# Error in Computational Tools
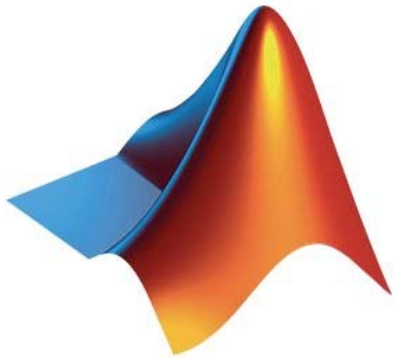
- Mathematically, we know that
$$10^{-16} = 1 + 10^{-16} - 1$$

# Error in Computational Tools

- Mathematically, we know that

$$10^{-16} = 1 + 10^{-16} - 1$$

```
>> 1+10^(-14)
ans =
    1.0000
>> 1+10^(-15)
ans =
    1.0000
>> 1+10^(-16)
ans =
    1
```

# Error in Computational Tools

- Mathematically, we know that
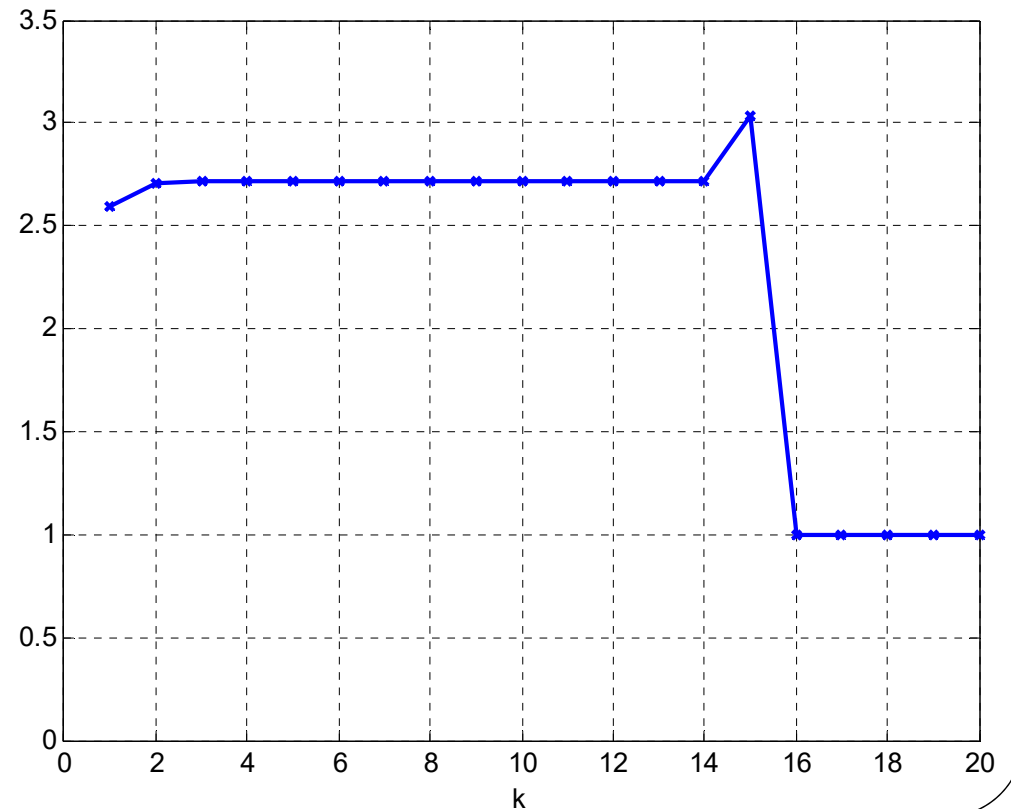$$10^{-16} = 1 + 10^{-16} - 1$$

```
>> log10(10^(-16))
ans =
    -16
>> log10(1+10^(-16)-1)
ans =
   -Inf
```

# Error in Computational Tools

- Now, let's work on expression of the form

$$\left(1 + \frac{1}{10^k}\right)^{10^k}$$

```
>> k = 14; (1+10^(-k))^(10^k)
ans =
    2.716110034087023
>> k = 15; (1+10^(-k))^(10^k)
ans =
    3.035035206549262
>> k = 16; (1+10^(-k))^(10^k)
ans =
    1
```



54

# Error in Computational Tools

- Now, let's work on expression of the form

$$\left(1 + \frac{1}{10^k}\right)^{10^k}$$

Theoretically, when $k$ is large, we know that this expression should converge to $e^1 = e \approx 2.7183$.

```
>> k = 14; (1+10^(-k))^(10^k)
ans =
    2.716110034087023
>> k = 15; (1+10^(-k))^(10^k)
ans =
    3.035035206549262
>> k = 16; (1+10^(-k))^(10^k)
ans =
    1
```
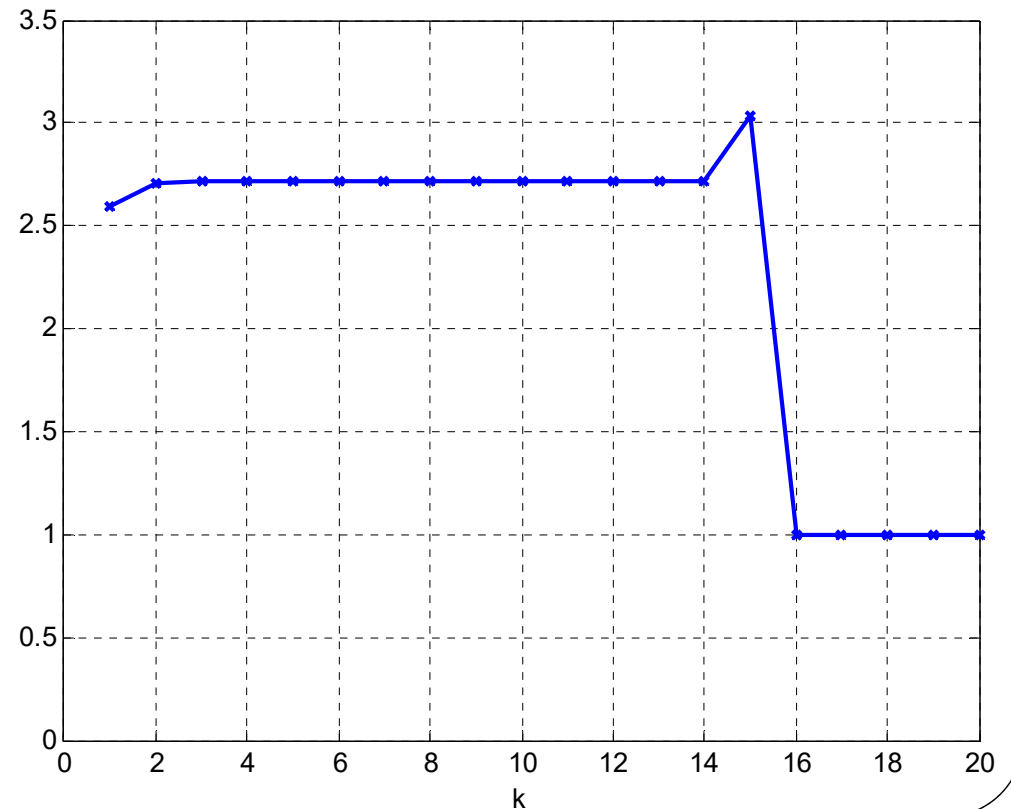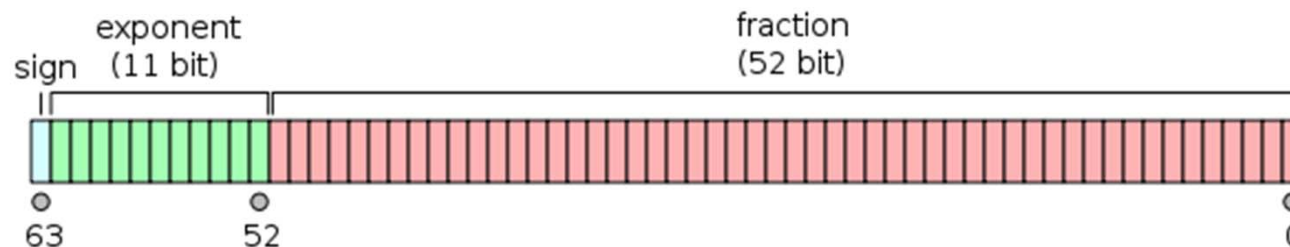
# Accuracy of Floating-Point Data

- Computers only represent numbers to a finite precision.

- Computations sometimes yield mathematically nonintuitive results.

- Almost all operations in MATLAB are performed in double-precision arithmetic conforming to the IEEE standard 754.



- The real value assumed by a given 64-bit double-precision datum with a given biased exponent $e$ and a 52-bit fraction is

$$(-1)^{\text{sign}} (1.b_{51}b_{50}...b_0)_2 \times 2^{e-1023}$$

# Double-Precision Accuracy

- Because there are only a finite number of double-precision numbers, you cannot represent all numbers in double-precision storage.

- On any computer, there is a small gap between each double-precision number and the next larger double-precision number.

- You can determine the size of this gap, which limits the precision of your results, using the **eps** function.
  - This **machine epsilon** gives an upper bound on the relative error due to rounding in floating point arithmetic.

- For example, to find the distance between 1 and the next larger double-precision number, enter

```
>> format long
>> eps(1)
ans =
      2.220446049250313e-16
```

# Probability Calculation for Binomial RV

- For binomial RV, probability are of the form
$$\binom{n}{x} p^x (1-p)^{n-x}.$$

- For example, when $x = 0$, we have
$$\binom{n}{0} p^0 (1-p)^{n-0} = (1-p)^n.$$

- When is $p$ small, the number $1 - p$ may not be accurately represented.

# Probability Calculation for Binomial RV

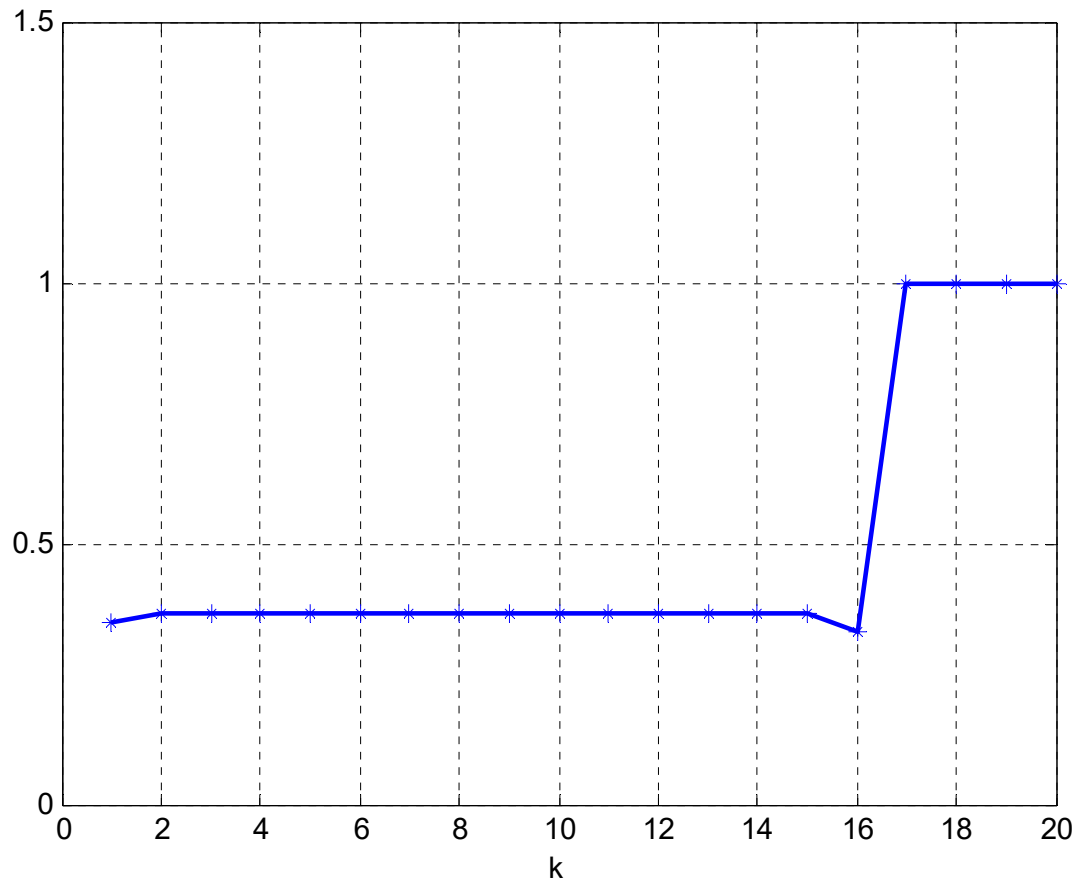- For binomial RV, probability are of the form

$$\binom{n}{x} p^x (1-p)^{n-x}.$$

- For example, when $x = 0$, we have

$$\binom{n}{0} p^0 (1-p)^{n-0} = (1-p)^n.$$

- When is $p$ small, the number $1-p$ may not be accurately represented.
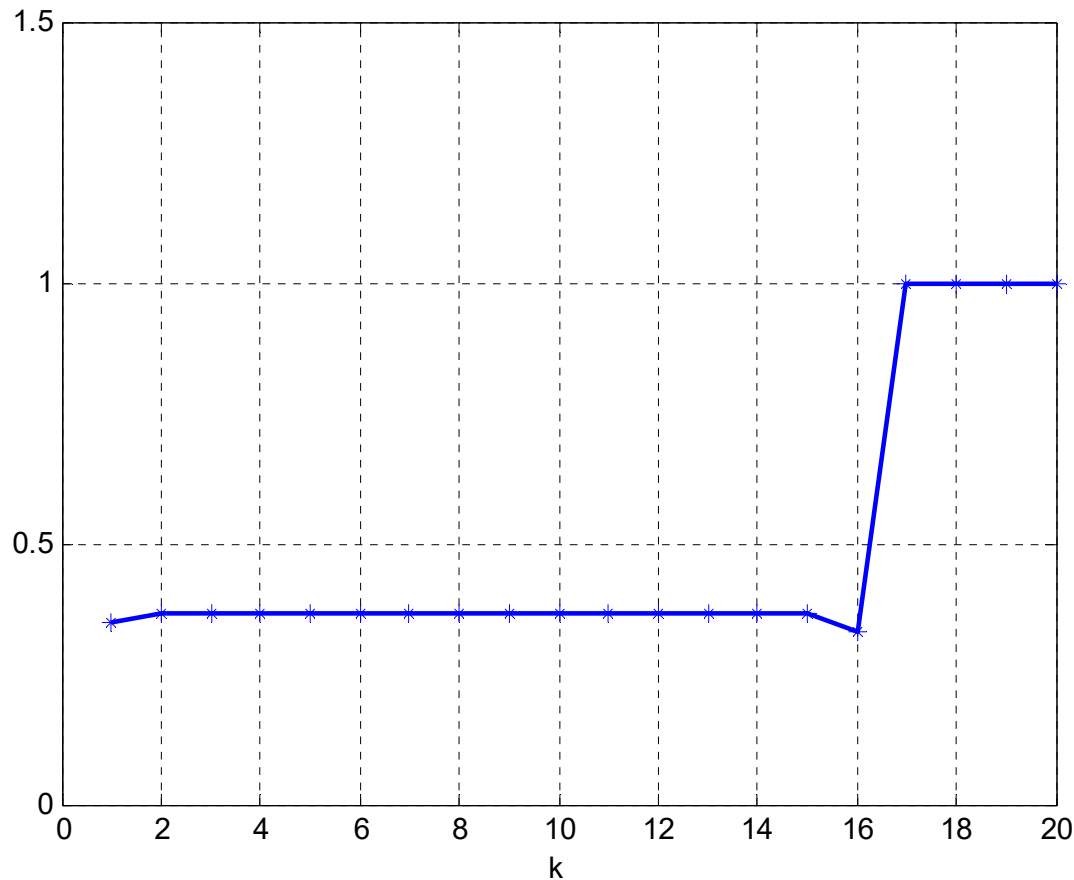
- Consider $p = \frac{1}{10^k}$ and $n = 10^k$.

# Probability Calculation for Binomial RV

- $(1-p)^n$ when $p = \dfrac{1}{10^k}$ and $n = 10^k$.

# Probability Calculation for Binomial RV

- $(1-p)^n$ when $p = \dfrac{1}{10^k}$ and $n = 10^k$.

Theoretically, when $k$ is large, we know that this expression should converge to $e^{-1} \approx 0.3679$.