APCK 23 Prob 1

1. prove $\left|\dfrac{x - rd(x)}{x}\right| \le 2^{-p} \to x = \pm\left(\sum\limits_{\ell=1}^{\infty} b_{-\ell}\, 2^{-\ell}\right) \cdot 2^{e}$

Now $\quad rd(x) = \begin{cases} \dfrac{\left(\pm\sum\limits_{\ell=1}^{p} b_{-\ell} 2^{-\ell}\right) \cdot 2^{e} \,\therefore\, b_{-(p+1)} = 0}{} \\[4mm] \left(\pm\left(\sum\limits_{\ell=1}^{p} b_{-\ell} 2^{-\ell}\right) + \dfrac{2^{-p}}{\uparrow}\right) 2^{e} \,\therefore\, b_{-(p+1)} = 0 \\ \quad\text{from rounding} \end{cases}$

$\to$ First consider $b_{-(p+1)} = 0$ (discarded bit is zero)

Then $x - rd(x) = \pm\left(\sum\limits_{\ell=1}^{\infty} b_{-\ell} 2^{-\ell} - \sum\limits_{\ell=1}^{p} b_{-\ell} 2^{-\ell}\right) \cdot 2^{e}$

$= \pm\left(b_{-(p+1)} \cdot 2^{-(p+1)} + \sum\limits_{\ell=p+2}^{\infty} b_{-\ell} 2^{-\ell}\right) \cdot 2^{e}$

① $\quad 0 = b_{-(p+1)}$

so $\quad x - rd(x) = \pm\sum\limits_{\ell=p+2}^{\infty} b_{-\ell} 2^{-\ell} \; 2^{e} \quad$ for $b_{-(p+1)} = 0$

Similarly, if $b_{-(p+1)} = 1$ then

$2^{-p} - 2^{-p} = -2^{-(p+1)}$

$x - rd(x) = \pm\left(\sum\limits_{\ell=1}^{\infty} b_{-\ell} 2^{-\ell} - \sum\limits_{\ell=1}^{p} b_{-\ell} 2^{-\ell} - 2^{-p}\right) 2^{e}$

$= \pm\left(b_{-(p+1)} \cdot 2^{-(p+1)} - 2^{-p} + \sum\limits_{\ell=p+2}^{\infty} b_{-\ell} 2^{-\ell}\right) 2^{e} \quad -2^{-(p+1)}$

$1 = b_{-(p+1)}$

$x - rd(x) = \pm\left(\left(\sum\limits_{\ell=p+2}^{\infty} b_{-\ell} 2^{-\ell}\right) - 2^{-(p+1)}\right) 2^{e} \quad$ for $b_{-(p+1)} = 1$

$\to$ returning to $b_{-(p+1)} = 0$, $\left|\dfrac{x - rd(x)}{x}\right|$ is maximized if all $b_{-\ell}$ from $p+2 \to \infty = 1$

so $\left|\dfrac{x - rd(x)}{x}\right| < \dfrac{\left|\sum\limits_{\ell=p+2}^{\infty} 2^{-\ell}\right| 2^{e}}{2^{e-1}}$

Min floating point
$x$ is $2^{e} - 1$ (see class notes)

$$= 1 + 2^{-1} + 2^{-2} \ldots = \underline{2}$$

$$\left| \frac{x - rd(x)}{x} \right| \leq \left| 2^{-p-2} \left( \sum_{l=0}^{\infty} 2^{-l} \right) \right| \cdot 2$$

$$\leq 2^{-p-1} \cdot 2$$

$$\therefore \left| \frac{x - rd(x)}{x} \right| \leq 2^{-p} \quad \text{for } b_{-(p+1)} = 0$$

now for $b_{-(p+1)} = 1$, $|x - rd(x)|$ is maximized if $b_{-(p+1)} = 1$ and $b_{-(1+2)} \to b_{-\infty} = 0$

Since then are just barely of the threshold to round up and thus furthest from from $rd(x)$

$\therefore$ for $b_{-(p+1)} = 1$ $\left| \frac{x - rd(x)}{x} \right| \leq \left| \frac{\sum_{l=p+2}^{\infty} 0 \cdot 2^{-l} - 2^{-(p+1)}}{2^{e-1}} \right| 2^e$

$$\leq \frac{2^{-(p+1)} 2^e}{2^{e-1}}$$

$$\therefore \left| \frac{x - rd(x)}{x} \right| \leq 2^{-p} \quad \text{for } b_{-(p+1)} = 1$$

$\Rightarrow$ have shown $\left| \frac{x - rd(x)}{x} \right| \leq 2^{-p}$ for all cases as required