

18.335 Problem Set 2

Due Friday, 2 October 2009.

Problem 1: Floating-point

- (a) Trefethen, problem 13.2. (For part c, you can use Matlab, which employs IEEE double precision by default.)
- (b) A generalization of Trefethen, problem 14.2: given a function $g(x)$ that is analytic (i.e., has a Taylor series) for $|x|$ sufficiently small, and $g'(0) \neq 0$, show that $g(O(\epsilon)) = g(0) + g'(0)O(\epsilon)$.

Problem 2: Addition

This problem is about the floating-point error involved in summing n numbers, i.e. in computing the function $f(x) = \sum_{i=1}^n x_i$ for $x \in \mathbb{F}^n$ (\mathbb{F} being the set of floating-point numbers), where the sum is done in the most obvious way, in sequence. In pseudocode:

```
sum = 0
for i = 1 to n
    sum = sum + x_i
f(x) = sum
```

For analysis, it is a bit more convenient to define the process inductively:

$$\begin{aligned} s_0 &= 0 \\ s_k &= s_{k-1} + x_k \text{ for } 0 < k \leq n, \end{aligned}$$

with $f(x) = s_n$. When we implement this in floating-point, we get the function $\tilde{f}(x) = \tilde{s}_n$, where $\tilde{s}_k = \tilde{s}_{k-1} \oplus x_k$, with \oplus denoting (correctly rounded) floating-point addition.

- (a) Show that $\tilde{f}(x) = (x_1 + x_2) \prod_{k=2}^n (1 + \epsilon_k) + \sum_{i=3}^n x_i \prod_{k=i}^n (1 + \epsilon_k)$, where the numbers ϵ_k satisfy $|\epsilon_k| \leq \epsilon_{\text{machine}}$.
- (b) Show that $\prod_{k=i}^n (1 + \epsilon_k) = 1 + \delta_i$ where $|\delta_i| \leq (n - i + 1)\epsilon_{\text{machine}} + O(\epsilon_{\text{machine}}^2)$.
- (c) Show that the error can be bounded as: $|\tilde{f}(x) - f(x)| \leq n\epsilon_{\text{machine}} \sum_{i=1}^n |x_i|$.
- (d) Suppose that the ϵ_k values are uniformly randomly distributed in $[-\epsilon_{\text{machine}}, +\epsilon_{\text{machine}}]$. Show that

the *mean* error can be bounded by $|\tilde{f}(x) - f(x)| = O(\sqrt{n}\epsilon_{\text{machine}} \sum_{i=1}^n |x_i|)$. (Hint: google “random walk”...you can just quote standard statistical results for random walks, no need to copy the proofs.)

- (e) Compare your error bounds above to numerical experiments in Matlab. Here, we will use an old trick to compute the floating-point errors: compare the results computed in one precision to the “exact” results computed in a higher precision. In particular, we will use the Matlab `single()` function to accumulate the sum in single precision, rather than Matlab’s default double precision. Plot the error $|\tilde{f}(x) - f(x)| / \sum_i |x_i|$ as a function of n on a log-log scale (Matlab’s `loglog` command), and explain your observation in terms of your results above.¹

This is implemented in the example file `loopsum.m`, posted on the course page, which computes the sum $f(x) = \text{loopsum}(x)$ via the above algorithm in single precision. For your numerical experiment, compute the sum of n random inputs $x \in [0, 1]^n$ via Matlab’s `rand(1, n)` function. You can then compute $|\tilde{f}(x) - f(x)| / \sum_i |x_i|$ for a given n via

```
x = rand(1, n);
err = abs(loopsum(x) - sum(x)) /
sum(abs(x));
```

Problem 3: Addition, another way

Here you will analyze $f(x) = \sum_{i=1}^n x_i$ as in problem 2, but this time you will compute $\tilde{f}(x)$ in a different way. In particular, compute $f(x)$ by a recursive divide-and-conquer approach, recursively dividing the set of values to be summed in two halves and then summing the halves:

$$\tilde{f}(x) = \begin{cases} 0 & \text{if } n = 0 \\ x_1 & \text{if } n = 1, \\ \tilde{f}(x_{1:\lfloor n/2 \rfloor}) \oplus \tilde{f}(x_{\lfloor n/2 \rfloor + 1:n}) & \text{if } n > 1 \end{cases}$$

¹Use enough n values to get a clear result. e.g. the command `n = round(logspace(2, 6, 100))` gives 100 logarithmically spaced n values from 10^2 to 10^6 .

where $\lfloor y \rfloor$ denotes the greatest integer $\leq y$ (i.e. y rounded down). In exact arithmetic, this computes $f(x)$ exactly, but in floating-point arithmetic this will have very different error characteristics than the simple loop-based summation in problem 2.

- (a) For simplicity, assume n is a power of 2 (so that the set of numbers to add divides evenly in two at each stage of the recursion). With an analysis similar to that of problem 2, prove that $|\tilde{f}(x) - f(x)| \leq \epsilon_{\text{machine}} \log_2(n) \sum_{i=1}^n |x_i| + O(\epsilon_{\text{machine}}^2)$. That is, show that the worst-case error bound grows *logarithmically* rather than *linearly* with n !
- (b) If the floating-point rounding errors are randomly distributed as in problem 2, estimate the average-case error bound.
- (c) Pete R. Stunt, a Microsoft employee, complains, “While doing this kind of recursion may have nice error characteristics in theory, it is ridiculous in the real world because it will be insanely slow—I’m proud of my efficient software and can’t afford to have a function-call overhead for every number I want to add!” Explain to Pete how to implement a slight variation of this algorithm with the same logarithmic error bounds (possibly with a worse constant factor) but roughly the same performance as a simple loop (hint: look at how I implemented recursive matrix multiplication in my cache-oblivious handout from lecture 3).²
- (d) On the course web page, I’ve posted a file `div2sum.m` that computes $\tilde{f}(x) = \text{div2sum}(x)$ by the above algorithm. Modify it to not be horrendously slow via your suggestion in (c), and then plot its errors for random inputs as a function of n as in problem 2. Are your results consistent with your error estimates above?
- (e) Suppose we now multiply two $m \times m$ random matrices A and B ($\in [0, 1]^{m \times m}$, uniformly distributed) to form $C = AB$. If you

look at any given entry c_{ij} of C , how quickly do you expect the errors to grow with m if you compute AB via the simple 3-loop row-column algorithm? What if you use the optimal cache-oblivious algorithm from class?

Problem 4: Stability

- (a) Trefethen, exercise 15.1.
- (b) Trefethen, exercise 16.1.

²In fact, there is a common real-world algorithm that does summation in precisely this recursive way: the Cooley-Tukey fast Fourier transform.