Ball arithmetic in MPC

Andreas Enge

CANARI project-team INRIA Bordeaux andreas.enge@inria.fr http://enge.math.u-bordeaux.fr/

MPFR/MPC/MPFI/FLINT Workshop, Bordeaux, June 17-19, 2024



naín

Dagstuhl 2018: New number type to implement

Complex ball arithmetic

- Used internally to replace hand calculations of error bounds
- As a tool to implement Taylor and Laurent series
- As a building block for polynomials, class polynomials, etc.
- Only four basic arithmetic operations and square root?
- Need for real ball arithmetic? interface with MPFI?
- Representation as rectangles? as balls?



Dagstuhl 2018: New number type to implement

Complex ball arithmetic

- Used internally to replace hand calculations of error bounds
- As a tool to implement Taylor and Laurent series
- As a building block for polynomials, class polynomials, etc.
- Only four basic arithmetic operations and square root?
- Need for real ball arithmetic? interface with MPFI?
- Representation as rectangles? as balls?

Release 1.3.0 "Ipomoea batatas", December 2022:

- Ball arithmetic (marked experimental)
- New function mpc_agm
- New function mpc_eta_fund



If any of r_1 , r_2 , $r_2 < r_2$ graphs 0, then the significant has ofther finished in the second trap, or the corresponding near its compared exactly always up reveals $a_2 < r_1$, a contradiction. So the case of interest is that these four input values are non-zero. Then the $z_1 = 2^{-1}(r_1^2 + 2^{-1}q_1^2)$ with integers r_1^2 and f_1 with remove $\delta_1 > 0$ since $y_1 > z_1$. We have $u_1 = 2^{2-1} \phi(u_1)$ with the integer u_1 and u_1 is compared in the range $u_1 = r_1 + 2^{-1} \phi(u_1)$ with the integer $u_1 = r_1 + 2^{-1} \phi(u_1)$ with the integer $u_1 = r_1 + 2^{-1} \phi(u_1)$.

So $d_1, d_2 > p$. The integer m_{ee} , and from the lowest to the most significant time solution of $(q_1^{(2)} = \cos d \sin 2q)$ with $(q_2^{(2)} = \sin d \sin 2q)$ and $(q_2^{(2)} = \sin d \sin 2q)$. (b) followed by $d_2 = 2p > 10$ to its followed by $(q_2^{(2)} = \sin d \sin 2q)$ with 2p that (with growthard) can be diagonally a significant of $(q_1^{(2)} = \sin q)$. (b) $(q_2^{(2)} = q)$, $(q_2^{(2)} = q)$, $(q_2^{(2)} = q)$, $(q_2^{(2)} = a)$,

3.10 mpc_agm

ZECCARD

Definition. Let a_i be incover complex numbers. Define sequences of arithmetic mans (a_i) and groundrift manus (b_i) and groundrift manus (b_i) and groundrift manus (b_i) and groundrift manus (b_i) and groundrift $(b_i) = b_{i+1}^{k-1} a_{i+1} b_{i+1} - b_{i+1}^{k-1} a_{i+1} b_{i+1} - b_{i+1}^{k-1} a_{i+1} b_{i+1} - b_{i+1}^{k-1} a_{i+1} b_{i+1}$. Unconstant) anguing different from 0 and a_i . The there f diff are two-dimensional pointer dimensional groundrift and groundrift (b_i) and (b_i) and (

It is immediate that AGM is symmetric, that is, AGM(a, b) = AGM(b, a), and homogeneous, that is, $AGM(\lambda a, bb) = \lambda AGM(a, b)$ for any non-zero complex number λ . So we may assume that $|a| \ge |b|$, and AGM(a, b) = a AGM(a, bb) with $a_0 = 1$.

b) we may assume that $|a| \ge |b|$, and $ACOU(a, b) = aACOU(a_0, b_0)$ $b_0 = b/a, |b_0| \le 1.$

We need to examine the corner cases. It one or both of a and b are zero, all geometric means are zeros, and $AGM_{(h, h)} = 0$. If the magnle between an and b is (thus his a positive real number, and $AGM_{(h, h)}$) may be computed with $m_{\rm FF}$. If the angle between and b is (thus his ||a| = 1, 0, 0). The arithmetic mean of 1 and -1 is over, as $AGM_{(h, -1)} = 0$. $AGM_{(h, 0)} = 0$. If the state of the state

So in the following, we analyse the computation of AGM(1, \tilde{b}_0) with $\tilde{b}_0 = o(b_0)$ for b_0 in the unit disk centered at the origin (except - 1 and 1), where the real and imaginary parts of \tilde{b}_0 are rounded towards 0, which ensures $|\tilde{b}_0|_{\delta_0}| < 1$ and $|\tilde{b}_0| < 1$ with absolute errors of the real and imaginary parts at most 11 up. In the following, relative errors will

31

be most convenient to work with; by Proposition 10 we have relevror $(\tilde{b_0}) \leq 2^{1-p}$, where $p \geq 2$ is the working precision.

Warm-up — a note on angle 0. As said, when the angle between a and b is 0, which is immediately detected from $\Im(\tilde{h}_0) = 0$ (computed at any precision), then $\operatorname{ACM}(a, b) =$ $a\operatorname{ACM}(1, b_0)$ with a real number $0 < b_0 < 1$. We need to consider the error induced by computing $(a \circ (\operatorname{ACM}(1, \tilde{b}_0)))$ instead.

Recall that $b_0 = (1 + \vartheta) \delta_0$ with $-\varepsilon \leqslant \vartheta \leqslant \varepsilon = 2^{1-p}$. Since for positive real numbers each step of the AGM is increasing in each of the two arguments, so is the AGM itself, and we have

$$AGM(1, b_0) = AGM\left(1, (1 + \vartheta)\tilde{b}_0\right) \leq AGM\left(1 + \varepsilon, (1 + \varepsilon)\tilde{b}_0\right) \leq (1 + \varepsilon) AGM\left(1, \tilde{b}_0\right)$$
.

By the same kind of argument

$$AGM(1, b_0) \ge (1 - \varepsilon) AGM(1, \tilde{b_0})$$
.

So the relative error of at most ε in the input \tilde{b}_0 is preserved by the AGM.

By Proposition 7 applied to non-representable numbers and Proposition 3 this relative error translates into an absolute error satisfying

 $|\operatorname{AGM}(1, b_0) - \operatorname{AGM}(1, \tilde{b_0})| \le 2 \cdot 2^{\operatorname{Exp}(\operatorname{AGM}(1, \tilde{b_0})) - p} \le 2 \cdot 2^{\operatorname{Exp}(\circ(\operatorname{AGM}(1, \tilde{b_0}))) - p}$

of at most 2 ulp of the rounded value, and the final rounding of the mpfr_agm function leads to a total error bounded by 3 ulp.

We multiply this value by the exact a; applying (2) with $k_1 = 3$ and $k_2 = 0$ and taking the final rounding into account leads to an error of at most 7 ulp for the real and for the imaginary parts of the complex AGM.

The first iteration — entering a quadrant. If $R(\tilde{b}_0) < 0$, then significant cancellation can occur for the arithmetic mean in the first iteration, which thus needs to be analysed separately.

From now on, we use arbitrary rounding modes and apply Proposition 11 with c = 1. We let $b_1 = \sqrt{b_0}$ and $\tilde{b}_1 = o\left(\sqrt{\tilde{b}_0}\right)$ with

$$\operatorname{relerror}(\overline{b_1}) \leq 2^{1-p} + (1 + 2^{1-p}) 2^{1-p} \leq (\sqrt{2})^3 \cdot 2^{1-p}$$

by (18) (where we use $\leq \epsilon_1 \lesssim \sin \epsilon_2 \lesssim 2^{1-\rho} \lesssim 1, \epsilon_1$ being the relative error on $\frac{h}{h}$) and Proposition 11, and where we have bounded $2 + 2^{1-\rho} \lesssim 2.5$ by $(\sqrt{2})^3 \approx 2.83$. We let $\alpha_1 = 1/2 + h_0/2$ and $\hat{\alpha}_1 = c(1/2 + \hat{h}_0/2)$. The imaginary part of $\hat{\alpha}_1$ has an error of at most 1 hap, and the same holds for the real part if $\Re(h_0) = 0$ or equivalently $\Re(h_0) = 0$; the real part of $\hat{\alpha}_1$ has an absolute error bounded by

 $\max\left(2, 2^{-\operatorname{Exp}(\Re(\widetilde{a_1}))-1}\right) \operatorname{ulp}\left(\Re(\widetilde{a_1})\right)$.



Andreas Enge

Ball arithmetic in MPC

Indeed, let z be the value z_{0} of z_{0} and y the real put of G_{1} , we have y=G(1/2+z/2). Remember that -1<e<1 b. I-1/2<<2< b. I-1/2<<2, I-1/2. Indeed, z=0<-1 the -1/2<<2, I-1/4 and y>1/2. Als and Eqs(y) 1. Segmet the theorem on y is at most 2 day is not found from the three reasons in $2/z_{1}$ and so the three haddings. It is the ender the hand -1<z<2, I/2, then 1/2+z/2, then 1/2>, I/2 is the sequence of the sequ

$$\operatorname{relerror}(\tilde{a}_1) \leq \max\left(2, 2^{-\operatorname{Exp}(\Re(\tilde{a}_1))-1}\right) \cdot 2^{1-p}.$$

If $\Re(b_0) \ge 0$, then $Exp(\Re(\tilde{a}_1)) \ge 0$, and we obtain a bound that is similar to the one above for referror(\tilde{b}_1). If \tilde{b}_0 approaches -1, the relative error in \tilde{a}_1 becomes arbitrarily bad, as \tilde{a}_1 becomes arbitrarily small.

Letting

$$k_1 = \max (3, -2 \exp(\Re(\tilde{a}_1)) - 2$$

(24)

(25)

we obtain an upper bound of

$$\operatorname{relerror}(\widetilde{a_1}), \operatorname{relerror}(\widetilde{b_1}) \leq (\sqrt{2})^{n_1} \cdot 2$$

on the loss of precision in the first iteration.

Notice that, independently of the rounding mode, \hat{a}_1 and \hat{b}_1 lie in the same complex quadrant of numbers having non-negative real part and an imaginary part of the same sign as that of b_0 (c of positive imaginary part of b_0 is real). During the remainder of the algorithm, we will not leave this quadrant and thus not see any more cancellation in the arithmetic mean.



Figure 1: Extremal (from above) values of a_1 (solid curve) and b_1 (dotted curve) after the first iteration of the AGM for $\Im(b_0) \ge 0$. The values are bounded from below by the *z*-axis. The bias of the analysis. Let us find give the basic blass of the following, rather technic and analysis. Assume a transport periodical of N kin has a karger state term of about quadratic and the state of the sta

Unlike Newton Iterations, AGM Iterations are not auto-correcting: due to rounding errors, we lose a constant number of bits per iteration; so to reach the desired precision of N bits, we need to carry out all computations at a working precision of $p \in N + O\left(B(N, b_0, \bar{\alpha}_1)\right)$. The following discussion provides explicit bounds for all these quantities:

Rounding error propagation. Let $a_n = \frac{a_{n-1}+b_{n-1}}{2}$, $c_n = a_{n-1}b_{n-1}$, $b_n = \sqrt{c_n}$. The computation of \hat{a}_1 and \hat{b}_2 and their error analysis have been given above. For $n \ge 2$, we compute the sequences

$$\widetilde{a}_{n} = \circ \left(\frac{\widetilde{a_{n-1}} + \widetilde{b}_{n-1}}{2} \right),$$

 $\widetilde{c}_{n} = \circ \left(\widetilde{a_{n-1}} \widetilde{b}_{n-1} \right),$
 $\widetilde{b}_{n} = \circ \left(\sqrt{\widetilde{c}_{n}} \right).$

Then one sees by induction that $\tilde{a_n}$ and $\tilde{b_n}$ lie in the same quadrant as $\tilde{a_1}$ and $\tilde{b_1}$ (or, for that matter, a_1 and b_1).

Let $\alpha_n = relerror(\tilde{a_n})$, $\gamma_n = relerror(\tilde{c_n})$ and $\beta_n = relerror(\tilde{b_n})$. By (7) and Proposition 11,

$$\alpha_n \le \sqrt{2} \max(\alpha_{n-1}, \beta_{n-1}) (1 + 2^{1-p}) + 2^{1-p}$$
. (26) eq: agnalph

By (10) and Proposition 11,

 $\gamma_n \leq (\alpha_{n-1} + \beta_{n-1} + \alpha_{n-1}\beta_{n-1})(1 + 2^{1-p}) + 2^{1-p}.$ (27) eq: agregamma

By (18) and Proposition 11,

$$\beta_n \leq \frac{1}{\sqrt{3}} \gamma_n (1 + 2^{1-p}) + 2^{1-p} \text{ if } \gamma_n \leq \frac{1}{4}.$$
 (28) eq: sgnber

Let $r_n = ((\sqrt{2})^{n+1} - 1) r_1$ for some integer r_1 such that $\alpha_1, \beta_1 \leq r_1 2^{1-p}$. We may use $r_1 = (\sqrt{2})^{k_1}$ with k_1 as in (24), so that $r_n \leq (\sqrt{2})^{n+k_1+1}$. We now show by induction



that $\alpha_n, \beta_n \leq r_n 2^{1-p}$ if the number n of iterations is not too large compared to the working precision p. For n = 1, this follows from (25) and the definition of r_1 . As a preparation for the induction step and taking the shape of (26) to (28) into account, we carry out the following computation for $n \ge 2$, $p \ge 2$, $\frac{290(201)}{512} \le \gamma \le 2$ (we will only need $\gamma \in \{\sqrt{2}, 2\}$ later) and $0 \le \delta \le 1$ (we will only need $\delta \in \{0, 1\}$):

$$\begin{split} R(n, \gamma, \delta) &:= \left(\gamma r_{n-1} + \delta r_{n-1}^2 2^{1-p}\right) \left(1 + 2^{1-p}\right) + 1 \\ &= \left(\gamma + \delta r_{n-1} 2^{1-p}\right) r_{n-1} + 1 + 2^{1-p} \left(\gamma r_{n-1} + \delta r_{n-1}^2 2^{1-p}\right) \\ &\leqslant \left(\gamma + \delta(\sqrt{2})^{a+4z+2-2p}\right) r_{n-1} + 1 + 2^{1-p} \left((\sqrt{2})^{a+4z+2} + 2^{a+4z+1-p}\right). \end{split}$$

So assuming $p \ge \frac{n+k_1+10}{2}$ we have

$$R(n, \gamma, \delta) \le (\gamma + \delta(\sqrt{2})^{-8})r_{n-1} + 289/256$$

$$\leq \frac{(\gamma + \delta(\sqrt{2})^{-n})}{\sqrt{2}}r_n + 289/266 - \frac{(\gamma + \delta(\sqrt{2})^{-n})}{\sqrt{2}}(r_n - \sqrt{2}r_{n-1}) \\ \leq \frac{(\gamma + \delta(\sqrt{2})^{-n})}{\sqrt{2}}r_n + 289/266 - \frac{(\gamma + \delta(\sqrt{2})^{-n})}{\sqrt{2}}(\sqrt{2} - 1)(\sqrt{2})^2 \\ \sin cr_1 \geq (\sqrt{2})^4 \\ \leq \frac{(\gamma + \delta_{10}/2)}{(\gamma + \delta_{10}/2)}r_n \sin c\delta \geq 0 \text{ and } \gamma \geq \frac{289(\sqrt{2} + 1)}{512}.$$

Using the induction hypothesis for $n \ge 2$ on (26) and (27) yields

$$\alpha_n \leqslant R(n, \sqrt{2}, 0) 2^{1-p} \leqslant r_n 2^{1-p},$$

 $\gamma_n \leqslant R(n, 2, 1) 2^{1-p} \leqslant \sqrt{2} \cdot \frac{33}{32} \cdot r_{n-1} 2^{1-p},$
 $\leqslant \sqrt{6} r_{n-1} 2^{1-p},$
 $\leqslant 4 \cdot (\sqrt{2})^{n+k_1} \cdot (\sqrt{2})^{-(a_1+k_1+6)} = 1/4.$

So (28) is valid, and substituting γ_0 by its upper bound $\sqrt{6}r_{n-1}2^{1-p}$ yields

$$\beta_n \le R(n, \sqrt{2}, 0) 2^{1-p} \le r_n 2^{1-p}$$
.

To summarise, we compute in iteration n approximations $\tilde{a_s}$ and $\tilde{b_s}$ to a_s and b_s with a relative error bounded above by

$$2^{\frac{n+k_2+3}{2}-\rho}$$
 assuming that the working precision sat

tisfies
$$p \ge \frac{n + k_1}{n} + 5$$
. (29) eq: agapter

Mathematical error. We also need to estimate the error made by carrying out only a finite number of iterations. Let z satisfy $\Re(z) \ge 0$ and $z \ne 0, 1$, and consider the optimal 35

AGM sequences (a'_n) and (b'_n) computed (with infinite precision) from $a'_n = 1$ and $b'_n = z$. Let $N' \in \mathbb{N}$ and

$$n \ge B'(N', z) := \max \left(1, \left\lceil \log_2 \left| \log_2 |z| \right\rceil\right) + \left\lceil \log_2(N' + 2) \right\rceil + 2$$

(where log₂ 0 is to be understood as -∞). Then by [6, Prop. 3.3, p. 88].

$$a'_{\alpha} = (1 + \vartheta_2) \operatorname{AGM}(1, z)$$
 with $|\vartheta_2| \le 2^{-N'}$.

(Notice that here, the relative error as defined in [6, Def. 1.2, p. 20] is taken with the roles of the correct and the approximated values reversed compared to our definition.)

In our context, we may have $\Re(b_0) < 0$, but after one iteration, $\Re(b_1/a_1) \ge 0$. So we consider $z = b'_0 = b_1/a_1$, so that by homogeneity, $a_{n+1} = a_1a'_n$, $b_{n+1} = a_1b'_n$ and $AGM(1, b_0) = AGM(a_1, b_1) = a_1 AGM(1, z)$. Thus we need bounds on

$$|z| = \left|\frac{b_1}{a_1}\right| = \frac{|\sqrt{b_0}|}{|a_1|}.$$

Since we have to decide on the number of iterations from approximations to b_0 and a_1 instead of the correct values, we may wish to replace b_0 by b_0 and a_1 by $\tilde{a_1}$. (In fact the arguments will be quite delicate and switch between the different quantities.) The relative errors derived above on \tilde{b}_{0} of 2^{1-p} and on \tilde{a}_{1} of $(\sqrt{2})^{k_{1}} 2^{1-p}$ with $p \ge k_{1}/2 + 5 \ge 5$ show that

$$\left(1-2^{-4}\right)\left| \widetilde{b_0} \right| \leqslant \left| b_0 \right| \leqslant \left(1+2^{-4}\right) \left| \widetilde{b_0} \right| \ \text{and} \ \left(1-2^{-4}\right) \left| \widetilde{a_1} \right| \leqslant \left| a_1 \right| \leqslant \left(1+2^{-4}\right) \left| \widetilde{a_1} \right| .$$

Since $\tilde{b}_0 = o(b_0)$ with both parts rounded towards zero, we even have $|\tilde{b}_0| \le |b_0|$. Altogether we have the (much coarser) following bounds:

$$|\tilde{b}_0| \le |b_0| \le \sqrt{2} |\tilde{b}_0|$$
 and $\frac{1}{\sqrt{2}} |\tilde{a}_1| \le |a_1| \le \sqrt{2} |\tilde{a}_1|$. (30) eq: agnabatic equation (31) (32)

The quantity of interest is a (double) logarithm, so it is helpful to consider the (easily available) exponents of the numbers. Since $\tilde{b}_0 = o(b_0)$ with both parts rounded towards zero we have

 $Exp(\Re(\tilde{b}_0)) = Exp(\Re(b_0))$ and $Exp(\Im(\tilde{b}_0)) = Exp(\Im(b_0))$.

(31) eq:agnempt

Recall that $\tilde{a}_1 = \circ \left(\frac{1+\tilde{a}_1}{n}\right)$ so that $\Im(\tilde{a}_1) = \Im(\tilde{b}_0)/2$ does not require any additional rounding: together with (31) this yields $Exp(\Im(\tilde{a}_1)) = Exp(\Im(\tilde{b}_0)) - 1 = Exp(\Im(b_0)) - 1 = Exp(\Im(a_1)).$

However $\Re(\tilde{a}_1) = o(1 + o(\Re(b_0)))/2$ is computed with two roundings, so it takes a bit more work to relate it to $\Re(a_1) = (1 + \Re(b_0))/2$, which we postpone to the following case distinction since we do not need a completely general result. 36



Andreas Enge

Ball arithmetic in MPC

The above bound B^{ϵ} for the required number of iterations n grows with $|\log_2 |z||$, so problematic situations may occur when |z| is either very large or very small; the former may happen when $|b_0|$ is large and/or $|a_1|$ is small, the latter when $|b_0|$ is small and/or $|a_1|$ is large. Remember that in our setting we have the bounds

$$|b_0| \le 1$$
, so that $|a_1| = \left|\frac{b_0 + 1}{2}\right| \le 1$.

thus interesting situations can occur only when either $|b_0|$ is small or $|\alpha_1|$ is small, that is, b_0 is close to -1; and both cannot happen simultaneously. This leads to the following case distinction:

(a) Assume Exp(ℜ(b˜₀)), Exp(ℑ(b˜₀)) ≤ −1.

By (31) this is equivalent to $Exp(\Re(b_0)), Exp(\Im(b_0)) \leq -1$, whence $\Re(b_0), \Im(b_0) < 1/2$ and

$$|b_0| \le \frac{\sqrt{2}}{2} < 1.$$

This implies

$$1 - \frac{\sqrt{2}}{2} \le 2|a_1| \le 1 + \frac{\sqrt{2}}{2}$$
 and $1/8 \le |a_1| \le 1$

and we know $|\sqrt{b_0}| \ge |\sqrt{b_0}|$ by (30). Collecting these upper and lower bounds, we obtain for $z = \sqrt{b_0}/a_1$ that

$$\sqrt{\left| \widetilde{b_0} \right|} \leqslant |z| \leqslant 8 \text{ and } \left| \log_2 |z| \right| \leqslant \max \left(3, -\frac{1}{2} \log_2 \left| \widetilde{b_0} \right| \right).$$

Using the estimate

$$|\tilde{b}_0| \ge \max \left(|\Re(\tilde{b}_0)|, |\Im(\tilde{b}_0)| \right) \ge 2^{\max\left(\operatorname{Exp}(\Re(\tilde{b}_0)), \operatorname{Exp}(\Im(\tilde{b}_0)) \right) - 1}$$

we finally obtain

$$|\log_2 |z|| \le \max \left(3, -\frac{1}{2} \max \left(\operatorname{Exp}(\Re(\widetilde{b_0})), \operatorname{Exp}(\Im(\widetilde{b_0}))\right) + \frac{1}{2}\right).$$
 (33)

(b) Assume Exp(ℑ(b̃₀)) ≤ −1 and Exp(ℜ(ã₁)) ≤ −2.

By (12) the first condition implies $\text{Exp}(2(\alpha_0)) \leq -2$ or $|2(\alpha_1)| \leq 1/4$. The second condition together with the first inequality of Proposition 3 implies that $\text{Exp}\left((1 + \Re(b_0))/2\right) \leq -2$, or $\left[1 + \Re(b_0)\right] < 1/2$. So in fact $-1 < \Re(b_0) < -1/2$, and then $-1 < \Re(b_0) < -1/2$ are well. This in turn implies $0 < \Re(a_1) < 1/4$. Putting these together, we obtain

$$1/2 \le |\Re(b_0)| \le |b_0|$$
 or $1/\sqrt{2} \le |\sqrt{b_0}|$, and $|a_1| \le \sqrt{2}/4$,

and we already know $|b_0|\leqslant 1 \text{ or } |\sqrt{b_0}|\leqslant 1$

from our general setting and

$$|a_1| \ge |a_1|/\sqrt{2}$$

by (30), whence

$$2 \le |z| \le \frac{\sqrt{2}}{|\tilde{a_1}|}$$
 and $|\log_2 |z|| \le -\log_2 |\tilde{a_1}| + \frac{1}{2}$.

Using the estimate

$$|\tilde{a}_1| \ge \max(|\Re(\tilde{a}_1)|, |\Im(\tilde{a}_1)|) \ge 2^{\max(\exp(\Re(\tilde{a}_1)), \exp(\Im(\tilde{a}_1)))-1}$$

and (32) we finally obtain

$$|\log_2 |z|| \le -\max \left(\exp(\Re(\tilde{a}_1)), \exp(\Im(\tilde{b}_0)) - 1 \right) + \frac{3}{2}.$$

(34) eq:agnla

(c) In the remaining case, since we are not in (a), at least one of Exp (ℜ(δ₀)) and Exp (ℑ(δ₀)), or equivalently by (31) at least one of Exp(ℜ(b₀)) and Exp(ℑ(b₀)) is 0 or larger, so that

$$\frac{1}{2} \leq \max(|\Re(b_0)|, |\Im(b_0)|) \leq |b_0| \leq 1 \text{ and } \frac{1}{\sqrt{2}} \leq \sqrt{|b_0|} \leq 1.$$

If $Exp\left(\Im(\tilde{b}_0)\right) \ge 0$, then $Exp(\Im(a_1)) \ge -1$ by (32) and

$|a_1| \ge |\Im(a_1)| \ge 1/4.$

Otherwise since we are not in (b), $Exp(\Re(\tilde{n})) \ge -1 \circ [\Re(\tilde{n})] \ge 1/4$, as it is known to be positive in fact $\Re(\tilde{n}) \ge 1/4$. In precision at heats 1, the values $1/4 \cdot 31/32$ is representable and smaller than 1/4 so that the value rounded to $\Re(\tilde{n}_1)$ satisfies $(1 + \Re(\tilde{n}_2)) < -33/64$. In precision at text S, the value -34/64 is representable and smaller than -33/64 so that the unrounded value satisfies $\Re(n) > -17/52$, which implies $\Re(n) > 157/64 > 1/8$ and

$|a_1| \ge |\Re(a_1)| \ge 1/8.$

We also know $|a_1| \leq 1$ in our general setting. Altogether

$$\frac{1}{\sqrt{2}} \le |z| \le 8$$
 and $|\log_2 |z|| \le 3$. (35) eq:age

38



Letting $L(\tilde{b}_0, \tilde{\alpha}_1)$ denote the above bound on $|\log_2|z||$ depending on the exponents occurring in \tilde{b}_0 and \tilde{a}_1 , and counting the additional first iteration to compute $\tilde{\alpha}_1$ and \tilde{b}_1 , not present in the above analysis, we fix a number of iterations n such that $n \ge B(N, \tilde{b}_0, \tilde{\alpha})$ with

Then

$$B(N, \tilde{b}_0, \tilde{a}_1) = \max \left(1, \left\lceil \log_2 L(\tilde{b}_0, \tilde{a}_1) \right\rceil \right) + \left\lceil \log_2(N + 4) \right\rceil + 3.$$

 $a_- = (1 + d_0) MGM(1 |b_0|) \text{ with } |d_0| \le 2^{-(N+2)}$

(35)

Total error and working precision. Combining with (29), we obtain for a sufficiently large precision p and sufficiently many iterations or that $AGM(1, b_0) = \frac{|k+2p|}{1+2p} \widetilde{\alpha}_n = (1+\delta)\widetilde{\alpha}_n$ with $|\theta_1| \leq 2^{-1(k+2-2)}$, so that

$$|\vartheta| \le \frac{|\vartheta_1| + |\vartheta_2|}{|1 - |\vartheta_2||} \le \frac{4}{3} \left(2^{\frac{n+k_1+3}{2}-p} + 2^{-(N+2)}\right)$$

for $N \ge 2$. So after $n = B(N, \widetilde{b_0}, \widetilde{\alpha}_1)$ steps of the AGM iteration at a working precision of $p \ge N + \frac{kk_0^{k-2}}{2} \ge \frac{kk_0^{k-2}}{2}$, we obtain $\widetilde{\alpha_n}$ which approximates AGM $(1, b_0)$ with a relative error bounded by $\frac{2}{2} \cdot 2^{-N}$.

Finally, we let $z = AGM^{-}(a, b) = a AGM(1, b_0)$ and $\tilde{z} = \circ(a\tilde{u}_{a})$, where a is known exactly. By (10) and Proposition 11, using that $N \ge 2$ and $k_1 \ge 3$ imply $n = B(N, \tilde{b}_0, \tilde{a}_1) \ge 7$ $ad p \ge N + 8 \ge 10$, this leads to a relative error bounded by 2^{-N} .

Summary. In our analysis, the working precision p depends on

 $k_1 = \max(3, -2 \operatorname{Exp}(\Re(\tilde{a}_1)) - 2)$

of (13), which is turn depends on only on the input data, but also in the working reprised and the first ACM trends. It is to complicate any of the physical computation at authinoidy to precision. Since h_0^{-1} is compared as $a^{-1}b$ with a depending on the respection of the experimental sector h_0^{-1} and h_0^{-1} and h_0^{-1} and h_0^{-1} and h_0^{-1} the expense of $\Re(h_0^{-1})$ may only be wrong if $\Re(h_0^{-1})$ is one fixed in b^{-1} that it is rounded near any other of $\Re(h_0^{-1})$ may only be wrong if $\Re(h_0^{-1})$ is one fixed h^{-1} with Π^{-1} and Π^{-1} Π^{-1} and Π^{-1} Π^{-1} and Π^{-1} Π^{-1} and Π^{-1} Π^{-1} and Π^{-1} Π^{-1} and Π^{-1} and Π^{-1} and Π^{-1} and Π^{-1} and Π^{-1} and Π^{-1} Π^{-1} and Π^{-1} Π^{-1} and $\Pi^$

Then we fix a desired accuracy N (around the target precision plus a safety margin), compute $L(\tilde{0}_0, \tilde{n}_1)$ by (33), (34) or (35) and the number of iterations $n = B(N, \tilde{0}_0, \tilde{n}_1)$ by (36) (more often than not, this will result in $n = \lceil \log_2 N \rceil + 5$). Then the working precision is given by

$$p = N + \left\lceil \frac{n + k_1 + 7}{2} \right\rceil$$

Using Propositions 9 and 7, the complex relative error of 2^{-N} may be translated into an error expressed in ulp. With $\tilde{x} = \tilde{x} + i\tilde{y}$ the computed approximation of z =AGM(a, b), let $k_R = \max(\text{Exp}(\tilde{y}) - \text{Exp}(\tilde{x}) + 1, 0) + 1$, and $k_I = \max(\text{Exp}(\tilde{x}) - \text{Exp}(\tilde{y}) + 1, 0) + 1$. Then we have error($\tilde{x}) \leq 2^{4n+p-N} ulp(\tilde{x})$ and error($\tilde{y}) \leq 2^{4n+p-N} ulp(\tilde{x})$.

In practice, one should take this additional loss into account. If rounding fills after the first computation, nevertheless the values of k_R and k_I will most likely not change with a larger precision. So one should let $k' = \max(k_R, k_I)$, replace N by N + k' and adapt the number of iterations and the working precision accordingly.

The number of iterations is also slightly pessimistic in practice, in particular the additional constant in (36). So the computations may be stopped earlier if the numbers occurring in the AGM iterations do not change any more, since then additional iterations will fix the result.

4 Complex ball arithmetic

We propose a simple implementation of complex balls, which here track of rounding errors over sever) operations. The originality of our implementation is that it uses complex relative errors as in §1.1.4. A complex ball of type mpch z_1 is defined by a nonzero centre of type mpc z_1 and a relative radius of type mpch z_1 is defined by a nonzero centre of type mpc z_1 and a relative radius of type mpch z_1 is defined by a complex numbers $z_1 = (1 + \eta)$ with $|0| \leq r$, or equivalently the closed circle with centre cand radius r(c) is the following: we see the notation (c - 1) for this complex ball.

The radius type represents the non-negative real half-axis from 0 to ∞ , including this special infinite number. It is implemented internally as a non-negative baseling point value with a signed 64 bit integer matrixes m_i normalised to 31 bits, so that two matrixes may be multiplied without rounding. This sign is used on encode $\tau = \infty$, a raw matrixes accels 0, otherwise a matrixes is always positive. The exponent is encoded as a 64 bit integer a with all (for finite radii) $\tau = m_i^2$ with $m_i^2 = m_i^2 = 2m_i^2$ in a most applications radius $\tau \ge 3$ will be meaningless, so that in practice we will almost always have $\epsilon < -1$, are $\tau = \infty$.

Mathematical functions are then understood to work on sets and to "round up": They return a complex hall containing the set obtained by applying the function to every combination of arguments from the input balls. Reasonable efforts are made to return small balls, but there is no guarantee that the returned ball is minimal.

In any case, the centre of the resulting ball is obtained by applying the corresponding MPC function on the centres of the input balls with rounding to nearest. The analyses of §3 then provide upper bounds on the radius.

Compared to a representation decomposed along the real and imaginary parts, with separate relative can absolute errors, which leads to retracking instead of circles, our representation simplifies multiplicative operations and makes additive operations more complicated. Diversing in read-custs, which show deposes on the decomposition of a constant of the relative structure of the relative structure of the second tracking operation of the relative structure of the relative constant of the largest drawline is that interval control in 0 may not be represented at all of these purersity, interval containing on an almost network they correspond to balls

Innía

Andreas Enge

Ball arithmetic in MPC

Algorithm for AGM using ball arithmetic

with $r \ge 1$, which means that even the most significant digit of the centre is uncertain.

4.1 Crossing the axes



To correctly evaluate functions at branch cuts, it may be useful to examine how complex balls are positioned with respect to the axes. As already mentioned above, a ball (c,r) contains the origin if and only if $r \ge 1$.

For c = x + iy, it crosses (or touches) the real axis if and only if it contains the point x, which means that

 $|x - c| \le r|c| \Leftrightarrow |y| \le r|c| \Leftrightarrow y^2 \le r^2 (x^2 + y^2) \Leftrightarrow (1 - r^2) y^2 \le r^2 x^2$.

When does a complex ball work there exists in 25 m issues that $(1-r)^2 \mu^2 = r^2 \sigma^2$, which approxolotion by when r = 1 and $\sigma = 0$, and $\mu = 0$, and the ball is a point coiled on the real axis, or when r = 1 and $\sigma = 0$, and the the ball is control on the imaginary axis and axis, $(\sigma + \mu) = 0$, and $(\sigma + \mu) = 0$. The second one is the second on the imaginary main structure is the second one is the

Symmetrically, the complex ball (x + iy, r) crosses (or touches) the imaginary axis if and only if

$$(1 - r^2)x^2 \le r^2y^2$$
.

It touches the imaginary axis if and only if r = 0 and x = 0, so that the ball is a point on the imaginary axis; or r = 1 and y = 0, so that the ball is centred on the real axis and touches the origin.

The ball has a common point with the negative real axis (including the origin) if and only if either $x \leq 0$ and the ball has a common point with the real axis (since then x is such a common point); or x > 0 and the ball contains the origin.

4.2 mpc_agm

Implementing the AGM (see §3.10) for complex balls would require to check whether the input crosses the negative real axis, where we have placed the branch cut (which is inherited from the branch cut of the complex square root). However if the input is a complex number, which can be considered to be exact, then an implementation using complex halls can obtain a correctly rounded result with a greatly simplified analysis compared to §3.10.

In a first step, assuming that $|a| \ge |b|$, we compute $b_0 = b/a$ as a complex ball centred around $\tilde{b_0}$.

If $3(\frac{1}{30}) = 0$, then the maple between *z* and *b* is 0 or *z*. Regardless of the rounding directions, $\Re(b_1)$ also the same size as *w* in $\Re(b_1) = \Re(b_1) - 2\Re(b_1) -$

After n iterations starting from 1 and a ball around b_0 , we end up with balls $(\widetilde{\alpha_n}, r_n)$ and (b_n, r_1) such that the exact values satisfy $a_n = (1 + \vartheta_n)\widetilde{\alpha_n}$ with $|\vartheta_n| \leq r_n$ and $b_n = (1 + \vartheta_n)\widetilde{b_n}$ with $|\vartheta| \leq r_1$.

By [6, p.87] we have $|AGM(1, b_0) - a_n| \le |a_n - b_n|$. Plugging the expressions for a_n and b_n into this inequality yields

 $|\operatorname{AGM}(1, b_0) - (1 + \vartheta_a)\widetilde{a_b}| \leq |\widetilde{a_b} - \widetilde{b_b}| + |\vartheta_a\widetilde{a_b}| + |\vartheta_b\widetilde{b_b}|$

The lower triangle inequality gives

$$|\operatorname{AGM}(1, b_0) - (1 + \vartheta_{\hat{a}})\widetilde{a_0}| \ge |\operatorname{AGM}(1, b_0) - \widetilde{a_{\hat{a}}}| - |\vartheta_{\hat{a}}| \cdot |\widetilde{a_{\hat{a}}}|,$$

and putting these inequalities together we obtain

$$|\operatorname{AGM}(1,b_0) - \widetilde{a_n}| \leqslant \left(\left| \frac{\widetilde{a_n} - \widetilde{b_n}}{\widetilde{a_n}} \right| + 2|\vartheta_s| \right) |\widetilde{a_n}| + |\vartheta_b| \cdot |\widetilde{b_n}|$$

Write $\tilde{b_n} = (1 + \vartheta_{a,b})\tilde{a_n}$ and $r_{a,b} = |\vartheta_{a,b}|$; then we obtain

 $|AGM(1, b_0) - \widetilde{a_n}| \le (r_{a,b} + 2r_a + r_b(1 + r_{a,b}))|\widetilde{a_a}| \le (2(r_{a,b} + r_a) + r_b)|\widetilde{a_a}|$ if $r_b \le 1$.

Otherwise said, AGM(1, b_0) is contained in the ball $(\bar{\alpha}_n, \tau_{a,b}+2(\tau_a+\tau_3))$, and multiplying this with the ball $(\alpha, 0)$ we obtain a ball containing AGM(α, b). If this can be rounded to a unique MPC number with the belief rounding model, then we have correctly computed mpc_age, otherwise we need to repeat the computations at a higher precision, and the exponent of the multis gives an indication on the messary precision increase.

In practice convergence of the sequences of \hat{u}_n and \hat{b}_n is often such that $\hat{u}_n = \hat{b}_n$ and thus $r_{n,k} = 0$. Otherwise the two generally differ by only 1 ulp, and a very coarse estimate of $r_{n,k}$ as a power of 2

 $r_{a,b} \leq 2^{\max\left(\operatorname{Exp}(\Re(\widetilde{a_n} - \widetilde{b_n})), \operatorname{Exp}(\Im(\widetilde{a_n} - \widetilde{b_n}))\right) + 1 - (\min(\operatorname{Exp}(\Re(\widetilde{a_n})), \operatorname{Exp}(\Im(\widetilde{a_n}))) - 1)}$

is enough, where Exp(0) is considered as $-\infty$.



Andreas Enge

Ball arithmetic in MPC

Complex balls

Centre and radius as relative error

$$\begin{aligned} (c,r) &= \{z \in \mathbb{C} : |z-c| \leq r |c|\} \\ &= \{z = c(1+\vartheta) : |\vartheta| \leq r\} \end{aligned}$$



Complex balls

Centre and radius as relative error

$$\begin{aligned} (c,r) &= \{z \in \mathbb{C} : |z-c| \leq r |c|\} \\ &= \{z = c(1+\vartheta) : |\vartheta| \leq r\} \end{aligned}$$

• Scaling is easy

$$s(c, r) \subseteq (sc, r)$$
 for $s \in \mathbb{R}$

Multiplication is easy

$$(c_1, r_1) + (c_2, r_2) \subseteq (c_1 c_2, r_1 + r_2 + r_1 r_2)$$

• Square root is easy

$$\sqrt{(\boldsymbol{c},\boldsymbol{r})} \subseteq (\sqrt{\boldsymbol{c}},\boldsymbol{r}/2)$$

• Addition is more difficult

$$(c_1, r_1) + (c_2, r_2) \subseteq (c_1 + c_2, (|c_1|r_1 + |c_2|r_2)/|c_1 + c_2|)$$



0 is not representable.

Complex balls in C

```
typedef struct {
   mpc_t c;
   mpcr_t r;
} mpcb_t;
```



Complex balls in C

```
typedef struct {
    mpc_t c;
    mpcr_t r;
} mpcb_t;
```

```
typedef struct {
    int64_t mant;
    int64_t exp;
} mpcr_t;
```



Complex balls in C

```
typedef struct {
    mpc_t c;
    mpcr_t r;
} mpcb_t;
```

- typedef struct {
 int64_t mant;
 int64_t exp;
- } mpcr_t;
 - radius 0: mant = 0
 - radius ∞ : mant < 0
 - 31 bit normalised mantissa: $2^{30} \leqslant \texttt{mant} < 2^{31}$



Surprisingly many and quirky functions...

Results are normalised and rounded up (unless exception).

- Predicates
 - mpcr_inf_p
 - mpcr_zero_p
 - mpcr_lt_half_p
 - mpcr_cmp
- Setters
 - mpcr_set_inf
 - mpcr_set_zero
 - mpcr_set_one
 - mpcr_set
 - mpcr_set_ui64_2si64
 - mpcr_max
- Output
 - mpcr_out_str

Functions on radii

- Arithmetic
 - mpcr_add, mpcr_mul, ...
 - mpcr_sub_rnd Takes MPFR_RNDU, MPFR_RNDD.
 - mpcr_c_abs_rnd
 Used for error of addition

$$(c_1, r_1) + (c_2, r_2) \subseteq (c_1 + c_2, (|c_1|r_1 + |c_2|r_2)/|c_1 + c_2|)$$



Functions on radii

- Arithmetic
 - mpcr_add, mpcr_mul, ...
 - mpcr_sub_rnd Takes MPFR_RNDU, MPFR_RNDD.
 - mpcr_c_abs_rnd
 Used for error of addition

$$(c_1, r_1) + (c_2, r_2) \subseteq (c_1 + c_2, (|c_1|r_1 + |c_2|r_2)/|c_1 + c_2|)$$

mpcr_add_rounding_error (r, p, rnd)
 Accounts for shift of centre by 1/2 ulp or 1 ulp depending on rnd.
 Adds (1 + r)2^{-p} or twice this to r.
 Called once at the end of each operation.



Principles of complex balls

- $(\mathbf{c}_1, \mathbf{r}_1) \circ (\mathbf{c}_2, \mathbf{r}_2) \subseteq (\mathbf{c}_1 \circ \mathbf{c}_2, \mathbf{r})$
- One precision for real and imaginary part
- Initialised without precision: mpcb_init (z), mpcb_clear (z)
- Precisions are tracked automatically.
 - $z_1 \circ z_2$ gets minimum precision of z_1 and z_2 .
 - Precision is not decreased when radius increases.

Setting complex balls

- mpcb_set_inf (z)
- mpcb_set (z, z1)

Challenges: Exact input and inputs with errors.



Setting complex balls

- mpcb_set_inf (z)
- mpcb_set (z, z1)

Challenges: Exact input and inputs with errors.

- mpcb_set_ui_ui (z, re, im, prec)
 Uses maximum of prec and sizeof (ulong).
- mpcb_set_c (z, c, prec, err_re, err_im)
 Assumes c has err_re and err_im half-ulp errors.
 - If prec large and err_re and err_im = 0: exact, r = 0.
 - Otherwise, r encodes err_re, err_im and rounding of c.

Computations with complex balls

- mpcb_neg (z, z1)
- mpcb_add (z, z1, z2)
- mpcb_mul (z, z1, z2)
- mpcb_sqr (z, z1, z2)
- mpcb_pow_ui (z, z1, e)
- mpcb_sqrt (z, z1)
- mpcb_div (z, z1, z2)
- mpcb_div_2ui (z, z1, e)



Rounding complex balls

• mpcb_round (c, z, rnd)

Rounds the centre of z to c (with its own precision).



Rounding complex balls

- mpcb_round (c, z, rnd)
 Rounds the centre of z to c (with its own precision).
- mpcb_can_round (z, prec_re, prec_im, rnd) true

if rounding any complex (mathematical) number in z to

a complex (floating point) number of precision (prec_re, prec_im) in direction rnd

yields the same result and ternary return value.

Rounding complex balls

- mpcb_round (c, z, rnd)
 Rounds the centre of z to c (with its own precision).
- mpcb_can_round (z, prec_re, prec_im, rnd) true

if rounding any complex (mathematical) number in z to

a complex (floating point) number of precision (prec_re, prec_im) in direction rnd

yields the same result and ternary return value.

• Beware of infinite loops for exact results.



"Normal" functions using balls

mpc_agm

Uses a priori error analysis.

Tests compare with mpc_mpcb_agm.

Ínnía

"Normal" functions using balls

 mpc_agm Uses a priori error analysis. Tests compare with mpc_mpcb_agm.

• mpc_eta_fund (rop, z, rnd)

1 10

$$q^{1/24} = \exp(2\pi i z/24)$$

$$q = \left(q^{1/24}\right)^{24}$$

$$\eta = q^{1/24} \left(1 - q - q^2 + q^5 + q^7 - q^{12} - q^{15} + \cdots\right)$$

 $q^{1/24}$ is computed as mpc with an a priori error analysis, then handled with mpcb_set_c.



What is next?

- Question the design choices.
- Handle 0?
 - Exact 0 is easily encoded, but needs case distinctions.
 - Balls around 0 need absolute radius encoding and more case distinctions.
- Use balls internally for series (special values of special functions)?
- Use balls in existing functions instead of error analysis?
- Implement functions on balls?

