac{f'}{2\sqrt{f}} \cdot \varphi(\zeta) - \theta \cdot \Phi(-\theta \cdot \zeta)$$
(2.10)

$$v''(K) = \frac{\varphi(\zeta)}{4K\sqrt{f}} \cdot \left[2f'' + \left(z \cdot \frac{f'}{f} - 2\right)^2 - f'^2 \left(\frac{1}{4} + \frac{1}{f}\right)\right]$$
(2.11)

$$v'''(K) = \frac{\varphi(\zeta)}{K^2 \sqrt{f}} \cdot \left[ -\frac{3}{2} - \frac{3f''}{4} + \frac{f'''}{2} - \frac{f'}{8} - \frac{3f'}{2f} - \frac{3f''f'}{16} - \frac{3f''f'}{4f} + \frac{3f'^2}{32} \right]$$
(2.12)

$$+ \frac{3f'^2}{8f} + \frac{f'^3}{128} + \frac{3f'^3}{8f^2} + \frac{f'^3}{16f} - \frac{z}{f} - \frac{3f''z}{2f} + \frac{3f'z}{2f} + \frac{9f'^2z}{4f^2} + \frac{3f'^2z}{16f} + \frac{3f'z}{2f^2} + \frac{3f''f'z^2}{4f^2} - \frac{3f'^2z^2}{8f^2} - \frac{3f'^3z^2}{4f^3} - \frac{f'^3z^2}{16f^2} - \frac{3f'^2z^3}{4f^3} + \frac{f'^3z^4}{8f^4} \right]$$

with

$$\zeta := \frac{z}{\sqrt{f}} + \frac{\sqrt{f}}{2} . \tag{2.13}$$

In order to obtain an intuition as to what the condition (2.8) means for asymptotically high and low strikes, we substitute the affine form

$$\lim_{|z| \to \infty} f(z) \approx a + b \cdot z \tag{2.14}$$

into (2.11). This results in the requirement

$$b \leq 2 \tag{2.15}$$

and thus we must have the asymptotic behaviour

$$\lim_{|z| \to \infty} |f'(z)| \leq 2 \tag{2.16}$$

which is one of the main results of R. Lee's moment formula [Lee04]. Based on these findings and the simpler form of (2.6) in comparison to (2.11), we generally prefer all interpolation of implied volatility to be done explicitly in (z, f)-coordinates by the aid of transformation (2.9) from the original data given in terms of  $(K, \sigma)$  pairs. As for extrapolation, linear extrapolation of f(z) gives satisfactory results, and is asymptotically sound as long as it is either flat or increasing with a linear coefficient of no more than 2 (in absolute value), though we will give the precise conditions for the extrapolation to be free of arbitrage later in section 8.

In order to avoid arbitrage, any given interpolation method of the implied variance function f(z) must satisfy

$$2f'' + \left(z \cdot \frac{f'}{f} - 2\right)^2 - f'^2 \left(\frac{1}{4} + \frac{1}{f}\right) \ge 0.$$
(2.17)

Whilst this is in terms of symbolic complexity significantly more manageable than demanding that the right hand side of (2.6) be non-negative, it is still in practice effectively intractable as a constraint to any interpolant. This is essentially the starting point of the research presented in this article.

### **3** Interpolation of option prices

Denote C(K) as the (undiscounted) price of a call option struck at K. For a given ordered set of strike/price pairs  $\{(K_i, C_i)\}$  for  $i = 1 \cdots n$ , augmented by the strike zero<sup>3</sup>

 $<sup>^{3}</sup>$ We explain the somewhat esoteric reason for the need to add the zero strike data point to the option price set in section 4.1.

with associated put price of zero (and call price of F), arbitrage is present in the data if for any i either of

$$C_i > C_{i+1} \tag{3.1}$$

$$\frac{C_{i-1}}{K_i - K_{i-1}} - C_i \cdot \left(\frac{1}{K_i - K_{i-1}} + \frac{1}{K_{i+1} - K_i}\right) + \frac{C_{i+1}}{K_{i+1} - K_i} \ge 0$$
(3.2)

fails to hold true. Any interpolation C(K) of the option price data generates spurious arbitrage if either of the two conditions

$$C' < 0 \qquad \text{and} \qquad C'' \ge 0 \tag{3.3}$$

is violated. In section 4, we discuss further conditions that pose a situation of arbitrage which need to be considered when input data are analysed and potentially filtered prior to even attempting an arbitrage-free interpolation.

**Remark 3.1.** A marginal situation of arbitrage is the case when the call option price C(K) for increasing K levels out at a positive number  $C(K^*) = C_{\min} > 0$  at some critical strike  $K^*$ , and remains constant thereafter, i.e.,

$$C(K) = C_{\min} \quad \forall \quad K \ge K^* . \tag{3.4}$$

This makes

$$C'(K) = 0 \quad \forall \quad K > K^* .$$
 (3.5)

which translates equation (2.10) to

$$\frac{f'}{2\sqrt{f}} \cdot \varphi(\zeta) - \Phi(-\zeta) = 0 \quad \forall \quad z \ge z^* := \ln(K^*) .$$
(3.6)

Substituting the asymptotic expression [AS84, equation 26.2.12] to first order

$$\Phi(-|\zeta|) \approx \frac{\varphi(\zeta)}{|\zeta|} \cdot \left(1 - \frac{1}{z^2} + \ldots\right)$$
(3.7)

for large z, we obtain the asymptotic ordinary differential equation

$$f' = 4f/(2z+f)$$
(3.8)

which has the general solution in terms of the inverse function z(f)

$$z(f) = \frac{f}{2} - c \cdot \sqrt{f} \tag{3.9}$$

for some constant c. This gives us the interesting result of the *limiting asymptotic form* 

$$\sigma(K) = c + \sqrt{c^2 + 2\ln(K)}$$
 (3.10)

for large K. Essentially the same result, apart from signs, can be obtained for  $K \to 0$  by the aid of expressing put options as rescaled call options on the reciprocal of the underlying. We do not use these results in the following, but mention it in aid of appreciating our choices of inter-, and particularly, extrapolation. **Remark 3.2.** Equation (3.9) implies

$$f' = \left[\frac{\mathrm{d}z(f)}{\mathrm{d}f}\right]^{-1} = \frac{2}{1 - c/\sqrt{f}}$$
 (3.11)

and equally (without derivation) for put options apart from the sign, whence

$$\lim_{z \to \pm \infty} f' = \pm 2 \tag{3.12}$$

for the limiting asymptotic form, as is of course to be expected from (2.16).

Returning to the consideration of option price interpolation, it is clear that any interpolation C(K) method which, given strictly monotone and at least marginally convex data  $\{(K_i, C_i)\}$ , preserves the conditions (3.3) will by construction be free of arbitrage. We can therefore, in principle, design an implied volatility interpolant by transforming all input implied volatilities first to call and, for strikes below the forward, put options (to avoid roundoff truncation), respectively, interpolate on prices, and transform back to volatilities by implication. This approach does of course require an efficient and accurate implied volatility function that works even extremely far away from the money, though, fortunately, that is readily available [Jäc13]. Secondly, we need an interpolation that preserves monotonicity and convexity. For this purpose, we employ the rational cubic method of Delbourgo and Gregory [DG85], specifically with their geometric mean method of choosing the slopes at interpolation interval boundaries given in their equations (3.25) and (3.26).

In principle, this route via option prices that are interpolated under observation of (3.3), is perfectly viable. We show an example in figure 1, along with a conventional interpolation of implied volatilities. Intriguingly, there appear to be some waves when comparing to a conventional smooth implied volatility interpolation. Extending the strike range in figure 2 demonstrates how bad this can get. We emphasize that the generated waves are not an artefact of the choice of price interpolation method. In fact, many other smooth (but not necessarily shape-preserving) interpolation methods such as Akima splines, Catmull-Rom splines, or natural cubic splines, generate even worse examples, and very easily result in arbitrageable output as shown in figure 3. Clearly, the spurious oscillations we see in figures 1 and 2 are undesirable, especially when implied volatility should really be nearly flat. In other words, while shape-preserving interpolation of prices is workable, it can be outright ugly, and would not be acceptable to any trading practitioner.

### 4 Input data filtering

In practical applications, the data sets provided for the interpolation of implied volatilities are frequently not filtered for stale data points, or marginal inaccuracies. Most frequently,



FIGURE 1: Interpolation of implied volatilities via transformation to call option prices and shapepreserving interpolation for an arbitrary data set. The smooth reference interpolation is done as variance over log-strike, i.e., as the function f(z), also with the rational cubic method of Delbourgo and Gregory.



FIGURE 2: Same as figure 1 for a wider range of strikes.

this causes problems in the wings of the implied volatility profile. It is therefore, subject to careful judgement as to the target usage, sometimes desirable to have a procedure that recognizes all data points that are not viable for any given time horizon and thus cannot be included in an arbitrage-free interpolation logic, and, depending on subjective preferences, are either excluded from the data set, or give rise to an application exception that flags a severe data error. A pragmatic approach is to distinguish between strictly intolerable data points, and those that can be remedied as follows.



FIGURE 3: Same data as in figures 1 and 2. This time, call option prices were interpolated with natural cubic splines. Interpolated prices have no implied volatility for  $K/F \in [11, 14]$  giving rise to the associated implied volatility curve dropping to zero in that range.

#### **Removable violations**

Starting with the leftmost strike, one may allow dropping all points  $\{(K_i, \sigma_i)\}$  to the left, at which

• put options have non-positive value, i.e.,

$$P_i \leq 0, \qquad (4.1)$$

with  $P_i := P(K_i)$ ,

• the zero-strike-complemented put butterfly, i.e., an asymmetric put option butterfly built over the strikes  $\{0, K_i, K_{i+1}\}$ , indicates that there is not no positive probability in  $(0, K_{i+1})$ , i.e.,

$$\frac{P(0)}{K_i - 0} - \frac{P(K_i)}{K_i - 0} - \frac{P(K_i)}{K_{i+1} - K_i} + \frac{P(K_{i+1})}{K_{i+1} - K_i} \le 0, \qquad (4.2)$$

which, using  $P(0) \equiv 0$ , is equivalent to

$$P_i K_{i+1} \geq P_{i+1} K_i,$$
 (4.3)

tested in sequence and stopped when for some i no violations are found.

Next, starting with the rightmost strike, one may allow dropping all points  $\{(K_i, \sigma_i)\}$  to the right, at which

• call options have non-positive value, i.e.,

$$C_i \leq 0, \qquad (4.4)$$

with  $C_i := C(K_i)$ ,

• call options are not decreasing, i.e.,

$$C_{i-1} \leq C_i, \qquad (4.5)$$

tested in sequence and stopped when for some i no violations are found.

#### Strictly intolerable violations

Any interior data point  $\{(K_i, \sigma_i)\}$  for  $i = 2 \cdots n - 1$  at which either

• put options are not increasing (for strikes below the forward), i.e.,

$$P_i \geq P_{i+1}$$
 for any *i* such that  $K_i < F$ , (4.6)

• or call options are not decreasing (for strikes above the forward), i.e.,

$$C_i \geq C_{i-1}$$
 for any *i* such that  $K_i > F$ , (4.7)

• or (asymmetric) put or call butterflies<sup>4</sup> are negative, i.e.,

$$\max\left(\frac{P_{i-1}}{\Delta K_{i-1}} - \frac{P_i}{\Delta K_{i-1}} - \frac{P_i}{\Delta K_i} + \frac{P_{i+1}}{\Delta K_i}, \frac{C_{i-1}}{\Delta K_{i-1}} - \frac{C_i}{\Delta K_{i-1}} - \frac{C_i}{\Delta K_i} + \frac{C_{i+1}}{\Delta K_i}\right) < 0 \quad (4.8)$$

with

$$\Delta K_i := K_{i+1} - Ki \tag{4.9}$$

can not be helped. Here, we must reject any attempt of arbitrage-free interpolation since it simply cannot work.

#### 4.1 A special case: put versus digital at the lowest strike

Subject to all of the above lateral removable violations having been filtered out, it is under rare circumstances still possible to have arbitrage induced by an interaction between the digital option price as implied by the (initially unmodified) interpolator of variance-overlog-strike f(z) and the vanilla put option price at the lowest strike  $K_1$ . Denoting the digital

<sup>&</sup>lt;sup>4</sup>For round-off reasons, it is advisable to test against the larger of the put and the call butterfly.

put option as  $P'_1$ , with  $P' \equiv v'$  to be computed via (2.10), this situation of arbitrage is given when

$$K_1 P_1' < P_1.$$
 (4.10)

In this case, when the digital put option is too cheap relative to the vanilla put option, one could construct a trade consisting of a short position in the put option struck at  $K_1$ and a long position in the digital struck at  $K_1$  with notional  $K_1$  such that the trade would have upfront negative cost (i.e., generates upfront cash in putting on the trade), and at maturity breaks even when the final spot  $S_T$  meets  $S_T = 0$  or  $S_T \ge K_1$ , and positive if  $0 < S_T < K_1$ . In terms of risk-neutral probability, this situation indicates that either there must be positive risk-neutral probability located at strikes *less than zero*, or, that there is negative probability somewhere in  $(0, K_1)$ . In this case, we should clearly switch to interpolation over (put) option prices on the interval  $[K_1, K_2]$ , and have the digital put option price at  $K_1$  governed by the price interpolator. We emphasize that this special case must be checked for before any overall arbitrage-free interpolation correction logic as described in the main text as of section 3, and this is why we list it here as part of the initial filtering logic.

This is not the whole story yet, though. The violating condition (4.10) above can in fact *still* arise for a perfectly good monotonicity and convexity preserving price interpolation algorithm for perfectly good and viable put option prices. The reason is that it is in this case left up to the interpolator's internal logic to come up with an estimate as to what  $v'(K_1)$  should be, and, generically, the employed price interpolation algorithm, e.g., Delbourgo and Gregory's rational cubic method, does not have the avoidance of this specific arbitrage situation built into it. If, however, we add the zero-strike put option with zero value to the set of put option prices over which we interpolate, then, the described situation can no longer arise if the price interpolator strictly preserves monotonicity and convexity! This is the reason why we mandated at the beginning of section 3 that all price interpolation must be augmented by the zero strike (and zero value for the put option struck at zero) in order to *anchor* the price interpolator sensibly.

## 5 Clamped interpolants

For there to be no spurious arbitrage to be generated by any interpolation, whether that is as interpolation of prices or as interpolation of variances, the interpolator should produce a function that is of class  $C^1$ , i.e., continuously differentiable. The second derivative, however, does not have to be continuous — it merely must not jump too much to violate the condition that the density must not be negative at any point. With this in mind, we now construct an interpolant that, preferably, interpolates variances over log-strikes smoothly, but on any interval where this leads to negative densities, switches to interpolation of prices over strikes. The crucial point is to connect the two interpolators at the transitions such that implied volatility is not only continuous, but also continuously differentiable. In terms of conventional literature, this means we not only pin down the interpolators at the abscissa/ordinate node pairs where a switch occurs, but we also *clamp* them in the sense that we enforce the local slope at the transition node.

Since natural cubic splines are of class  $C^2$ , giving rise to continuous Bronzin-Breeden-Litzenberger densities, we prefer to use those when we interpolate f(z), i.e., variances over log-strikes. Traditionally, natural cubic splines are usually only documented with the possible external specification of imposed slope values at their end points [PTVF92]. It is of course straightforward to extend this concept to the clamping at interior nodes by viewing the full cubic spline with internal node slope conditions as a sequence of subsplines, clamped at the internal nodes at which the slope needs to be explicitly specified, instead of automatically generated by the spline. We show an example for clamped natural cubic and rational cubic splines in figures 4, 5, and 6 for some arbitrary interpolation data. The clamping causes the second derivative locally to jump, but we will later ensure that





FIGURE 5: The first derivative of the splines in figure 4.



FIGURE 6: The second derivative of the splines in figure 4. Note the discontinuity at the clamping location (0.6) for both the clamped cubic and the clamped rational cubic spline.

the associated Bronzin-Breeden-Litzenberger density never jumps to negative values. As for the interpolation of option prices, we use the aforementioned shape-preserving rational cubic method of Delbourgo and Gregory. We emphasize that all option price interpolation is always done in duplicate: once for calls, and once for puts, and on any segment where interpolation in prices is to be queried, the choice for whether to use the calls or the puts is given by the moneyness of the segment to avoid the roundoff truncation that would be incurred when deeply in the money options are used to imply volatilities.

### 6 Continuous arbitrage detection

Given a set of discrete normalised log-strikes  $\{z_i\}$  with  $z_i = \ln \frac{K_i}{F}$ , and associated variances  $\{f_i\}$  with  $f_i = \sigma(K_i)^2$ , and a chosen interpolation method for f(z) which we prefer to be natural cubic splines, we need to identify all intervals  $I_j := [z_j, z_{j+1}]$  on which the interpolator f(z) implies arbitrage. Since it is analytically intractable, or at least inefficient, to compute a closed form condition that ensures that the density does not dip below zero inside any interval, even when we have the interpolator f(z) given in explicit form as a cubic polynomial, we use local approximations for the sake of robustness and numerical expediency. For any interval  $I_i$ , we carry out the following procedure.

a) From equation (2.10), ensure that we have

$$0 \leq v'(K) \leq 1$$
, (6.1)

with v(K) denoting put option prices (meaning,  $\theta = -1$ ), at both ends of the interval, i.e., at  $K_i$ , and  $K_{i+1}$ , by analytical evaluation of f(z) and f'(z). The violation of (6.1) represents negative digital option prices which can happen even when the density at the same strike is positive. In fact, this happens almost as a rule for the SABR parametric form [HKL02, equation (2.17a)] when calibrated to

interest rate swaption prices for some small interest rate strike range near 1 basis point.

b) Compute the left-hand-side and right-hand-side limits of the density

$$\psi_l := \lim_{z \searrow z_i} (\psi_{\text{BBL}}) \tag{6.2}$$

$$\psi_r := \lim_{z \nearrow z_{i+1}} (\psi_{\text{BBL}}) \tag{6.3}$$

by the aid of  $\psi = v''(K)$  and equation (2.11). This is to be done by analytical evaluation of the respective limits for f(z), f'(z), and, importantly, f''(z), inside the interval  $I_i$ . We emphasize the last point since, even though the interpolator f(z) will always remain of class  $C^1$ , it is possible that the value f''(z) jumps at the interval boundary, either because the original interpolation method was not  $C^2$ to start with, or because an earlier correction stage already resulted in a clamping of f(z) at either boundary. We will come back to this latter issue at the end of section 7.

c) Compute the log-mid-point density

$$\psi_m := \psi_{\text{BBL}}(K_m) \tag{6.4}$$

with

$$K_m := F \cdot e^{\frac{1}{2}(z_i + z_{i+1})} \tag{6.5}$$

also analytically from equation (2.11).

d) Denoting  $K_l := K_i$  and  $K_r := K_{i+1}$  for easier association to the left and right hand side ends of the interval, compute the slope of the density at  $K_l$ ,  $K_m$ , and  $K_r$  using

$$\psi'_l = v'''(K_l) , \quad \psi'_m = v'''(K_m) , \quad \psi'_r = v'''(K_r)$$
(6.6)

and (2.12) from the variance-over-log-strike interpolator f(z). Next, assess if there is a local minimum of the density, and estimate its location, separately on both halves of the interval. First, fit a cubic form for the density function to the locations, levels, and slopes for the left half given by  $(K_l, \psi_l, \psi'_l)$  and  $(K_m, \psi_m, \psi'_m)$ . Compute the locations of the extremums of this cubic form. If there is a minimum, and its location  $K_l^*$  is inside the subinterval  $[K_l, K_m]$ , then evaluate the density  $\psi_{\text{BBL}}(K_l^*)$ analytically from equation (2.11). If the density  $\psi_{\text{BBL}}(K_l^*)$  is positive (and all is well still), proceed to the right half interval with exactly the same check based on the locations, levels, and slopes  $(K_m, \psi_m, \psi'_m)$  and  $(K_r, \psi_r, \psi'_r)$ . If the respective cubic form for the right half interval suggests the existence of a local minimum at  $K_r^*$ , evaluate the density  $\psi_{\text{BBL}}(K_r^*)$  analytically from equation (2.11).

If any of the explicitly computed density check values are negative, the interval  $I_i$  is marked as *defective* for subsequent corrective action by interpolator switching as is discussed in the next section.

### 7 Smooth connections

Given an initial smooth (at least  $C^1$ ) interpolator f(z) of variances over log-strikes, for any interval  $I_i := [z_i, z_{i+1}]$  (with  $z_i := \ln \frac{K_i}{F}$  and  $z_{i+1} := \ln \frac{K_{i+1}}{F}$ ) that was identified as defective in the sense that the density implied by f(z) is negative somewhere in  $I_i$ , we decide whether for price interpolation on this interval we wish to use call or put option prices. In practice, we use call options if  $\frac{1}{2}(K_i + K_{i+1}) \ge F$ , i.e, if the mid-point is at or above the forward. In notation, we define the call/put indicator flag  $\theta_i$  for  $I_i$  as

$$\theta_i := 2 \cdot \mathbf{1}_{\left\{\frac{1}{2}(K_i + K_{i+1}) \ge F\right\}} - 1 \tag{7.1}$$

and set the vanilla price v(K) on  $K \in [K_i, K_{i+1}]$  as

$$v(K) := B(F, K, \sigma(K), \theta_i) \quad \forall \quad K \in [K_i, K_{i+1}].$$

$$(7.2)$$

We now compute the vanilla price slope values on the boundaries of the interval implied by the variance-over-log-strike interpolator f(z) as

$$\tilde{d}_j := \frac{f'(z_j)}{2\sqrt{f(z_j)}} \cdot \varphi(\zeta_j) - \theta_i \cdot \Phi(-\theta_i \cdot \zeta_j)$$
(7.3)

with

$$\zeta_j := \frac{z_j}{\sqrt{f(z_j)}} + \frac{\sqrt{f(z_j)}}{2}.$$
 (7.4)

for j = i and j = i + 1. At this point, we must pay attention to the fact that the above computed interval boundary price slope quantities  $\tilde{d}_i$  and  $\tilde{d}_{i+1}$  are not guaranteed to preserve the convexity condition

$$\tilde{d}_i \leq \Delta_i \leq \tilde{d}_{i+1} \tag{7.5}$$

with

$$\Delta_i := \frac{v(K_{i+1}) - v(K_i)}{K_{i+1} - K_i} \tag{7.6}$$

in the notation of Delbourgo and Gregory. Pragmatically, we allow falling back to the slope that can be computed directly from the (so far unclamped) price interpolator v(K) by setting

$$d_i := \begin{cases} \tilde{d}_i & \text{if } \tilde{d}_i \le \Delta_i \\ v'(K_i) & \text{else} \end{cases}$$
(7.7)

and

$$d_{i+1} := \begin{cases} \tilde{d}_{i+1} & \text{if } \Delta_i \leq \tilde{d}_{i+1} \\ v'(K_{i+1}) & \text{else} \end{cases}$$

$$(7.8)$$

We then add the conditions

$$v'(K_i) \stackrel{!}{=} d_i \tag{7.9}$$

$$v'(K_{i+1}) \stackrel{!}{=} d_{i+1}$$
 (7.10)

to the price interpolator for v(K), thus clamping it in  $K_i$  and  $K_{i+1}$  to slope values that make it consistent with the variance interpolator f(z) in both price and slope at the interval boundaries. This is of course only strictly true if the convexity condition (7.5) held to start with: if either  $\tilde{d}_i \leq \Delta_i$  or  $\Delta_i \leq \tilde{d}_{i+1}$  was not met, then, at the respective interval boundaries, the price interpolator just retained its own original slope! To remedy this not uncommon situation when the original price interpolator slope prevailed at either side, we adjust the variance interpolator's slope at the respective boundary instead. Specifically,

• if  $\tilde{d}_i > \Delta_i$ , we add the condition

$$f'(z_i) \stackrel{!}{=} 2 \cdot (d_i + \theta_i \Phi(-\theta_i \cdot \zeta_i)) \frac{\sqrt{f(z_i)}}{\varphi(\zeta_i)}$$
(7.11)

and

• if  $\Delta_i > \tilde{d}_{i+1}$ , we add the condition

$$f'(z_{i+1}) \stackrel{!}{=} 2 \cdot (d_{i+1} + \theta_i \Phi(-\theta_i \cdot \zeta_{i+1})) \frac{\sqrt{f(z_{i+1})}}{\varphi(\zeta_{i+1})}$$
(7.12)

to the variance interpolator f(z), thus clamping the variance-over-log-strike interpolator f(z) to the price interpolator.

It will not have escaped the attentive reader that the above reverse clamping of the variance-over-log-strike interpolator f(z) at either boundary  $z_i$  or  $z_{i+1}$  can give rise to the interpolator f(z) now implying negative densities on some other intervals that previously passed the checks described in section 6. For all interpolation types, this can immediately happen in the intervals adjacent to  $I_i$  if those were previously deemed fit for interpolation via f(z). For non-local interpolation methods, i.e., those that attain class  $C^2$  by the aid of global calculations, the change of slope in any point can give rise to changes in any other interval which, in turn, now may imply negative densities. It does indeed happen that other intervals end up being affected by negative densities following the first correction stage. This, however, is easily taken care of by applying an outer iteration to our construction logic to repeat the stages of sections 6 and 7 in sequence until no new intervals are identified requiring adjustments because of negative densities. This sounds worse than it is in practice. Even for extreme data we have not seen the need for more than one extra iteration of the correction stage after the initial pass. In theory, it could of course lead to the whole sequence of intervals all being interpolated via prices, but if this does indeed happen for some extreme data, we still have an arbitrage-free interpolation, which is so much better than the alternative in this case, namely, an interpolation that implies negative densities somewhere in all its interior intervals.

### 8 Extrapolation

We mentioned previously that linear extrapolation of f(z) gives satisfactory results. In this section, we want to qualify this statement further. We focus on extrapolation for high strikes, i.e., the extrapolation in z beyond some endpoint  $z_r$ . Linear extrapolation beyond the last interpolation node  $(z_r, f_r)$  can be written as

$$f(z) = a + b \cdot z \quad \forall \quad z > z_r .$$

$$(8.1)$$

If we consider the last value for in the endpoint, i.e.,  $f_r = f(z_r)$ , as non-negotiable, then the only free variable is the extrapolation slope b. The parameter a is then given by

$$a = f_r - b \cdot z_r . aga{8.2}$$

Substituting (8.1) into (2.11), and demanding that the Bronzin-Breeden-Litzenberger density be positive gives

$$\omega(z) > 0 \tag{8.3}$$

with

$$\omega(z) := [b \cdot z - 2 \cdot f(z)]^2 - b^2 \cdot [1 + f(z)/4] \cdot f(z)$$
(8.4)

where we have made the assumption that f(z) > 0, which is probably safe. The function  $\omega(z)$  is a second order polynomial in z whose quadratic coefficient is  $b^2(4-b^2)$ . Since we already know from equation (2.15) that we must have b < 2, the quadratic form in (8.4) will always have a global minimum at

$$z_{\min} = \frac{(2+a)b^2 - 8a}{b(4-b^2)}$$
(8.5)

which is finite because 0 < b < 2. Substituting (8.2), we can compute that the location of the minimum of  $\omega(z)$  is to the left of  $z_r$  if

$$b < \tilde{b}_r \tag{8.6}$$

with

$$\tilde{b}_r := \frac{\sqrt{z_r^2 + 2f_r^2 + 4f_r} - z_r}{1 + f_r/2} , \qquad (8.7)$$

which is always positive. When  $z_{\min} > z_r$ , and the minimum of  $\omega(z)$  is therefore in the extrapolation domain, then

$$\omega(z_{\min}) = \left(a \cdot (4-a) - b^2\right) \cdot \frac{b^2}{4-b^2}$$
(8.8)

is positive if

$$4f_r - f_r^2 + 2z_r(f_r - 2) \cdot b - (1 + z_r^2) \cdot b^2 > 0$$
(8.9)

where we have again substituted (8.2). Condition (8.9) can be reexpressed as

$$\Delta_r > 0 \qquad \text{and} \qquad b < \hat{b}_r \tag{8.10}$$

with

$$\Delta_r := 4z_r^2 - f_r^2 + 4f_r \tag{8.11}$$

and

$$\hat{b}_r := \frac{(f_r - 2)z_r + \sqrt{\Delta_r}}{1 + z_r^2} .$$
(8.12)

We are now in a position to state the condition on b concisely. Linear extrapolation beyond the endpoint  $(z_r, f_r)$ , i.e., for  $z > z_r$ , is free of arbitrage if the digital option price and the Bronzin-Breeden-Litzenberger density at  $z_r$  are free of arbitrage, and

$$0 \leq b < \min(b_{\max}, 2) \tag{8.13}$$

with

$$b_{\max} := \begin{cases} \max(\tilde{b}_r, \hat{b}_r) & \text{if } \Delta_r > 0 \\ \tilde{b}_r & \text{else} \end{cases}$$
(8.14)

With respect to extrapolation to the left hand side in z below the point  $(z_l, f_l)$ , i.e., for  $z < z_l$ , we obtain with

$$f(z) = a + b \cdot z \quad \forall \quad z < z_l .$$
(8.15)

the formulae

$$\tilde{b}_l := -\frac{\sqrt{z_l^2 + 2f_l^2 + 4f_l} + z_l}{1 + f_l/2}$$
(8.16)

$$\Delta_l := 4z_l^2 - f_l^2 + 4f_l \tag{8.17}$$

$$\hat{b}_{l} := \frac{(f_{l} - 2)z_{l} - \sqrt{\Delta_{l}}}{1 + z_{l}^{2}}$$
(8.18)

and

$$b_{\min} := \begin{cases} \min(\tilde{b}_l, \hat{b}_l) & \text{if } \Delta_l > 0 \\ \tilde{b}_l & \text{else} \end{cases}$$

$$(8.19)$$

Extrapolation for  $z < z_l$  is free of arbitrage if the digital option price and the Bronzin-Breeden-Litzenberger density at  $z_l$  are free of arbitrage, and

$$\max(b_{\min}, -2) < b \le 0$$
 (8.20)

In practice, we have observed that, when the respective discrimant  $\Delta_{(\cdot)}$  is positive, it is rare for the associated value  $\hat{b}_{(\cdot)}$  to exceed the corresponding value  $\tilde{b}_{(\cdot)}$  (in absolute value). Since  $\tilde{b}_{(\cdot)}$  is the limiting value for the slope at the extrapolation boundary for which the density is already negative at the very boundary edge itself, the fact that linear interpolation in f(z) has its limits (due to the fact that the density may become negative for some much higher strike) is in practice not really an issue. It just doesn't happen. The main limits on the extrapolation slope of f(z) are given by the restrictions on f'(z) on the boundary to extrapolation in the form of the conditions (3.3) with (2.10) and (2.11).

#### 8.1 Log-linear factor extrapolation

There is an alternative to linear extrapolation in variance-over-log-strike that is free of arbitrage by construction, and typically leads to a more rapid levelling out of implied volatilities for strikes far away from the money. This extrapolation method is based on the representation of the underlying financial variable as given by a *quantile map* [Jäc05] of a standard normal variate.

Consider that the distribution of the spot variable S is generated by virtue of a mapping function  $\chi(y)$  that maps a standard normal variate y to the logarithm of S, i.e.:

$$S = F \cdot e^{\chi(y)} . \tag{8.21}$$

The quantile map  $\chi(y)$  must of course satisfy the forward for the underlying:

$$1 = \int e^{\chi(y)} \varphi(y) dy . \qquad (8.22)$$

The quantile map representation is a common approach to generate a distribution for S that matches a given implied volatility smile by construction [Jäc05]. In any practical implementation, the quantile map  $\chi(y)$  is represented by an interpolation methodology over a chosen discrete set of values for y, and, typically, combined with *linear extrapolation*.

We write for the linear extrapolation rule of  $\chi(y)$  beyond the last strike the form

$$S/F = e^{\alpha + \beta \cdot y} . \tag{8.23}$$

For a put option struck at a strike K below the extrapolation point  $K_l = F \cdot e^{\alpha_l + \beta_l \cdot y_l}$ , we readily compute from (8.23) the analytical value

$$v(K) = K \cdot \Phi(y_K) - F \cdot e^{\alpha_l + \frac{1}{2}\beta_l^2} \cdot \Phi(y_K - \beta)$$
(8.24)

with

$$K = F \cdot e^{\alpha_l + \beta_l \cdot y_K} \iff y_K = \left[ \ln(K/F) - \alpha_l \right] / \beta_l$$
(8.25)

In order to obtain the extrapolation coefficients  $\alpha_l$  and  $\beta_l$ , we make use of the price  $v_l$  of the put option struck at the extrapolation boundary  $K_l$  as given by the variance-over-logstrike interpolator f(z) and the Black formula. Also, from (2.10), we can compute the digital put option price which, by definition of the quantile map, is equal to  $\Phi(y_l)$ , i.e.,

$$v'(K_l) = \Phi(y_l) \tag{8.26}$$

and thus we also have

$$y_l = \Phi^{-1}(v'(K_l)) . (8.27)$$

Given the put option price  $v_l$ , the digital put price  $\Phi(y_l)$ , and the threshold value  $y_l$ , we need to solve the system of equations (8.24) and (8.25) for  $\alpha_l$  and  $\beta_l$ . This reduces to solving

$$\Phi(y_l) - \frac{v_l}{K_l} = e^{\frac{1}{2}\beta_l^2 - \beta_l \cdot y_l} \cdot \Phi(y_l - \beta_l) .$$
(8.28)

Substituting  $\eta_l := y_l - \beta_l$ , and taking the logarithm, we arrive at

$$\frac{1}{2}\eta_l^2 + \ln\left(\Phi\left(\eta_l\right)\right) = \ln\left(\Phi(y_l) - \frac{v_l}{K_l}\right) + \frac{1}{2}y_l^2.$$
(8.29)

Noting that all the terms on the right hand side are known, and subsuming them into some constant c, finding  $\eta_l$  amounts to solving an equation of the form

$$\frac{1}{2}x^2 + \ln(\Phi(x)) = c \tag{8.30}$$

which we call the quantile slope equation. We show in appendix C how this can be done semi-analytically with little numerical effort. Once we have found  $\eta_l$ , we obtain

$$\beta_l = y_l - \eta_l \tag{8.31}$$

$$\alpha_l = \ln(K_l/F) - \beta_l \cdot y_l . \tag{8.32}$$

For extrapolation above the last interpolation node at  $K_r$ , we omit the derivation and merely state that, given the call option price  $v_r$  struck at  $K_r$ , the digital call price  $-v'(K_r) = \Phi(-y_r)$ , and the threshold value  $y_r = -\Phi^{-1}(-v'(K_r))$ , we need to solve

$$\frac{1}{2}\eta_r^2 + \ln\left(\Phi\left(\eta_r\right)\right) = \ln\left(\Phi(-y_r) + \frac{v_r}{K_r}\right) + \frac{1}{2}y_r^2$$
(8.33)

for  $\eta_r$ , and have

$$\beta_r = y_r + \eta_r \tag{8.34}$$

$$\alpha_r = \ln(K_r/F) - \beta_r \cdot y_r . \tag{8.35}$$

Now that we have the log-linear quantile map extrapolation coefficients, to price a put option struck at  $K < K_l$ , we have from (8.24) and (8.25)

$$F \cdot \left[ e^{z} \cdot \Phi\left(\frac{z-\alpha_{l}}{\beta_{l}}\right) - e^{\alpha_{l} + \frac{1}{2}\beta_{l}^{2}} \cdot \Phi\left(\frac{z-\alpha_{l}}{\beta_{l}} - \beta_{l}\right) \right]$$
(8.36)

where we have substituted  $z := \ln(K/F)$ . In complete analogy, for call options struck above  $K_r$ , we obtain

$$F \cdot \left[ e^{\alpha_r + \frac{1}{2}\beta_r^2} \cdot \Phi\left(\frac{\alpha_r - z}{\beta_r} + \beta_r\right) - e^z \cdot \Phi\left(\frac{\alpha_r - z}{\beta_r}\right) \right] .$$
(8.37)

For any given strike in the extrapolation domains, we can now obtain the implied volatility for that strike by first computing the respectively out-of-the money option, and inverting it to the associated Black volatility as usual.

Before we conclude this section, we wish to provide some analytical understanding of the log-linear factor extrapolation method, and its asymptotic behaviour for large log-strikes z. For this purpose, we use the extrapolation option price formulae (8.36) and (8.37) to derive an expansion for the implied volatility, or, to be precise, for the variance-over-log-strike function f(z). We start this with the postulation that f(z) can be approximated by a finite a power series  $\hat{f}(z)$  of  $\frac{1}{z}$  for large |z|:

$$\hat{f}(z) := \sum_{i=0}^{N} c_i \cdot z^{-i}$$
 (8.38)

We then substitute this into the Black formula (2.1) with  $\sigma = \sqrt{\hat{f}(z)}$ . We equate this with (8.36) for puts when  $\theta = -1$ , and with (8.37) for calls when  $\theta = +1$ . Note that the resulting equations do not permit Taylor expansions in  $\frac{1}{z}$ . Instead, we use the asymptotic expression [AS84, equation 26.2.12] given in equation (3.7) to replace all occurrences of the cumulative normal function. These equations, then, still don't permit any conventional Taylor expansion and we need to match the coefficients of the leading orders of z in the respective exponents. Up to second order in  $\frac{1}{z}$ , for both call and put options, we so obtain and match

$$\frac{F \cdot \beta^3}{z^2 \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \frac{\alpha^2}{\beta^2} + \frac{\alpha + \beta^2}{\beta^2} z - \frac{1}{2\beta^2} z^2} = \frac{F \cdot \sqrt{c_0^3}}{z^2 \sqrt{2\pi}} \cdot e^{-\frac{c_0^4 + 4c_1^2 + 4c_0c_2}{8c_0^3} + \frac{c_0^2 + c_1}{\beta^2} z - \frac{1}{2c_0} z^2}.$$
 (8.39)

This holds for all z when

$$c_0 = \beta^2 \tag{8.40}$$

$$c_1 = \beta^2 (\beta^2 + 2\alpha) \tag{8.41}$$

$$c_2 = \frac{1}{4}\beta^2(\beta^2(5\beta^2 + 16\alpha) + 12\alpha^2) .$$
(8.42)

As we go higher in orders of  $\frac{1}{z}$ , we can use standard Taylor expansion and match powers. Up to order N = 5, the coefficients are:-

$$c_3 = \frac{1}{4}\beta^2(\beta^2(\beta^2(7\beta^2 + 30\alpha - 4) + \alpha(40\alpha - 8)) + 16\alpha^3)$$
(8.43)

$$c_4 = \frac{1}{8}\beta^2 (\beta^2 (\beta^2 (\beta^2 (21\beta^2 + 112\alpha - 30) + \alpha(210\alpha - 96))) + \alpha^2 (160\alpha - 72)) + 40\alpha^4)$$
(8.44)

$$c_{5} = \frac{1}{24}\beta^{2}(\beta^{2}(\beta^{2}(\beta^{2}(\beta^{2}(99\beta^{2} + 630\alpha - 242) + \alpha(1512\alpha - 1044) + 72) + \alpha(\alpha(1680\alpha - 1416) + 144)) + \alpha^{3}(840\alpha - 592)) + 144\alpha^{5})$$

$$(8.45)$$

To demonstrate the shape of the implied volatility profile from the log-linear factor extrapolation methodology in comparison to linear extrapolation of f(z),, we show in figure 7 an example for the 1Y USDRUB smile as observed on the 24th of September 2013. Two interpolators were calibrated to exactly the same data and forced to have the same slopes in implied volatility in the input data node locations. Log-linear factor extrapolation clearly generates a much gentler increase in the wings of the smile. Also noteworthy is the excellent accuracy of the asymptotic expansion for implied volatility given by the square root of  $\hat{f}(z)$  for  $|z| \gg 0$ . For |z| near zero, or even for |z| near the point of extrapolation, however, the asymptotic form is of no use.

For applications where an explicit analytical extrapolation formula is desired instead of the path to compute a price from (8.36) or (8.37) and implying volatility from it, it is possible to use (once again) the rational cubic form of Delbourgo and Gregory [DG85]. We write the extrapolation as the function

$$f(z) = \check{g}(w) \tag{8.46}$$



FIGURE 7: A quantile map interpolator with log-linear factor extrapolation is shown as the solid green line. Extrapolation applies to all abscissa values outside the range spanned by the input volatilities (purple crosses). Compare this with the extrapolation of implied volatility generated by linear extrapolation for f(z), i.e., variance-over-log-strike as shown in the brown dashed line. The latter was interpolated with a cubic spline that was clamped in all of the five input data nodes to a slope that matches that of the quantile map interpolator in the respective node.

with

$$w := 1/z$$
. (8.47)

It follows that

$$\check{g}'(w) = -f'(z) \cdot z^2 \tag{8.48}$$

$$\check{g}''(w) = [2f'(z) + zf''(z)] \cdot z^3 .$$
(8.49)

For extrapolation beyond the boundary  $z_r := \ln(K_r/F)$ , we use a rational cubic form for  $\check{g}(w)$  on the interval  $[0, w_r]$  with  $w_r = 1/z_r$ . The rational cubic form is specified to match

$$\check{g}(0) = c_0 \qquad \check{g}(w_r) = f(z_r)$$
(8.50)

$$\check{g}'(0) = c_1 \qquad \check{g}'(w_r) = -f'(z_r) \cdot z_r^2 \qquad (8.51)$$

and

$$\check{g}''(w_r) = [2f'(z_r) + z_r f''(z_r)] \cdot z_r^3 .$$
(8.52)

Note that the values  $f(z_r)$  and  $f'(z_r)$  are the very same numbers that we previously used to compute the vanilla and digital option prices at  $K_r$  in order to calibrate  $\alpha_r$  and  $\beta_r$ , whence they are already available in our calculations. As for  $f''(z_r)$ , we want to ensure that it is computed from the log-linear factor extrapolation, *not* from the interpolator for f(z) that we use on the inside of the interpolation domain. To this end, we note that the Bronzin-Breeden-Litzenberger density resulting from (8.36) at the boundary  $z_r$  is given by

$$v''(K_r) = \frac{\varphi(y_r)}{K_r \beta_r} \tag{8.53}$$

with

$$y_r = \frac{z_r - \alpha_r}{\beta_r} , \qquad (8.54)$$

which we can set equal to (2.11) in order to solve for  $f''(z_r)$ . Since the rational cubic form (A.1) matches the levels and slopes at both ends of its specification interval by construction, the only thing left to do is to compute the rational cubic shape parameter r such that the second derivative at the right hand side of the interval  $[0, w_r]$  is indeed matched to (8.52), and this is discussed in appendix B. For extrapolation on the put option side, the respective rational cubic form is of course to be defined on  $[w_l, 0]$  with  $w_l = 1/z_l$ . All else follows in complete analogy whence we omit the specifics. The net effect for the shape of implied volatility extrapolation from the rational cubic extrapolation formula derived from the log-linear factor extrapolation method is shown in figure 8. We



FIGURE 8: The rational cubic form of log-linear factor extrapolation for the same data as in figure 7.

argue that the result is highly satisfactory. It remains to be mentioned that the rational cubic extrapolation form for log-linear factor extrapolation is, alas, not guaranteed to be free of arbitrage, though, given the quality of its fit to the analytically guaranteed exact form (which involves implying volatilities from prices), it is unlikely that negative densities or out-of-bounds digitals would ensue. After all, the rational cubic form is accurate to second order near the extrapolation boundary, and accurate asymptotically in 1/|z| to first order for  $|z| \to \infty$ . As always with approximation methods, the proof of its viability will lie in its use in practical applications.

We conclude this section by identifying the following properties of *log-linear factor extrapolation* for implied volatility:-

- For  $|z| \to \infty$ , implied volatility converges to the respective extrapolation side's log-linear factor coefficient  $\beta$ , which we would have expected intuitively given the stochastic *factor* nature of the extrapolation.
- As  $|z| \to \infty$ , implied volatility, to lowest order, converges to  $\beta$  like 1/|z|, i.e., inverse in the log-strike.

## 9 Numerical examples

We show in figures 9 and 10 an example for the subtleties that can make all the difference between the occurrence of negative densities and arbitrage-free interpolation. On the



FIGURE 9: Three different interpolation methods giving rise to almost identical implied volatility profiles.

presented scale, the three interpolation methods seem to result in indistinguishable implied volatility curves. The main difference is that linear extrapolation of option prices below the lowest input strike results in a rapidly decaying implied volatilities, as one would expect. However, if we investigate the Bronzin-Breeden-Litzenberger densities associated with the three interpolation methods in figure 11, the difference becomes painstakingly obvious. The first thing we notice is the pronounced oscillation of the density generated by shape-preserving rational cubic interpolation of option prices. These oscillations do of course correspond to the virtually invisible gentle undulations of the implied volatility profile which, in a more pronounced form, we had already seen in figures 1 and 2. More subtly, but of more financial engineering concern, is the fact that the solid blue line in figure 11, which corresponds to the density generated from the natural cubic spline interpolation of f(z), actually dips into the negative domain for strikes near  $K/F \approx 0.05$ . We show an enlargement of this area in figure 12. The respective strike range may appear



FIGURE 10: The same data as in figure 9 on a logarithmic scale for the abscissa.



FIGURE 11: The Bronzin-Breeden-Litzenberger density associated with the three implied volatility interpolations of figures 9 and 10.

far from the money, but for the valuation of exotic contracts that need to be calibrated to all strikes, this is a real problem. What's more, the fact that the negative density is in this case very far from the money is just a matter of coincidence. After all, since the direct interpolation of f(z) provides no guarantee whatsoever that it does not generate spurious arbitrage: this could happen anywhere! Reassuringly, however, we see that the spurious negative density is indeed remedied by the clamped interpolation methodology presented in this article, as can be seen in the green long-dashed line in figures 11 and 12. We show how small the required adjustment for implied volatility is in figure 13. The fact that the adjusted density displays some positive spikes is of no practical or financial concern. After all, in order to reproduce the given implied volatilities in the interpolation nodes perfectly, the underlying probability has to be distributed somewhere, and if the



FIGURE 12: An enlargement of the area of negative density implied by cubic spline inerpolation of f(z) shown in figure 11.



FIGURE 13: The area where the adjusted implied volatility differs from the original cubic spline inerpolation of f(z). Only two intervals needed to be switched to interpolation of prices over strikes.

only solution that is free of arbitrage involves pinned probabilities, so be it! Also, the situation of pinned probabilities is not even unusual in options markets, for a number of reasons ranging from mere popularity of certain strikes and the effect of hedging, as well as technical traders' views on economically or psychologically critical levels of the underlying. In this case, we notice not only two spikes right next to the input data node locations of K/F = 0.0351 and K/F = 0.0686, but also a rather moderate positive bump near the input node location K/F = 0.049. This is uncannily consistent with the fact the we had originally produced the presented data from a finite-differencing discrete approximation of a stochastic volatility model from which we had computed the input data as implied volatilities from the so numerically obtained option prices. Uncanny, because these original option prices were by construction computed from a discrete set of

pinned probabilities!

In our second numerical example, we have constructed a set of implied volatilities that generate the call option prices for the three strikes  $K/F \in \{2.73, 3.82, 5.34\}$  to form a straight line when charted over strikes. This situation does of course imply that there is a region of zero density between the first and the last of those three strikes. This case is *marginal* in the sense that even the slightest numerical deviation in the wrong direction in that region will give rise to arbitrage, and thus comprises an extremely difficult case for any volatility interpolation methodology. First, we show in figure 14and an enlargement



FIGURE 14: A second example. Once again, the three different interpolation methods appear as almost identical implied volatility profiles.



FIGURE 15: The Bronzin-Breeden-Litzenberger density associated with the three implied volatility interpolations of figure 14.

for the region of particular interest of this test, i.e.,  $K/F \in [2.73, 5.34]$ , in figure 16. the global shape of the implied volatility smiles with the three different interpolation methods,



FIGURE 16: An enlargement of the critical region of zero Bronzin-Breeden-Litzenberger density for the input data of figure 14.



FIGURE 17: The same data as in figure 14, but on a logarithmic scale. The solid light blue line is the absolute difference between the clamped spline method (long-dashed dark green line) and implied volatility via natural cubic splines for f(z) (solid blue line).

namely, via natural cubic splines for f(z), via rational cubic interpolation for option prices, and via the presented methodology of clamped interpolators. Next, we show the associated Bronzin-Breeden-Litzenberger densities in figure 15, Unsurprisingly, the direct interpolation of implied volatility as a natural cubic spline for the variance-over-log-strike function f(z) results in negative density near the middle node of the set of three that should really form an area of exactly zero density. The clamped alternating interpolator, however, as was intended, avoids the negative density by switching to interpolation via option prices. As a consequence, we end up with positive lump sums of probability at the end points of the zero density interval [2.73, 5.34], which is the correct behaviour in this case. Finally, we show in figure 17 the absolute magnitude of the implied volatility adjustment, which highlights that the region of zero density is not the only area that required modification to remain free of arbitrage in this example, albeit that that region clearly was adjusted by the largest amount. Even though the adjustment is at its peak as much as 22 basis points (absolute) of implied volatility, it is important to remember that it is *exactly zero* at all of the original input data nodes. The adjustment only applies to *actual interpolation*!

### 10 Summary

We have described a procedure for the clamped concatenation of variance-over-log-strike interpolators with price interpolators for the sake of smooth interpolation of implied volatility without spurious arbitrage. In short, the process entails:

- 1. Build an underlying implied volatility interpolator based on interpolation of actual variance in the sense  $\sigma^2 \cdot T$  over  $z := \ln(K/F)$  as f(z) with  $\sigma(K) = \sqrt{f(\ln(K/F))/T}$ . This interpolator f(z) is preferably a natural cubic spline with linear extrapolation.
- 2. Detect out-of-bounds digitals and negative density occurrences in all intervals. First, by computing the left-hand-side and right-hand-side limits of the density from the  $C^1$  interpolator f(z) analytically. Then, by checking the midpoint of the interval analytically. Approximate g(K) := v'(K) as a rational cubic matching its value and slopes at either side of the interval, and matching its value in the mid-point of the interval. Note that the slope of g(K) is the density which was previously already computed at both ends and the mid-point of the interval. From this rational cubic fit of g(K) to the boundary values, the mid-point, and the boundary slopes, compute approximations for the slope of the density given by g''(K) at the interval boundaries and the mid-point. Then, from the so computed values for the density  $\psi$  and  $\psi'$ , for both the left-half, and the right-half sub-interval, individually, by the aid of a local cubic fit, assess whether the density has a local minimum and estimate its location. If there is a minimum inside the respective half-interval, check the actual density function analytically from the original  $C^1$  interpolator f(z) at that estimate of the location of the minimum.
- 3. Within each interval that was so identified as being density-defective, implied volatilities will then be computed by interpolation of prices. To do this, build rational cubic interpolators of (out-of-the-money) prices with the *geometric mean* method of choosing the slopes at interpolation interval boundaries to preserve monotonicity and convexity.
- 4. For any density-defective interval  $I_i$ , compute analytically the boundary slopes  $d_i$ and  $d_{i+1}$  of the price function (i.e., the digitals v'(K) at the boundaries of the interval) from the variance interpolator f(z), and impose those slopes on the price

interpolator. Note that this configures the subsequent price interpolation to have slopes on the interval boundaries consistent with the slopes of the original volatility/variance interpolator. This is the key to the appearance of smoothness. Also store the straight line interval slope values

$$\Delta_i := \frac{v(K_{i+1}) - v(K_i)}{K_{i+1} - K_i} . \tag{10.1}$$

These should satisfy the convexity condition

$$d_i \leq \Delta_i \leq d_{i+1} . \tag{10.2}$$

If the convexity condition is not met, say,  $d_i > \Delta_i$ , then set  $d_i$  to the value from the unconstrained price interpolator. The geometric mean calculation method for the interval boundary slopes mentioned above should ensure that those values meet the convexity condition.

- 5. On interval  $I_i$ , now interpolate prices with the so chosen lateral slope values  $d_i$  and  $d_{i+1}$  with rational cubic interpolation. Imply volatilities from the so obtained prices.
- 6. Rebuild the variance interpolator f(z), if it supports internal clamping, with an override of boundary slopes at those locations where the above procedure effectively resulted in overrides of the respective  $d_i$  values. When this happens, we need to repeat the check for defective intervals since this last step can sometimes reintroduce negative densities on other intervals. If any additional intervals are found to be defective, the above ironing-out algorithm needs to be repeated. In principle, it is possible for this to result in an iteration until all intervals are interpolated by prices, but this should only happen for totally corrupt input data.

In addition, we have also analyzed linear extrapolation in variance over logarithmic strike, and derived simple closed form conditions that are necessary and sufficient for the extrapolation to be free of arbitrage. Further, we have explained an alternative for extrapolation based on a log-linear normal factor representation of the underlying, and provided asymptotic expansions to understand and explain the behaviour of this extrapolation methodology in the limit of high or low strikes. Finally, we have given a rational cubic extrapolation form in terms of variance-over-log-strike which provides the compromise of very fast evaluation with a very close fit to the analytically exact log-linear factor extrapolation methodology.

### Acknowledgement

The author is grateful to Charles-Henri Roubinet, Head of Quantitative Research at VTB Capital, for authorizing the release of this note (originally from September 2013) into the public domain.

### A Rational cubic interpolation fit to an interior point

Given an abscissa interval  $[x_l, x_r]$  with boundary function values  $\{y_l, y_r\}$  and associated slope values  $\{s_l, s_r\}$ , the rational cubic interpolation formula of Delbourgo and Gregory reads

$$f(x) = \frac{y_r t^3 + (ry_r - hs_r)t^2(1-t) + (ry_l + hs_l)t(1-t)^2 + y_l(1-t)^3}{1 + (r-3)t(1-t)}$$
(A.1)

with  $h := x_r - x_l$  and  $t := \frac{x - x_l}{h}$  with a suitably chosen cubic shape parameter r > -1. If we want to choose r in order to match an interior point (x, y), equation (A.1) can be solved for r to obtain

$$r_{(x,y)} = \frac{y \cdot (1 - 3t(1 - t)) - y_r t^3 + h s_r t^2 (1 - t) - h s_l t (1 - t)^2 - y_l (1 - t)^3}{y_r t^2 (1 - t) - y t (1 - t) + y_l t (1 - t)^2} .$$
 (A.2)

# B Rational cubic interpolation fit to the second derivative at one end

The rational cubic form (A.1) matches a given second derivative  $f''_l$  at  $x_l$  when

$$r = \frac{\frac{1}{2} \cdot h \cdot f_l'' + (s_r - s_l)}{\Delta - s_l}$$
(B.1)

with  $h := x_r - x_l$  and  $\Delta := \frac{y_r - y_l}{h}$ . The second derivative  $f''_r$  at  $x_r$  is matched when

$$r = \frac{\frac{1}{2} \cdot h \cdot f_r'' + (s_r - s_l)}{s_r - \Delta} .$$
 (B.2)

# C Solving the quantile slope equation

Solving the quantile slope equation (8.30) amounts to finding the root of

$$f(x) = c \tag{C.1}$$

with

$$f(x) := \frac{1}{2}x^2 + \ln(\Phi(x))$$
 (C.2)

We note that

$$e^{f(x)} = \frac{1}{2} \operatorname{erfcx}(-\frac{x}{\sqrt{2}}) \tag{C.3}$$

where  $\operatorname{erfcx}(\cdot)$  is the scaled complementary error function [Cod69], and thus a formal solution to f(x) = c is

$$x = -\sqrt{2} \cdot \operatorname{erfcx}^{-1}(2 \cdot e^c) . \qquad (C.4)$$

However, in contrast to the error function  $\operatorname{erf}(\cdot)$  and the complementary error function  $\operatorname{erfc}(\cdot)$ , unfortunately, no standard implementations for the inverse of the scaled complementary error function are readily available. We therefore proceed with our description of a solution that is accurate to standard IEEE 64 bit floating point precision with two iterations, i.e., two evaluations of f(x).

The quantile slope function f(x) is strictly monotonic, convex, with growth at most of quadratic order, and thus readily amenable for standard iterative root finding procedures. For large |x|, it rapidly converges to the invertible forms

$$\lim_{x \to -\infty} f(x) \approx -\ln\left(-x\sqrt{2\pi}\right) \tag{C.5}$$

$$\lim_{x \to +\infty} f(x) \approx \frac{1}{2}x^2 \,. \tag{C.6}$$

Near the origin, its inverse is well approximated by the rational form

$$x_0^{\text{mid}}(c) := -\Delta_0 \cdot \sqrt{\frac{\pi}{2}} \frac{1 + \frac{\Delta_0}{2} \cdot \left(1 - \frac{\pi}{2}\right)}{1 + \Delta_0 \cdot \left(1 - \frac{\pi}{2} + \frac{\Delta_0}{3} \cdot \left(1 - \frac{\pi}{4}\right)\right)}$$
(C.7)

with

$$\Delta_0 := f(0) - c = \ln(1/2) - c.$$
 (C.8)

This suggests that a good initial guess can be formed by a suitable choice of branch switches on c to combine the asymptotics with the rational form. Choosing the lower branch to be for c < -2.25 and the upper branch for c > 1.4, we obtain the initial guess

$$x_{0}(c) = \begin{cases} -\frac{e^{-c}}{\sqrt{2\pi}} & \text{when} \quad c < -2.25 \\ x_{0}^{\text{mid}}(c) & \text{when} \quad -2.25 \le c \le 1.4 \\ \sqrt{2c} & \text{when} \quad c > 1.4 \end{cases}$$
(C.9)

For the root finding procedure, we recommend the third order Householder method (which is of fourth order in convergence)

$$x_{n+1} = x_n + \nu(x_n) \cdot \frac{1 + \nu(x_n) \cdot h_2(x_n)/2}{1 + \nu(x_n) \cdot [h_2(x_n) + \nu(x_n) \cdot h_3(x_n)/6]}$$
(C.10)

with

$$\nu(x) = -\frac{f(x) - c}{f'(x)}, \qquad h_2(x) = \frac{f''(x)}{f'(x)}, \qquad h_3(x) = \frac{f'''(x)}{f'(x)}.$$
(C.11)

The required derivatives of f(x) are given by the simple expressions

$$q := \frac{\varphi(x)}{\Phi(x)} \qquad f''(x) = 1 - q \cdot (x+q)$$

$$f'(x) = x + q \qquad f'''(x) = q \cdot [x^2 - 1 + q \cdot (2 \cdot q + 3 \cdot x)].$$
(C.12)



FIGURE 18: The quantile slope function (C.2) [right axis], the initial guess for its inverse [right axis], and the (decadic logarithm of the) relative errors of the initial guess and the first two iterations [left axis].

Combining the initial guess (C.9) with the iteration procedure (C.10), the function f(x) can be inverted to standard (64 bit) floating point machine accuracy with two iterations, as is shown in figure 18. Note that the range for x in that figure is [-10, 10] and consider that the function  $\Phi(x)$  becomes indistinguishable from 1 for  $x \ge 7.5$  on standard 64 bit IEEE floating point hardward. That point is in fact easy to identify in the figure as the abscissa where the blue lined labelled "log<sub>10</sub>(|relative error of initial guess|)" denoting the decadic locarithm of the absolute value of the relative error of the initial guess suddenly levels out at about  $2 \cdot 10^{-16}$  (which is approximately the machine accuracy defined as DBL\_EPSILON). For negative x, the residual relative error is essentially just the limit of the accuracy of the cumulative normal function  $\Phi(x)$  that is employed, and this accuracy diminishes as  $x \to -\infty$ , mainly due to the loss of accuracy of the exponential function of  $\Phi(x)$ . Having said that, it is in fact possible to implement f(x) via the relationship

$$f(x) = \ln\left(\operatorname{erfcx}\left(-\frac{x}{\sqrt{2}}\right)\right) - \ln 2$$
 (C.13)

by defining the auxiliary function

$$\operatorname{lnerfcx}(x) := \operatorname{ln}(\operatorname{erfc}(x))$$
 (C.14)

and implementing lnerfcx() directly without the involvement of any avoidable exponentials or logarithms. This is straightforward<sup>5</sup> with the source code of erfcx() given by Cody [Cod69], and this is how figure 18 was produced.

<sup>&</sup>lt;sup>5</sup> With the caveat that some terms of the form  $\ln(1-\epsilon)$  will need to be expanded for small  $\epsilon$  to avoid the otherwise desastrous loss of accuracy.

**Remark C.1.** If we define the inverse function  $\xi(c)$  by the aid of the implicit function theorem formally as

$$c \to \xi = \xi(c)$$
  $\frac{1}{2}\xi^2 + \ln(\Phi(\xi)) = c$ , (C.15)

then the above semi-analytical implementation for  $\xi(c)$  given by the initial guess (C.9) combined with two iterations of (C.10) effectively gives us also an implementation for the *inverse scaled complementary error function* via the simple relationship

$$\operatorname{erfcx}^{-1}(y) = -\frac{1}{\sqrt{2}}\xi(\ln(y/2))$$
 (C.16)

This makes the result in this section useful in other applications that require  $\operatorname{erfcx}^{-1}(\cdot)$ .

#### C.1 An optimised initial guess

It is possible to derive an even better initial guess function than (C.9) by the aid of *rational Chebyshev approximations* [PTVF92] of the inverse. Without going into the details of its calculation, we give the improved version

$$x_{0}(c) = \begin{cases} \frac{1}{\eta \cdot (1 + \eta \cdot R_{1}(\eta))} & \text{when} \quad c < -2.565 \\ \sqrt{\frac{\pi}{2}} \cdot y \cdot (1 - y \cdot R_{2}(c)) & \text{when} \quad -2.565 \le c \le 7.258 \\ \sqrt{2c} & \text{when} \quad c > 7.258 \end{cases}$$
(C.17)

with

$$\eta := -\sqrt{2\pi} \cdot e^c \tag{C.18}$$

$$y := c + \ln(2) \tag{C.19}$$

$$R_1(\eta) := \frac{5.997240089442634677 \text{E} - 05 + \eta \cdot 1.003256824584751072}{1 - \eta \cdot (0.05267912997789834378 + \eta \cdot 0.2842087186525299458)} \tag{C.20}$$

$$R_2(c) := \frac{P_2(c)}{Q_2(c)} \tag{C.21}$$

 $P_2(c) := 0.2327141995027849075 + c \cdot (0.1438410800713837423 + c \cdot$ 

 $(0.03897036447332576786 + c \cdot (0.004555561056678753343 + c \cdot 0.0001955212496637105635))) (C.22)$ 

 $Q_{2}(c) := 1 + c \cdot (0.8885805957115674447 + c \cdot (0.3392936144039581658 + c \cdot (0.06608206677330506762 + c \cdot (0.006268281524752193634 + c \cdot 0.0002126843801522933914)))))$  (C.23)

which has a relative accuracy of better than  $10^{-5}$  for all input values. With this initial guess, we attain standard IEEE 64 bit floating point precision (DBL\_EPSILON) with a single iteration of the Householder(3) step (C.10). In fact, in perfect precision, a single step



FIGURE 19: The quantile slope function (C.2) [right axis], the improved initial guess (C.17) for its inverse [right axis], and the (decadic logarithm of the) relative errors of the improved initial guess (C.17) and the first iteration [left axis], all in perfect precision.

gives a relative accuracy of better than  $10^{-22}$  as is shown in figure 19. The actual accuracy you obtain depends of course as always on the noise and accuracy of the implementation of the root function f(x).

In principle, it is of course also possible to compute highly accurate rational approximations for the inverse function that require no further refinement. This, however, has the side effect that the inverse function is accurate with respect to the reference data that were used for its own construction, which may or may not be perfectly consistent with *your* implementation of the root function f(x). The consequence could be that a round-trip calculation is not a closed loop, i.e.  $f^{-1}(f(x)) \neq x$ , and this is often more important than perfect theoretical accuracy of  $f^{-1}(\cdot)$  and f(x). It is for this reason, that we favour the approach of a good initial guess with at least one "polishing" iteration, bringing the inverse in line with the root function f(x).

# References

- [AS84] M. Abramowitz and I.A. Stegun. Pocketbook of Mathematical Functions. Harri Deutsch, 1984. ISBN 3-87144-818-4.
- [Ave98] M. Avellaneda. Minimum-relative-entropy calibration of asset pricing models. International Journal of Theoretical and Applied Finance, 1(4):442–472, 1998.
- [BK96] P. W. Buchen and M. Kelly. The maximum entropy distribution of an asset inferred from option prices. Journal of Financial and Quantitative Analysis, 31(1):143–159, 1996.
- [BL78] D. T. Breeden and R. H. Litzenberger. Prices of state-contingent claims implicit in option prices. Journal of Business, 51(4):621–651, 1978.

- [Bla76] F. Black. The pricing of commodity contracts. *Journal of Financial Economics*, 3:167–179, 1976.
- [Bro08] V. Bronzin. Theorie der Prämiengeschäfte. Verlag Franz Deuticke, 1908. https://fedora. phaidra.univie.ac.at/fedora/get/o:49048/bdef:Asset/view.
- [BS73] F. Black and M. Scholes. The Pricing of Options and Corporate Liabilities. Journal of Political Economy, 81:637–654, 1973.
- [Cod69] W. J. Cody. Rational Chebyshev approximations for the error function. Mathematics of Computation, pages 631-638, 1969. www.ams.org/journals/mcom/1969-23-107/ S0025-5718-1969-0247736-4/S0025-5718-1969-0247736-4.pdf.
- [DG85] R. Delbourgo and J.A. Gregory. Shape preserving piecewise rational interpolation. SIAM Journal of Scientific and Statistical Computing, 6(4):967-976, 1985. dspace.brunel.ac.uk/ bitstream/2438/2200/1/TR\_10\_83.pdf.
- [HB04] L. Hughston and D. Brody. Entropic Calibration Revisited. Technical report, Imperial College London, 2004. www.imperial.ac.uk/research/theory/people/brody/DCB/sa10.pdf.
- [HKL02] P. Hagan, D. Kumar, and A. S. Lesniewski. Managing Smile Risk. Wilmott Magazine, pages 84–108, September 2002.
- [Jäc05] P. Jäckel. A practical method for the valuation of a variety of hybrid products. In ICBI Global Derivatives Conference Paris, 2005. www.jaeckel.org/ APracticalMethodForTheValuationOfAVarietyOfHybridProducts.pdf.
- [Jäc13] P. Jäckel. Let's be rational. www.jaeckel.org/LetsBeRational.pdf, November 2013.
- [Kah04] N. Kahalé. An arbitrage-free interpolation of volatilities. Risk, pages 102–106, May 2004. www.risk.net/data/Pay\_per\_view/risk/technical/2004/0504\_tech\_option2.pdf.
- [Lee04] R. Lee. The Moment Formula for Implied Volatility at Extreme Strikes. *Mathematical Finance*, 14:469–480, 2004.
- [Mar04] G. Marsaglia. Evaluating the Normal Distribution. Journal of Statistical Software, 1:1-11, 2004. www.jstatsoft.org/v11/a05/paper.
- [PTVF92] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery. Numerical Recipes in C. Cambridge University Press, 1992. www.nrbook.com/a/bookcpdf.php.