

Optimal Floating-Point Realizations of Finite-Precision Digital Controllers

J. Wu[†], S. Chen[‡], J.F. Whidborne[§] and J. Chu[†]

[†] National Key Laboratory of Industrial Control Technology
Zhejiang University, Hangzhou, 310027, P. R. China

[‡] Department of Electronics and Computer Science
University of Southampton, Southampton SO17 1BJ, U.K.

[§] Department of Mechanical Engineering
King's College London, Strand, London WC2R 2LS, U.K.

Presented at 41st IEEE CDC, December 10-13, 2002
Las Vegas, Nevada, USA

Support of U.K. Royal Society is gratefully acknowledged

Floating-Point Representation

- Floating-point processor of bit length $\beta = 1 + \beta_w + \beta_e$ represents $x \in \mathcal{R}$: one bit for sign, β_w bits for mantissa, and β_e bits for exponent of x .
- Given β_e bits, the lower and upper limits of exponents are \underline{e} and \bar{e} , with $\bar{e} - \underline{e} = 2^{\beta_e} - 1$. Denote the set of integers $\underline{e} \leq e \leq \bar{e}$ as $\mathcal{Z}_{[\underline{e}, \bar{e}]}$.
- If the exponent of x , $e = \lfloor \log_2 |x| \rfloor + 1$, is within $\mathcal{Z}_{[\underline{e}, \bar{e}]}$, there is no underflow or overflow. In such a case, x is perturbed to

$$Q(x) = x(1 + \delta), \quad |\delta| < 2^{-(\beta_w + 1)}$$

The perturbation is multiplicative, unlike the additive perturbation resulting from fixed-point arithmetic.

- β_e determines the dynamic range, and β_w the precision of representation.

Motivations and Background

Finite precision controller implementation can seriously influence closed-loop performance.

- Two types of finite word length errors: roundoff errors in arithmetic operations – controller signal errors, and controller coefficient representation errors – controller parameter errors. This work is concerned with the latter, which has critical influence on closed-loop stability.
- Two strategies: direct and indirect. This work adopts an indirect approach.
- Most works deal with fix-point implementation. This work is for floating-point implemented controllers.
- A main contribution of this work is dealing with not only precision but also dynamic range of a numerical representation scheme.

Problem Definition

- Plant: $P(z) \sim (\mathbf{A}_P, \mathbf{B}_P, \mathbf{C}_P)$; $\mathbf{A}_P \in \mathcal{R}^{m \times m}$, $\mathbf{B}_P \in \mathcal{R}^{m \times l}$, $\mathbf{C}_P \in \mathcal{R}^{q \times m}$.
- Controller: $C(z) \sim (\mathbf{A}_C, \mathbf{B}_C, \mathbf{C}_C, \mathbf{D}_C)$; $\mathbf{A}_C \in \mathcal{R}^{n \times n}$, $\mathbf{B}_C \in \mathcal{R}^{n \times q}$, $\mathbf{C}_C \in \mathcal{R}^{l \times n}$, $\mathbf{D}_C \in \mathcal{R}^{l \times q}$.

Denote an initially designed controller realization as \mathbf{X}_0 and a generic realization \mathbf{X} . Let $\bar{\mathbf{A}}(\mathbf{X})$ be the closed-loop transition matrix with \mathbf{X} .

- Controller realization set

$$\mathcal{S}_C \triangleq \{ \mathbf{X} : \mathbf{A}_C = \mathbf{T}^{-1} \mathbf{A}_C^0 \mathbf{T}, \mathbf{B}_C = \mathbf{T}^{-1} \mathbf{B}_C^0, \mathbf{C}_C = \mathbf{C}_C^0 \mathbf{T}, \mathbf{D}_C = \mathbf{D}_C^0 \}$$

where $\mathbf{T} \in \mathcal{R}^{n \times n}$ is an arbitrary non-singular matrix

- All $\mathbf{X} \in \mathcal{S}_C$ are equivalent in infinite precision implementation: an identical set of closed-loop eigenvalues $\lambda_i(\bar{\mathbf{A}}(\mathbf{X}))$, $1 \leq i \leq m + n$, which are all within the unit disk.

Dynamic Range Measure

- An dynamic range (exponent) measure for floating-point realization \mathbf{X} :

$$\gamma(\mathbf{X}) \triangleq \log_2 \left(\frac{4\|\mathbf{X}\|_{\max}}{g(\mathbf{X})} \right)$$

where $\|\mathbf{X}\|_{\max} \triangleq \max_{j,k} |x_{j,k}|$ and $g(\mathbf{X}) \triangleq \min_{j,k} \{ |x_{j,k}| : x_{j,k} \neq 0 \}$.

- \mathbf{X} can be represented in floating-point format of β_e exponent bits without underflow or overflow, if $2^{\beta_e} \geq \gamma(\mathbf{X})$.
- Let β_e^{\min} be the smallest exponent bit length for \mathbf{X} without underflow and overflow. Then, $\beta_e^{\min} = -\lfloor -\log_2(\lfloor \log_2 \|\mathbf{X}\|_{\max} \rfloor - \lfloor \log_2 g(\mathbf{X}) \rfloor + 1) \rfloor$.
- $\gamma(\mathbf{X})$ provides an estimate of β_e^{\min} as:

$$\hat{\beta}_e^{\min} \triangleq -\lfloor -\log_2 \gamma(\mathbf{X}) \rfloor$$

Tractable Precision Measure

- A tractable precision (mantissa) measure is:

$$\mu_1(\mathbf{X}) \triangleq \min_{i \in \{1, \dots, m+n\}} \frac{1 - |\lambda_i(\overline{\mathbf{A}}(\mathbf{X}))|}{\left\| \frac{\partial |\lambda_i|}{\partial \Delta} \right\|_{\Delta=0} \Big|_{\text{sum}}}$$

where $\left\| \frac{\partial |\lambda_i|}{\partial \Delta} \right\|_{\text{sum}} \triangleq \sum_{j,k} \left| \frac{\partial |\lambda_i|}{\partial \delta_{j,k}} \right|$.

- Under some mild conditions, $|\lambda_i(\overline{\mathbf{A}}(\mathbf{X} + \mathbf{X} \circ \Delta))| < 1$ if $\|\Delta\|_{\max} < \mu_1(\mathbf{X})$.
- Let β_w^{\min} be the smallest mantissa bit length that guarantees closed-loop stability for floating-point implemented \mathbf{X} .
- $\mu_1(\mathbf{X})$ provides an estimate of β_w^{\min} as: $\hat{\beta}_{w1}^{\min} \triangleq -\lfloor \log_2 \mu_1(\mathbf{X}) \rfloor - 1$.

Finite Precision Stability Consideration

- Even without underflow or overflow, due to finite β_w , $\mathbf{X} \Rightarrow \mathbf{X} + \mathbf{X} \circ \Delta$, with perturbation matrix Δ satisfying $\|\Delta\|_{\max} < 2^{-(\beta_w+1)}$.
- With Δ , $\lambda_i(\overline{\mathbf{A}}(\mathbf{X})) \Rightarrow \lambda_i(\overline{\mathbf{A}}(\mathbf{X} + \mathbf{X} \circ \Delta))$: Will any of which become outside the unit disk? Or how robust closed-loop stability is to Δ ?
- It is critical to know how large Δ will cause closed-loop instability for realization \mathbf{X} . Or we would like to know the largest open hypercube in perturbation space, within which closed-loop system remains stable.
- The size of this open hypercube is defined by

$$\mu_0(\mathbf{X}) \triangleq \inf \{ \|\Delta\|_{\max} : \overline{\mathbf{A}}(\mathbf{X} + \mathbf{X} \circ \Delta) \text{ is unstable} \}$$

However, we do not know how to calculate $\mu_0(\mathbf{X})$ given \mathbf{X} .

FWL Closed-Loop Stability Measure

- Goodness of \mathbf{X} can be measured by a large value of $\mu_1(\mathbf{X})$ and a small value of $\gamma(\mathbf{X}) \Rightarrow$ FWL closed-loop stability measure:

$$\rho_1(\mathbf{X}) \triangleq \mu_1(\mathbf{X}) / \gamma(\mathbf{X})$$

- Define the minimum total bit length required in floating point implementation: $\beta_w^{\min} = \beta_e^{\min} + \beta_w^{\min} + 1$. $\rho_1(\mathbf{X})$ provides an estimate of β_w^{\min} as:

$$\hat{\beta}_1^{\min} \triangleq -\lfloor \log_2 \rho_1(\mathbf{X}) \rfloor + 1$$

- $\rho_1(\mathbf{X})$ takes into account both the dynamic range and precision considerations.
- Given a controller realization \mathbf{X} , the value of $\rho_1(\mathbf{X})$ can be computed.

Optimal Realization Problem

- An optimal controller realization problem is defined as

$$v \triangleq \max_{\mathbf{X} \in \mathcal{S}_C} \rho_1(\mathbf{X})$$

- With respect to a given initial realization \mathbf{X}_0 , $\mathbf{X} = \mathbf{X}(\mathbf{T})$. By defining

$$f(\mathbf{T}) \triangleq \rho_1(\mathbf{X}(\mathbf{T}))$$

the optimal realization is posed as the optimization problem:

$$v = \max_{\substack{\mathbf{T} \in \mathcal{R}^{n \times n} \\ \det \mathbf{T} \neq 0}} f(\mathbf{T})$$

- With an optimal transformation matrix \mathbf{T}_{opt} , the optimal realization \mathbf{X}_{opt} can readily be computed.

Example One

Example from (Gevers and Li 1993): $m = 4$, $n = 4$ and $l = q = 1$ with an initially design controller \mathbf{X}_0 .

Realization	\mathbf{X}_0	\mathbf{X}_s	\mathbf{X}_{opt}
ρ_1	2.6644e-9	4.7588e-6	9.5931e-6
$\hat{\beta}_1^{\min}$	30	19	18
μ_1	8.5182e-8	8.7907e-5	1.5229e-4
$\hat{\beta}_{w_1}^{\min}$	23	13	12
γ	3.1971e+1	1.8473e+1	1.5875e+1
$\hat{\beta}_e^{\min}$	5	5	4
β^{\min}	26	15	13
β_w^{\min}	20	9	8
β_e^{\min}	5	5	4

Design Experiments

- MATLAB optimization routine *fminsearch.m* is used to solve the optimization problem numerically.

The resulting optimal controller realization is denoted as \mathbf{X}_{opt} .

- Compare with an existing work (Whidborne and Gu 2002, IFAC World Congress), which minimizes a weighted closed-loop eigenvalue sensitivity index.

This is the only existing work we can find that deals with FWL closed-loop stability of floating-point implemented controller.

Note that this is effectively a precision measure only.

The resulting “optimal” controller realization is denoted as \mathbf{X}_s .

Example Two

Example from (Whidborne *et al.* 2001, IEEE Trans. AC, Vol.46): $m = 2$, $n = 3$ and $l = q = 1$ with an initially design controller \mathbf{X}_0 .

Realization	\mathbf{X}_0	\mathbf{X}_s	\mathbf{X}_{opt}
ρ_1	2.6767e-11	3.1047e-6	5.8446e-6
$\hat{\beta}_1^{\min}$	37	20	19
μ_1	2.8122e-10	7.6679e-5	8.2771e-5
$\hat{\beta}_{w_1}^{\min}$	31	13	13
γ	1.0506e+1	2.4697e+1	1.4162e+1
$\hat{\beta}_e^{\min}$	4	5	4
β^{\min}	30	15	12
β_w^{\min}	25	9	7
β_e^{\min}	4	5	4

Conclusions

- A new computationally tractable FWL closed-loop stability measure has been derived for floating-point controller realizations, which takes into account both the exponent and mantissa of floating-point representation.
- This new measure yields a more accurate estimate for the FWL robustness of closed-loop stability for given controller realization.
- Based on this FWL closed-loop stability measure, the optimal controller realization problem has been formulated, which can be solved for using standard numerical optimization algorithms.