

Distributions

Outline

- Discrete random variables and probability distributions
- The Binomial probability distribution
- The Poisson probability distribution
- The normal probability distribution

Discrete Random Variables and Probability Distributions

- **Random variables**

- value of which is the result of a random event, e.g. the number of laptops sold on a *randomly selected* day, the age of a student *randomly selected* on campus.

- **Discrete variables**

- One type of ratio variable, can take only a limited number of possible values within a given range, e.g. number of laptops.
- As opposed to **continuous variables** which can take unlimited number of possible values in the same range, e.g. the distance between two locations.

Discrete Random Variables and Probability Distributions

- **Probability distribution** for a discrete random variable:

$$P(A) = \lim_{n \rightarrow \infty} \frac{\text{Frequency}}{n}$$

- Relative frequency distribution constructed for the entire population of measurements
- Represents possible values of x and the probability associated with each value of x .
- Possible values of x are mutually exclusive events.
- $\sum_{i=1}^n p(x_i) = 1$
- Example –The probability distribution of number of cups of coffee taken a day

x	0	1	2	3	4	5
$p(x)$.28	.37	.17	.12	.05	.01

Discrete Random Variables and Probability Distributions

- **Probability distribution** for a discrete random variable:

- Mean (aka the expected value): $\mu = \sum_{i=1}^n x_i p(x_i)$

- Standard variation: $\sigma = \sqrt{\sum_{i=1}^n (x_i - \mu)^2 p(x_i)}$

- Example:

x	0	1	2	3	4	5
$p(x)$.28	.37	.17	.12	.05	.01

- What is the probability of a random coffee drinker taking no coffee a day?
- What is the probability of a random coffee drinker taking more than 2 coffees a day?
- What is the number of coffees a random coffee drinker expected to drink a day?
- What is the probability of the number of coffees a day fall in between $\mu \pm 2\sigma$?

The Binomial Probability Distribution

- **Binomial random variable** – has only two possible values
- **Binomial experiment**
 - Contains n *identical* trials.
 - Each trial results in *one of two outcomes*, e.g. Success or Failure.
 - The probability p of an outcome, e.g. Success, *remains the same* for all trials.
 - Trials are *independent*.
 - We are interested in x , the *number of Successes* observed in n trials.
- Example of a binomial experiment:
 - Tossing a coin 1000 times and observing the number of heads
 - Test of a new drug and counting the number of successful cases
 - Purchase lottery tickets many (many) times and count the number of wins.

The Binomial Probability Distribution

- The **probability** of $x = k$ successes (p is the probability of success) in n trials is

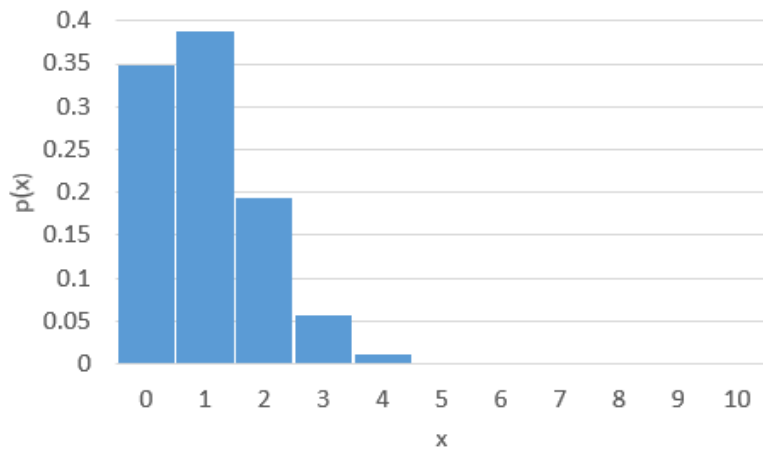
$$P(x = k) = C_k^n p^k (1 - p)^{n-k} = \frac{n!}{k!(n-k)!} p^k (1 - p)^{n-k} \quad \text{for } k = 0, 1, \dots, n$$

where $n! = n(n - 1)(n - 2) \dots (2)(1)$ and $0! = 1$

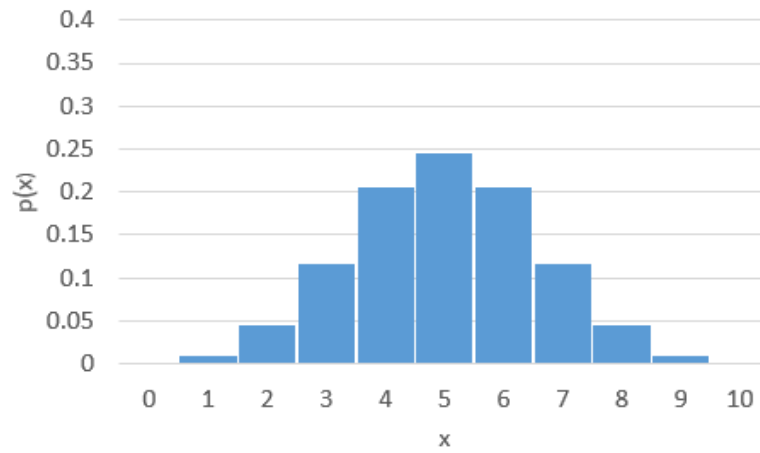
- Distribution of random variable x (the number of successes in n trials) has
 - **Mean** $\mu = np$
 - **Standard deviation** $\sigma = \sqrt{np(1 - p)}$

The Binomial Probability Distribution

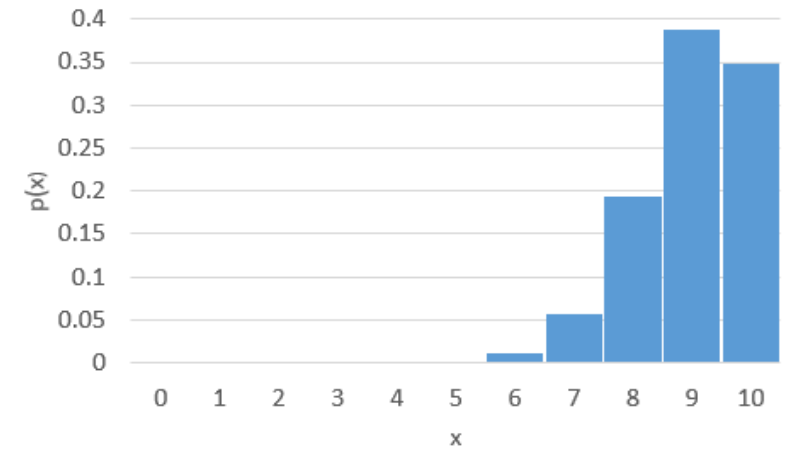
- Examples of binomial probability distribution



$n = 10, p = .1$
mean $\mu = np = 1$
std $\sigma = .95$



$n = 10, p = .5$
mean $\mu = np = 5$
std $\sigma = 1.58$



$n = 10, p = .9$
mean $\mu = np = 9$
std $\sigma = .95$

- Notice the shape of the distribution and expected value (mean) in each case.

The Binomial Probability Distribution - Examples

Example 1: What is the probability of tossing a coin 10 times and seeing 6 heads?

$$n = 10, p = 0.5, k = 6 \Rightarrow P(6) = \frac{10!}{6!(4)!} 0.5^6 (0.5)^4 = 0.205$$

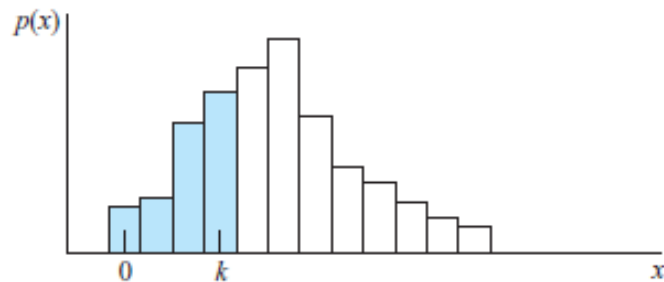
Example 2: 80% of pet insurance owners are women. If we randomly pick 9 pet insurance owners, what is the probability of have 5 females?

The Binomial Probability Distribution - Examples

Example 3: 60% of sport car buyers are men. If we randomly pick 25 of sport car buyers, what is the probability of have 10 men?

Solution hints

Consult the cumulative binomial probability table for $n = 25$, $p = 0.6$



$n = 25$

k	p												k		
	.01	.05	.10	.20	.30	.40	.50	.60	.70	.80	.90	.95		.99	
0	.778	.277	.072	.004	.000	.000	.000	.000	.000	.000	.000	.000	.000	.000	0
1	.974	.642	.271	.027	.002	.000	.000	.000	.000	.000	.000	.000	.000	.000	1
2	.998	.873	.537	.098	.009	.000	.000	.000	.000	.000	.000	.000	.000	.000	2
3	1.000	.966	.764	.234	.033	.002	.000	.000	.000	.000	.000	.000	.000	.000	3
4	1.000	.993	.902	.421	.090	.009	.000	.000	.000	.000	.000	.000	.000	.000	4
5	1.000	.999	.967	.617	.193	.029	.002	.000	.000	.000	.000	.000	.000	.000	5
6	1.000	1.000	.991	.780	.341	.074	.007	.000	.000	.000	.000	.000	.000	.000	6
7	1.000	1.000	.998	.891	.512	.154	.022	.001	.000	.000	.000	.000	.000	.000	7
8	1.000	1.000	1.000	.953	.677	.274	.054	.004	.000	.000	.000	.000	.000	.000	8
9	1.000	1.000	1.000	.983	.811	.425	.115	.013	.000	.000	.000	.000	.000	.000	9
10	1.000	1.000	1.000	.994	.902	.586	.212	.034	.002	.000	.000	.000	.000	.000	10
11	1.000	1.000	1.000	.998	.956	.732	.345	.078	.006	.000	.000	.000	.000	.000	11
12	1.000	1.000	1.000	1.000	.983	.846	.500	.154	.017	.000	.000	.000	.000	.000	12
13	1.000	1.000	1.000	1.000	.994	.922	.655	.268	.044	.002	.000	.000	.000	.000	13
14	1.000	1.000	1.000	1.000	.998	.966	.788	.414	.098	.006	.000	.000	.000	.000	14
15	1.000	1.000	1.000	1.000	1.000	.987	.885	.575	.189	.017	.000	.000	.000	.000	15
16	1.000	1.000	1.000	1.000	1.000	.996	.946	.726	.323	.047	.000	.000	.000	.000	16
17	1.000	1.000	1.000	1.000	1.000	.999	.978	.846	.488	.109	.002	.000	.000	.000	17
18	1.000	1.000	1.000	1.000	1.000	1.000	.993	.926	.659	.220	.009	.000	.000	.000	18
19	1.000	1.000	1.000	1.000	1.000	1.000	.998	.971	.807	.383	.033	.001	.000	.000	19
20	1.000	1.000	1.000	1.000	1.000	1.000	1.000	.991	.910	.579	.098	.007	.000	.000	20
21	1.000	1.000	1.000	1.000	1.000	1.000	1.000	.998	.967	.766	.236	.034	.000	.000	21
22	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	.991	.902	.463	.127	.002	.000	22
23	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	.998	.973	.729	.358	.026	.000	23
24	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	.996	.928	.723	.222	.000	24
25	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	25

The Poisson Probability Distribution

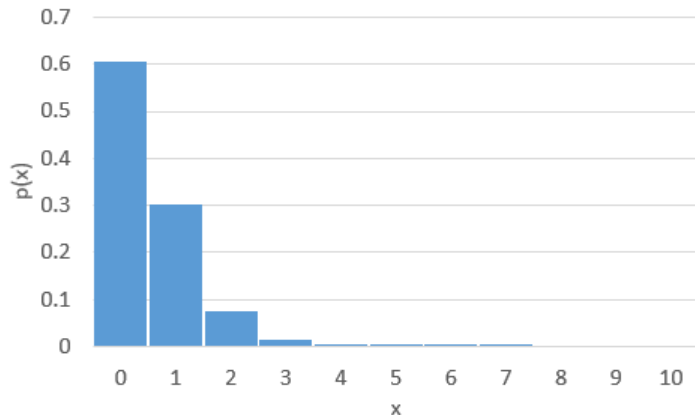
- Poisson's probability distribution
 - a good model for representing *the number of events* over a unit of *time* or *space*.
 - The events must occur *randomly* and *independently* (i.e. not at the same time)
 - The average time (distance) between events is *known* but their exact timing (location) is *unknown*.
- Examples:
 - The number of machine breakdowns during a given day
 - The number of traffic accidents at a given intersection during a given time period
 - The number passengers arriving at a bus stop in a given time window.

The Poisson Probability Distribution

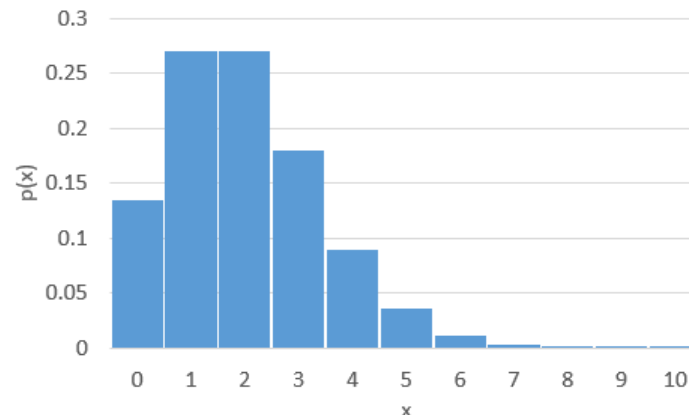
- The probability of k occurrences for the average number of occurrences μ

$$P(x = k) = \frac{\mu^k e^{-\mu}}{k!} \quad \text{for } k = 1, 2, 3, \dots$$

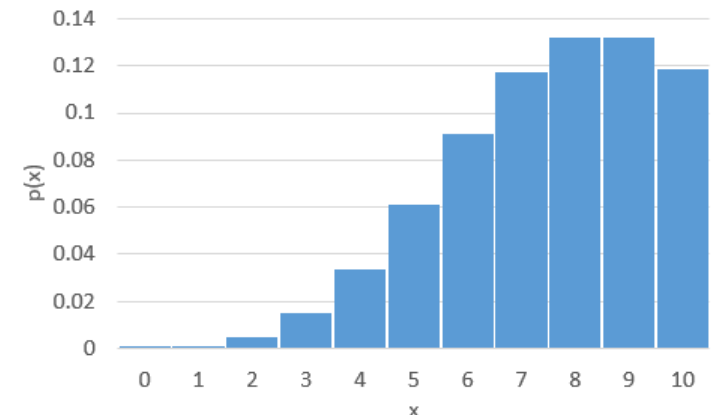
- Mean μ
- Standard deviation $\sigma = \sqrt{\mu}$



$\mu=0.5, \sigma=0.707$



$\mu=2, \sigma=1.414$



$\mu=9, \sigma=3$

- Notice the mode = expected value (mean) and that $p(x = \mu) = p(x = \mu - 1)$

The Poisson Probability Distribution

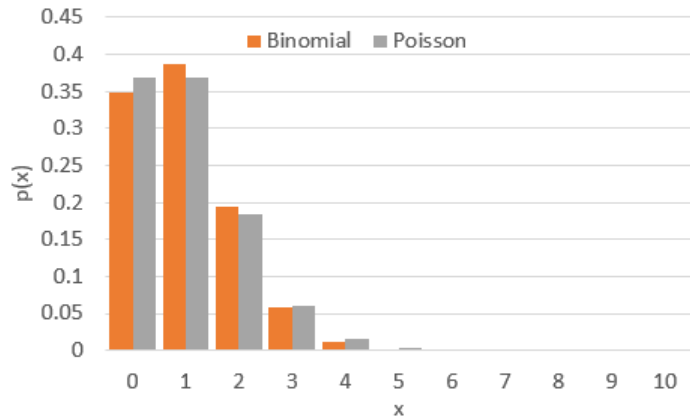
Example 1. The average number of traffic accidents on certain section of highway is 2 per week. What is (a) the probability of having at most 3 accidents and (b) the probability of having exactly 3 accidents a week?

k	μ											
	2.0	2.5	3.0	3.5	4.0	4.5	5.0	5.5	6.0	6.5	7.0	
0	.135	.082	.055	.033	.018	.011	.007	.004	.003	.002	.001	
1	.406	.287	.199	.136	.092	.061	.040	.027	.017	.011	.007	
2	.677	.544	.423	.321	.238	.174	.125	.088	.062	.043	.030	
3	.857	.758	.647	.537	.433	.342	.265	.202	.151	.112	.082	
4	.947	.891	.815	.725	.629	.532	.440	.358	.285	.224	.173	
5	.983	.958	.916	.858	.785	.703	.616	.529	.446	.369	.301	
6	.995	.986	.966	.935	.889	.831	.762	.686	.606	.563	.450	
7	.999	.996	.988	.973	.949	.913	.867	.809	.744	.673	.599	
8	1.000	.999	.996	.990	.979	.960	.932	.894	.847	.792	.729	
9		1.000	.999	.997	.992	.983	.968	.946	.916	.877	.830	
10			1.000	.999	.997	.993	.986	.975	.957	.933	.901	
11				1.000	.999	.998	.995	.989	.980	.966	.947	
12					1.000	.999	.998	.996	.991	.984	.973	
13						1.000	.999	.998	.996	.993	.987	
14							1.000	.999	.999	.997	.994	
15								1.000	.999	.999	.998	
16									1.000	1.000	.999	
17												1.000

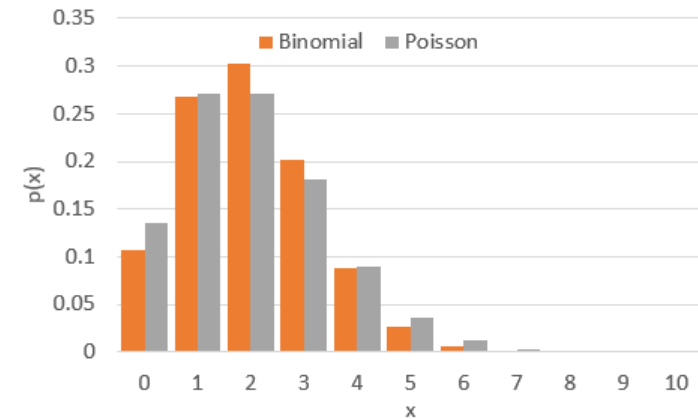
Solution hints. Assuming the number of weekly accidents follows a Poisson distribution. Consult the cumulative Poisson probability table with $\mu=2$.

The Poisson Probability Distribution

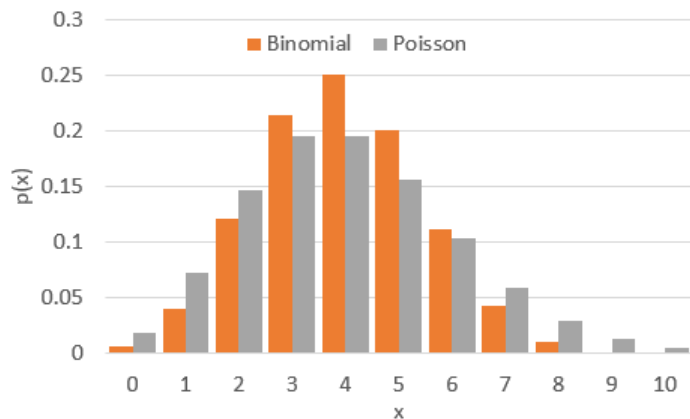
- Poisson distribution can be a good approximation of a binomial distribution that has small $\mu=np$ and preferably *large* n .



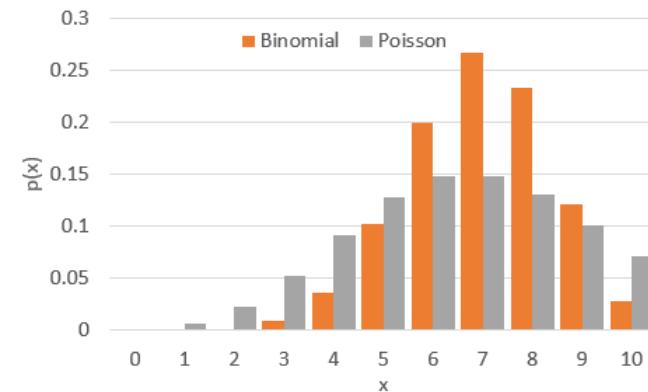
$n = 10, p = 0.1, \mu=1, \text{MAD} = 0.058$



$n = 10, p = 0.2, \mu=2, \text{MAD}=0.104$

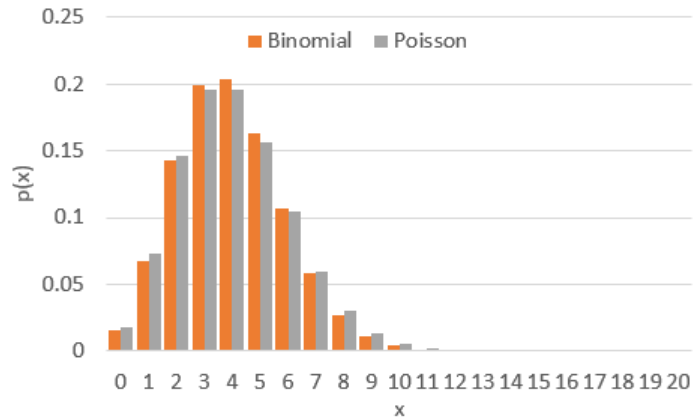


$n=10, p=0.4, \mu=4, \text{MAD}=0.251$

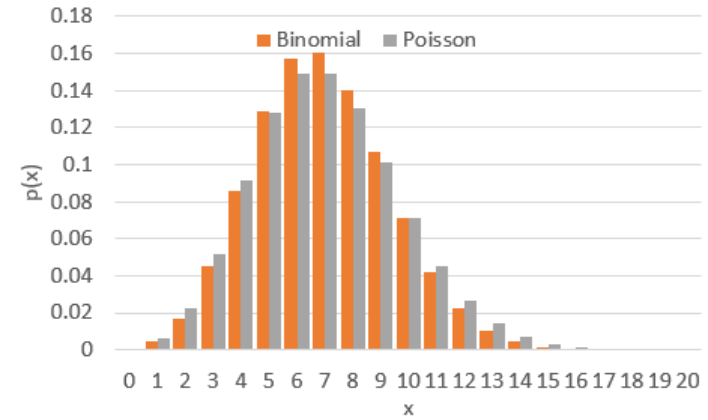


$n = 10, p = 0.7, \mu=7, \text{MAD}=0.485$

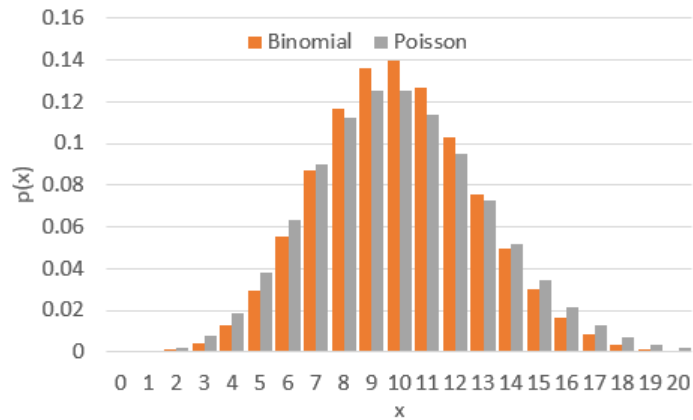
The Poisson Probability Distribution



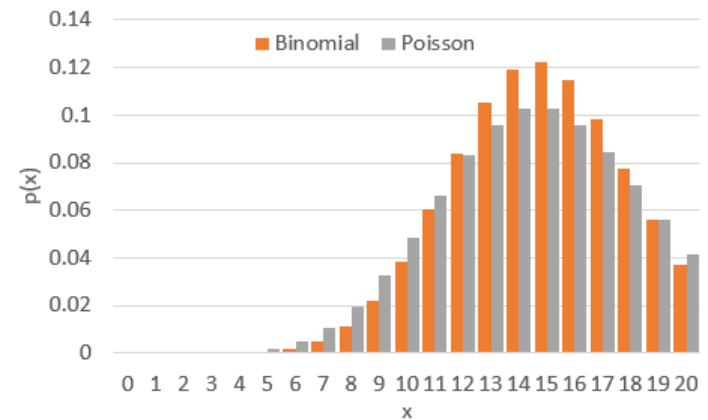
$n = 50, p = 0.08, \mu=4, \text{MAD}=0.042$



$n = 50, p = 0.14, \mu=7, \text{MAD}=0.073$



$n = 50, p = 0.2, \mu=10, \text{MAD}=0.108$



$n = 50, p = 0.3, \mu=15, \text{MAD}=0.136$

The Poisson Probability Distribution

Example 2. Assume the probability of a defective engine is $p=0.001$. Given a batch of 1000 engines, what is the probability of having 4 defective engines?

Solutions. This is a binomial experiment with $n = 1000$, $p = 0.001$, probability of having 4

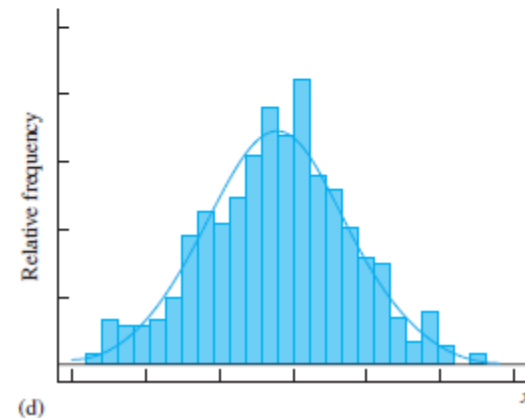
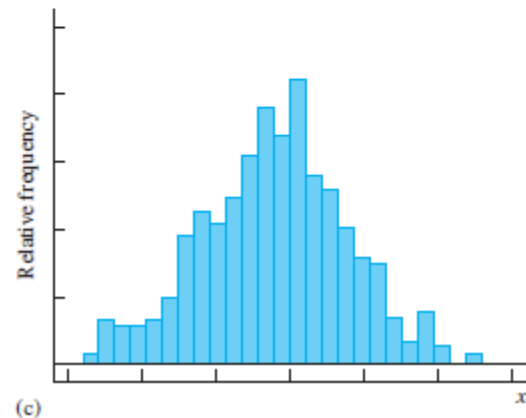
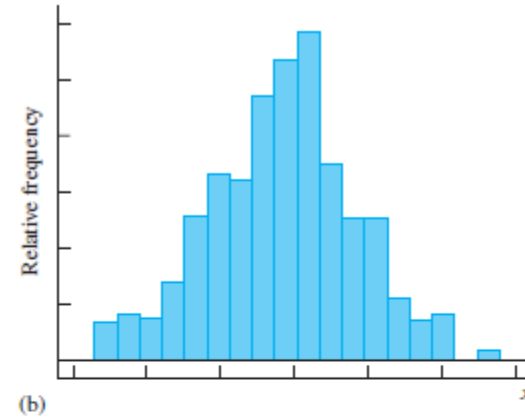
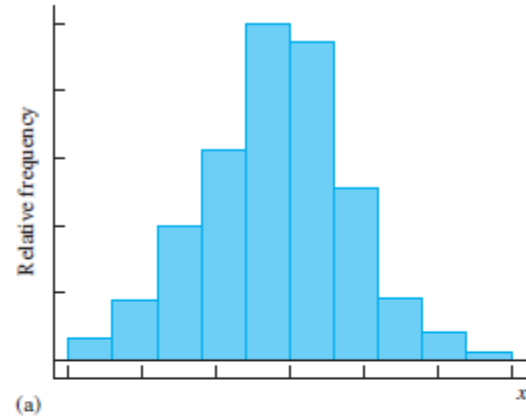
$$P(k = 4) = \frac{1000!}{4! (996)!} (0.001)^4 (1 - 0.001)^{996} = 0.01529$$

Alternatively because mean of this binomial distribution is small $\mu = np = 1$, we can approximate it with a Poisson distribution with $\mu = 1$.

$$P(x = 4) = \frac{1^4 e^{-1}}{4!} = 0.01533$$

Probability Distributions for Continuous Random Variables

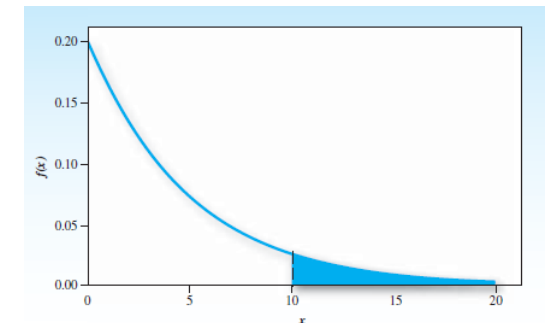
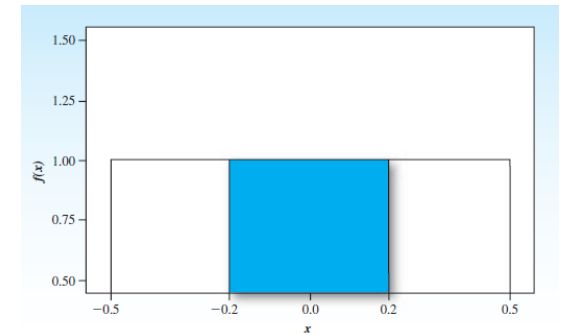
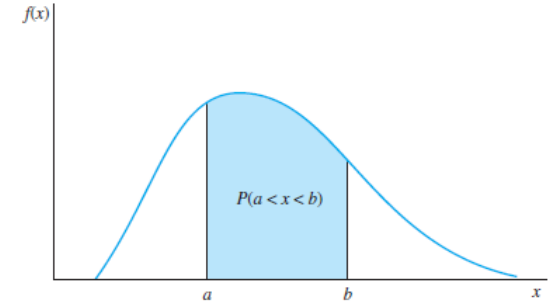
A **continuous variable** can take unlimited number of values in a given range.



Relative frequency histograms for increasingly large number of samples of a continuous random variable

Probability Distributions for Continuous Random Variables

- Characteristics of a probability distribution $f(x)$
 - The area under the distribution equals to 1
 - $P(a < x < b)$ equals to the area between a and b
 - $P(x = a) = 0$ because there is no area above $x = a$
 - $P(x \geq a) = P(x > a)$ and $P(x \leq a) = P(x < a)$
- Examples of continuous random variables
 - Rounding error x to nearest integer of values between -0.5 and 0.5 has a uniform distribution (because they *all* become 0), $f(x)=1$
 - The wait time x at a supermarket checkout may follow an exponential distribution, $f(x) = .2e^{-.2x}$

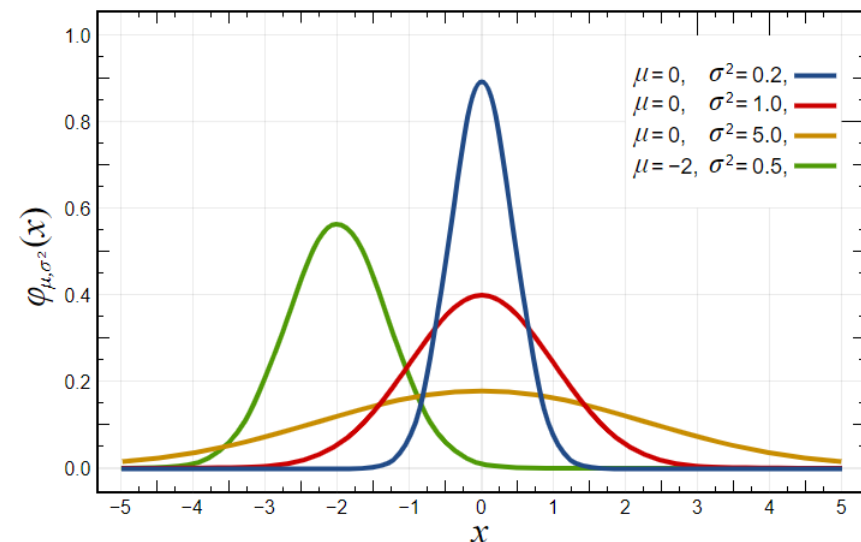
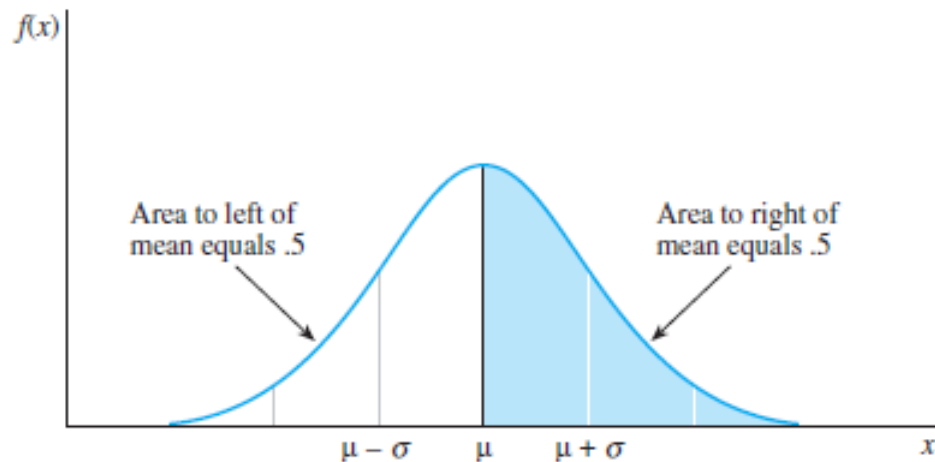


The Normal Probability Distribution

- Many continuous random variables in nature (weight, height, time) can be well described by **normal probability distribution** (thus the name *normal*)

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)} \quad \text{for } -\infty \leq x \leq \infty$$

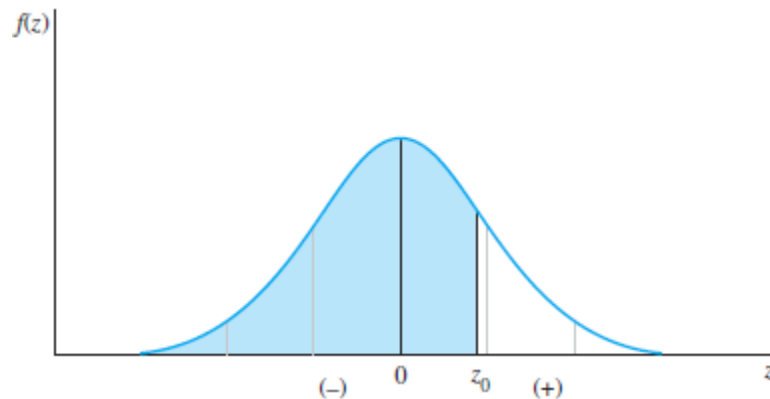
- The distribution is symmetric about the mean μ , which is also the mode and median.
- The shape of the curve is determined by the population standard deviation σ .
- The Empirical Rule!



Source: https://en.wikipedia.org/wiki/Normal_distribution

Tabulated Areas of the Normal Probability Distribution

- **Standardized normal random variable** z is defined as $z = \frac{x-\mu}{\sigma}$, essentially the number of σ the variable x lies to the left or right of μ .
- The probability distribution of z is called the **standardized normal distribution** because the mean is 0 and standard deviation is 1.
- Value of $P(z < z_0)$ is the shaded area and is tabulated, an extract of which is given below.



z	.00	.01	.02	.03	.0409
0.0	.5000	.5040	.5080	.5120	.5160		
0.1	.5398	.5438	.5478	.5517	.5557		
0.2	.5793	.5832	.5871	.5910	.5948		
0.3	.6179	.6217	.6255	.6293	.6331		
0.4	.6554	.6591	.6628	.6664	.67006879
0.5	.6915	.	.	.			
0.6	.7257	.	.	.			
0.7	.7580	.	.	.			
0.8	.7881						
0.9	.8159						
.	.						
.	.						
.	.						
2.0	.9772						

Example: $P(z < 0.41) = ?$

Example. Let x be a normal distributed variable with $\mu = 10$ and $\sigma = 2$. Find the probability of x lies between 9.4 and 10.6.

Solution hints

Calculate the standardized normal random variable z_1 corresponding to $x_1 = 9.4$

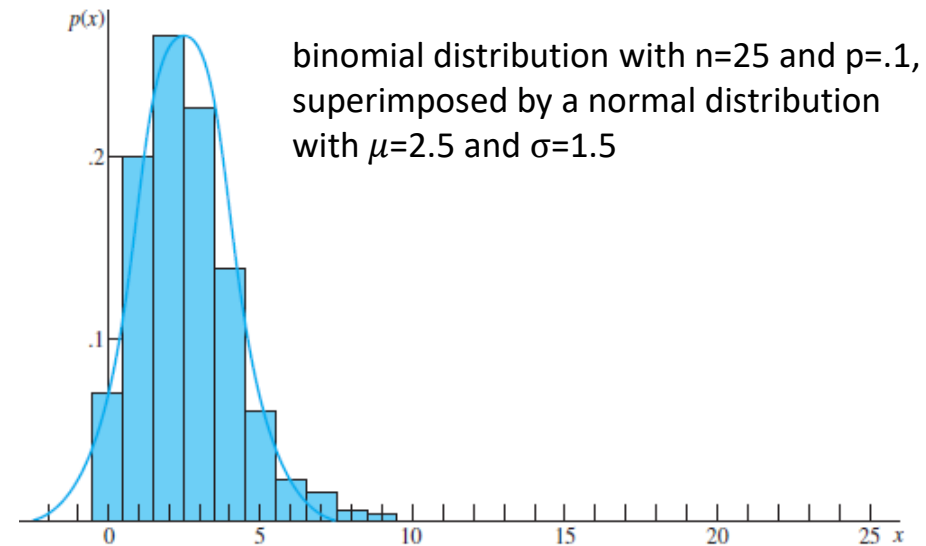
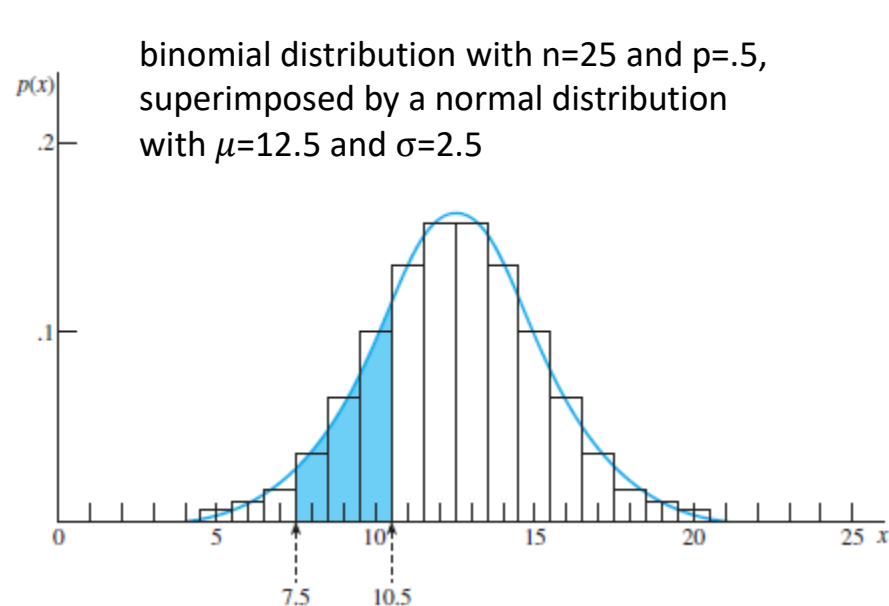
Calculate the standardized normal random variable z_2 corresponding to $x_2 = 10.6$

$$P(x_1 < x < x_2) = P(z_1 < z < z_2)$$

z	.00	.01	z	.00	.01
-3.4	.0003	.0003	0.0	.5000	.5040
-3.3	.0005	.0005	0.1	.5398	.5438
-3.2	.0007	.0007	0.2	.5793	.5832
-3.1	.0010	.0009	0.3	.6179	.6217
-3.0	.0013	.0013	0.4	.6554	.6591
-2.9	.0019	.0018	0.5	.6915	.6950
-2.8	.0026	.0025	0.6	.7257	.7291
-2.7	.0035	.0034	0.7	.7580	.7611
-2.6	.0047	.0045	0.8	.7881	.7910
-2.5	.0062	.0060	0.9	.8159	.8186
-2.4	.0082	.0080	1.0	.8413	.8438
-2.3	.0107	.0104	1.1	.8643	.8665
-2.2	.0139	.0136	1.2	.8849	.8869
-2.1	.0179	.0174	1.3	.9032	.9049
-2.0	.0228	.0222	1.4	.9192	.9207
-1.9	.0287	.0281	1.5	.9332	.9345
-1.8	.0359	.0351	1.6	.9452	.9463
-1.7	.0446	.0436	1.7	.9554	.9564
-1.6	.0548	.0537	1.8	.9641	.9649
-1.5	.0668	.0655	1.9	.9713	.9719
-1.4	.0808	.0793	2.0	.9772	.9778
-1.3	.0968	.0951	2.1	.9821	.9826
-1.2	.1151	.1131	2.2	.9861	.9864
-1.1	.1357	.1335	2.3	.9893	.9896
-1.0	.1587	.1562	2.4	.9918	.9920
-0.9	.1841	.1814	2.5	.9938	.9940
-0.8	.2119	.2090	2.6	.9953	.9955
-0.7	.2420	.2389	2.7	.9965	.9966
-0.6	.2743	.2709	2.8	.9974	.9975
-0.5	.3085	.3050	2.9	.9981	.9982
-0.4	.3446	.3409	3.0	.9987	.9987
-0.3	.3821	.3783	3.1	.9990	.9991
-0.2	.4207	.4168	3.2	.9993	.9993
-0.1	.4602	.4562	3.3	.9995	.9995
-0.0	.5000	.4960	3.4	.9997	.9997

The Normal Approximation to the Binomial Distribution

- Probability of a binomial variable x can be calculated via
 - the binomial formula or corresponding binomial tables, or
 - the Poisson probabilities for $np < 7$.
- When $\mu = np$ of a binomial distribution is large, the normal probability with $\mu = np$ and standard deviation $\sigma = \sqrt{np(1-p)}$ can be used as an approximation.
- For this approximation to hold, n must be large and p is not too close to 0 or 1.



The Normal Approximation to the Binomial Distribution

Rule of thumb. A normal distribution approximates well a binomial distribution if both $np > 5$ AND $n(1-p) > 5$ (because the binomial distribution is *fairly symmetric*)

Example. A random sample of 1000 fuses were tested. Assuming defect probability is 0.02. What is the probability of having more than 27 fuses defected.

Solution. This is a binomial experiment with $\mu = np = 20$ and $\sigma = \sqrt{np(1-p)} = 4.43$.

However because both np and $n(1-p)$ is larger than 5, we can approximate it by a normal distribution with $\mu = 20$ and $\sigma = 4.43$.

Because of **continuity correction**, the normal area corresponding to $P(x \geq 27)$ is the area to the right of $x=26.5$.

Standardised value of $x=26.5$ is $z_0 = \frac{26.5-20}{4.43} = 1.47$

Therefore $P(x \geq 27) \approx P(z > 1.47) = 1 - .9292 = .0708$