**Problem 2.83 Solution:**

This problem helps students appreciate the property of IEEE floating point that the relative magnitude of two numbers can be determined by viewing the combination of exponent and fraction as an unsigned integer. Only the signs and the handling of $\pm 0$ requires special consideration.

*code/data/floatcomp-ans.c*

```
1 int float_le(float x, float y) {
2     unsigned ux = f2u(x);
3     unsigned uy = f2u(y);
4     unsigned sx = ux >> 31;
5     unsigned sy = uy >> 31;
6
7     return
8         (ux<<1 == 0 && uy<<1 == 0) ||  /* Both are zero */
9         (sx && !sy) ||                 /* x < 0, y >=  0 */
10        (!sx && !sy && ux <= uy) ||    /* x >= 0, y >= 0 */
11        (sx && sy && ux >= uy);        /* x <  0, y <  0 */
12 }
```

**Problem 2.88 Solution:**

This problem requires students to think of the relationship between `int`, `float`, and `double`.

A. `(float) x == (float) dx`. Yes. Converting to `float` could cause rounding, but both `x` and `dx` will be rounded in the same way.

B. `dx - dy == (double) (x-y)`. No. Let $x = 0$ and $y = TMin_{32}$.

C. `(dx + dy) + dz == dx + (dy + dz)`. Yes. Since each value ranges between $TMin_{32}$ and $TMax_{32}$, their sum can be represented exactly.

D. `(dx * dy) * dz == dx * (dy * dz)`. No. Let $dx = TMax_{32}$, $dy = TMax_{32} - 1$, $dz = TMax_{32} - 2$. (Not detected with Linux/GCC)

E. `dx / dx == dz / dz`. No. Let `x = 0, z = 1`.

**Problem 2.89 Solution:**

This problem helps students understand the relation between the different categories of numbers. Getting all of the cutoff thresholds correct is fairly tricky. Our solution file contains testing code.

—————————————————————————————————————— *code/data/fpwr2-ans.c*

```
1 /* Compute 2**x */
2 float fpwr2(int x) {
3
4     unsigned exp, frac;
5     unsigned u;
6
7     if (x < -149) {
8         /* Too small.  Return 0.0 */
```

```
 9          exp  = 0;
10          frac = 0;
11      } else if (x < -126) {
12          /* Denormalized result */
13          exp  = 0;
14          frac = 1 << (x + 149);
15      } else if (x < 128) {
16          /* Normalized result. */
17          exp  = x + 127;
18          frac = 0;
19      } else {
20          /* Too big.  Return +oo */
21          exp  = 255;
22          frac = 0;
23      }
24      u = exp << 23 | frac;
25      return u2f(u);
26 }
```

**Problem 3.56 Solution:**

One way to analyze assembly code is to try to reverse the compilation process and produce C code that would look "natural" to a C programmer. For example, we wouldn't want any goto statements, since these are seldom used in C. Most likely, we wouldn't use a do-while statement either. This exercise forces students to reverse the compilation into a particular framework. It requires thinking about the translation of for loops.

A. We can see that result must be in register %edi, since this value gets copied to %eax at the end of the function as the return value (line 13). We can see that %esi and %ebx get loaded with the values of x and n (lines 1 and 2), leaving %edx as the one holding variable mask (line 4.)

B. Register %edi (result) is initialized to $-1$ and %edx (mask) to 1.

C. The condition for continuing the loop (line 12) is that mask is nonzero.

D. The shift instruction on line 10 updates mask to be mask << n.

E. Lines 6–8 update result to be result ^ (x&mask).

F. Here is the original code:

*code/asm/for.c*

```
1 int loop(int x, int n)
2 {
3     int result = -1;
4     int mask;
5     for (mask = 0x1; mask != 0; mask = mask << n) {
6         result ^= (x & mask);
7     }
8     return result;
9 }
```

*code/asm/for.c*