



STATISTICS FOR DATA SCIENCE

Continuity Correction

D. Uma

Department of Computer Science and Engineering

STATISTICS FOR DATA SCIENCE

Continuity Correction

D. Uma

- ✓ Continuity Correction and Why do we need it?
- ✓ Continuity Correction Factor.
- ✓ Normal Approximation to Binomial.
- ✓ Normal Approximation to Poisson.

Assume that, a surgeon is very skillful because his surgeries are 90% success and let us assume he performs the procedure on 12 patients.

If we have to find the probability of exactly four successful surgeries. It becomes more easier and plausible to do.

$$P(X = 4)$$

X : No of success

What if we have to find the probability of more than 200 surgeries?

$$P(X > 200)$$

Here we end up using binomial formula for 200 times which is not practical of course.

$$\begin{aligned} &= 1 - P(X \leq 200) \\ &= 1 - [P(X=0) + \dots + P(X=200)] \\ &= 1 - \end{aligned}$$

A quick approach to make it more efficient is to use normal distribution to approximate the binomial distribution resulting in efficiency of the results.

Bernoulli X : getting successes

Problem: Find the probability that the number of heads is greater than 60 out of 100 trials.

Bin

$$P(X > 60) = 1 - P(X \leq 60) \quad X: \text{No. of successes in } n \text{ trials}$$

$$\underline{\text{Bin}(n, p)} \approx \text{Normal}$$

$$\approx N(\text{Mean}, \text{Var})$$

Discrete

Continuous

$$\begin{aligned} &= 1 - \left[P(X=0) + P(X=1) + \dots + P(X=60) \right] \\ &= 1 - \left[\underbrace{100C_0 (0.5)^0 (0.5)^{100}} + \dots + \dots \right] \end{aligned}$$

Problem: Use normal curve to approximate the probability that the number of heads is greater than 60. $P(X \geq 60)$

By computing probability that corresponds to $X \sim \text{Bin}(100, 0.5)$

$$P(X = x) = \begin{cases} \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} & x = 0, 1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$$P(X \geq 60) = P(X = 60) + \dots + P(X = 100)$$

$$\begin{aligned} &= \frac{100!}{60!(100-60)!} (0.5)^{60} (1-0.5)^{100-60} + \dots \\ &\quad + \frac{100!}{100!(100-100)!} (0.5)^{100} (1-0.5)^{100-100} \\ &= \underline{0.0284} \end{aligned}$$

The actual probability of $P(X \geq 60)$ is 0.0284.

- If we want to **employ a continuous** (normal) **distribution** to approximate any discrete distribution (like binomial and Poisson), **continuity correction** should be used.
- It is used **to make adjustments** and it can improve the **accuracy of the approximation**.

Why do we need Continuity Correction?

- The **discrete** random variables can take **only integer values**.
- The **continuous** random variable can take **real values** and can be **used to approximate** any **discrete values** within the interval around specified values.
- **More accurate approximations** can be obtained by using continuity correction.



| Probabilities | Discrete | Continuity Correction | Continuous |
|---------------|---------------|------------------------|--------------------|
| $P(X = n)$ | $P(X=5)$ | $P(n-0.5 < X < n+0.5)$ | $P(4.5 < X < 5.5)$ |
| $P(X > n)$ | $P(X > 5)$ | $P(X > n+0.5)$ | $P(X > 5.5)$ |
| $P(X \geq n)$ | $P(X \geq 5)$ | $P(X \geq n-0.5)$ | $P(X \geq 4.5)$ |
| $P(X < n)$ | $P(X < 5)$ | $P(X < n-0.5)$ | $P(X < 4.5)$ |
| $P(X \leq n)$ | $P(X \leq 5)$ | $P(X \leq n+0.5)$ | $P(X \leq 5.5)$ |

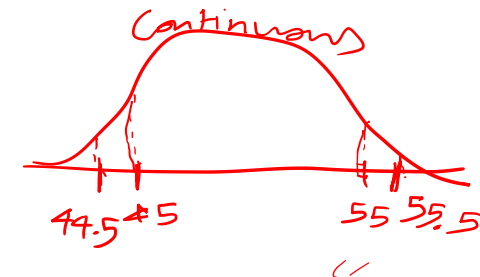
Given: $P(45 \leq X \leq 55)$ where X is a DRV.
(include end points)

Correction required when we approximate with continuous:
 $P(44.5 \leq X \leq 55.5) = P(X \leq 55.5) - P(X \leq 44.5)$
 (to include end points)

■ Note: Equality makes no difference

where X is a CRV.

$< \leq$
 $n-0.5$ $n+0.5$
 n
 x is discrete r.v.
 $P(X=n)$ Discrete can be found
 $P(X=x) = 0$
 where x is a continuous r.v.



| Probabilities | Discrete | Continuity Correction | Continuous |
|---------------|---------------|----------------------------|--------------------|
| $P(X = n)$ | $P(X = 5)$ | $P(n - 0.5 < X < n + 0.5)$ | $P(4.5 < X < 5.5)$ |
| $P(X > n)$ | $P(X > 5)$ | $P(X > n + 0.5)$ | $P(X > 5.5)$ |
| $P(X \geq n)$ | $P(X \geq 5)$ | $P(X > n - 0.5)$ | $P(X > 4.5)$ |
| $P(X < n)$ | $P(X < 5)$ | $P(X < n - 0.5)$ | $P(X < 4.5)$ |
| $P(X \leq n)$ | $P(X \leq 5)$ | $P(X < n + 0.5)$ | $P(X < 5.5)$ |

- Note: Equality makes no difference

Let X be a Binomial r.v with
 $X \sim \text{Bin}(n, p)$

$$X \sim \text{Bin}(100, 0.5)$$

$$X \sim \text{Bin}(n, p) \approx X \sim N(np, np(1-p))$$

$$p = 0.5$$

when n is large & p is small.

$$X \sim N(np, npq)$$

Binomial \approx Normal $\mu = n * p$

When $np > 10$ and $nq > 10$ Bin \approx Normal

$np > 5$
 Mean no. of successes

$nq > 5$
 Mean no. of failures



Actual Dist'n $n=100$ $x \sim \text{Bin}(100, 0.5)$

① Binomial

$$P(X \geq 60) = \boxed{0.0284}$$

② Normal without CC $Z = \frac{x-\mu}{\sigma}$ Nor

$$P(X \geq 60) = P\left(\frac{x-\mu}{\sigma} \geq \frac{60-50}{5}\right)$$

$$= P(Z \geq 2)$$

$$= 1 - P(Z \leq 2)$$

$$P(X \geq 60) = 1 - 0.9772$$

$$= \boxed{0.0228}$$

$$\Rightarrow X \sim N(np, npq)$$

$$X \sim N(100 \times 0.5, 100 \times 0.5 \times 0.5)$$

$$X \sim N(50, 5^2)$$

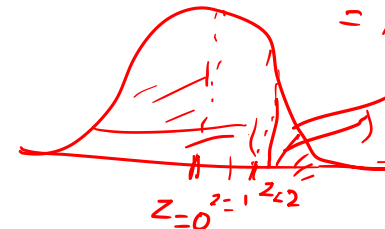
μ σ^2

③ Normal Distribution with CC

$$P(X \geq 59.5)$$

$$= P\left(\frac{x-\mu}{\sigma} \geq \frac{59.5-50}{5}\right)$$

$$= P(Z \geq 1.9)$$



$$= 1 - P(Z \leq 1.9)$$

$$= 1 - 0.9713$$

$$= \underline{\underline{0.0287}}$$

Discrete \Rightarrow Continuous
Distribution \Rightarrow Dist'n

If a fair coin is tossed 100 times, use the normal curve to approximate the probability that the number of heads is between 45 and 55 inclusive.

$$X \sim \text{Bin}(\overset{\text{Discrete}}{100}, 0.5) \Rightarrow \overset{\text{Continuous}}{X} \sim N(50, 5^2)$$

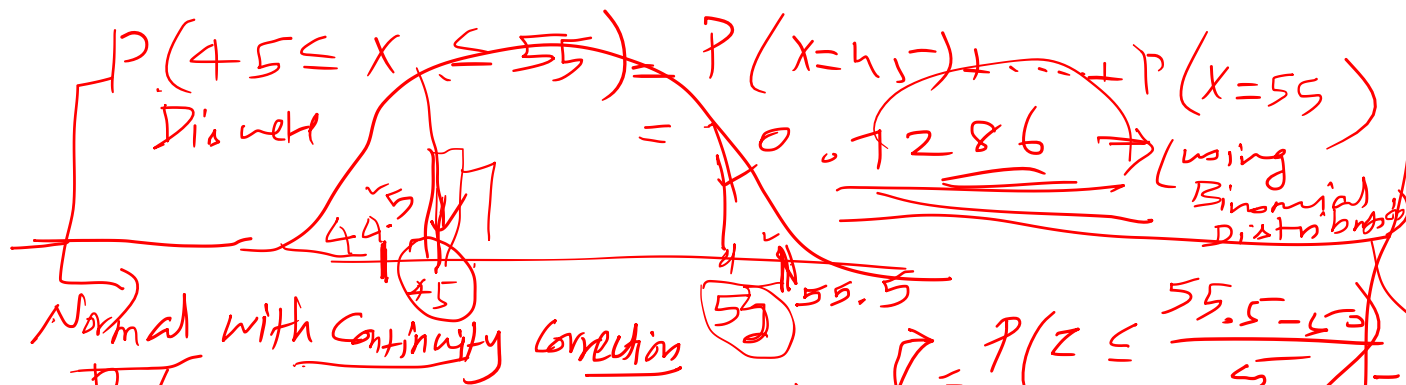
$\mu = np$ $\sigma^2 = np(1-p)$ $np > 5$ & $nq > 5$
 $np > 10$ & $nq > 10$

$P(45 \leq X \leq 55)$ → Discrete Distribution

Discrete

$$P(45 \leq X \leq 55) = P(X=45) + \dots + P(X=55)$$

$= 0.7286$ (using Binomial Distribution)



$$P(44.5 \leq X \leq 55.5)$$

$$= P(X \leq 55.5) - P(X \leq 44.5)$$

$$= P\left(Z \leq \frac{55.5 - 50}{5}\right) - P\left(Z \leq \frac{44.5 - 50}{5}\right)$$

$$= P(Z \leq 1.1) - P(Z \leq -1.1)$$

$$= 0.8643 - 0.1357$$

$$= 0.7286$$

Normal $Z = \frac{x - \mu}{\sigma}$
without Continuity Correction

$$P(45 \leq X \leq 55)$$

$$= P(X \leq 55) - P(X \leq 45)$$

$$= P\left(Z \leq \frac{55 - 50}{5}\right) - P\left(Z \leq \frac{45 - 50}{5}\right)$$

$$= P(Z \leq 1) - P(Z \leq -1)$$

$$= 0.6826$$

(using normal distribution without CC)

(using Normal distribution with CC)

STATISTICS FOR DATA SCIENCE

Solution for $P(X \geq 60)$ after continuity correction

By computing probability that corresponds to $X \sim N(50, 25)$

Binomial

$$P(X \geq 60) = P(X=61) + \dots + P(X=100) \text{ or } 1 - P(X < 60)$$

Normal without CC

$59.5 - X = 60$

$0 \quad z=2$

$$= 1 - [P(X=0) + \dots + P(X=59)]$$
$$= 0.0284$$

$$P(X \geq 60) = P\left(Z \geq \frac{60 - 50}{5}\right) = P(Z \geq 2) = 1 - P(Z < 2) = 0.0228$$

Normal with CC

$$P(X \geq 59.5) = P\left(Z \geq \frac{59.5 - 50}{5}\right) = P(Z \geq 1.9)$$
$$= 1 - P(Z < 1.9)$$
$$= 0.0287$$



$$\left. \begin{array}{l} X \sim \text{Bin}(n, p) \\ n - \text{large} \\ p - \text{small} \\ np > 10 \quad \& \quad nq > 10 \end{array} \right\} \approx$$

$$\begin{aligned} X &\sim N(\mu, \sigma^2) \\ &X \sim N(np, np(1-p)) \\ X &= Y_1 + Y_2 + \dots + Y_n \\ \text{where } Y_i &\sim \text{Bernoulli}(p) \end{aligned}$$

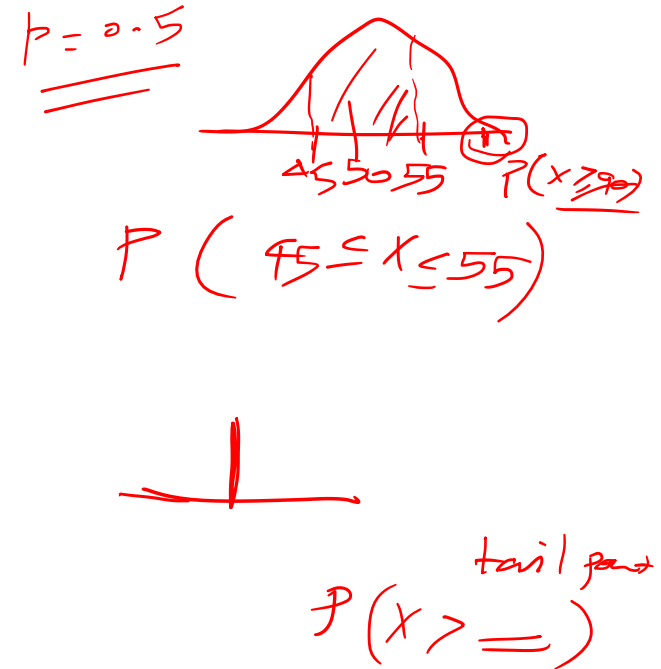
When p is unknown

$$\hat{p} = \frac{X}{n} = \frac{Y_1 + Y_2 + \dots + Y_n}{n} = \bar{Y} \Rightarrow \text{sample mean}$$
$$\sqrt{\hat{p}} \sim N\left(p, \frac{p(1-p)}{n}\right)$$



- The continuity correction improves the accuracy of the normal approximation to the binomial distribution when p is small and n is large.
- The continuity correction can in some cases reduce the accuracy of the normal approximation.
- It occurs when there is some degree of skewness in the distribution and when p is not equal to 0.5 and computing probability that corresponds to an area in the tail of the distribution.

$$p \neq 0.5$$



STATISTICS FOR DATA SCIENCE

Normal Approximation to Poisson



PES
UNIVERSITY
ONLINE

$$X \sim \text{Poisson}(\lambda) \stackrel{\text{mean rate}}{\sim} X \sim N(\lambda, \lambda) \quad \text{when } \lambda > 10$$

Mean = λ
Variance = λ

λ : Avg. rate at which

X : No. of events occurring

$$P(X \geq 60) = P(X=60) + P(X=61) + \dots + \left[\frac{e^{-\lambda} \cdot \lambda^x}{x!} \right] \approx 1 - P(X < 60)$$
$$X \sim N(15, 15)$$

- For areas that include the central part of the curve, the continuity correction generally improves the normal approximation.
- But, for areas in the tails, the continuity correction sometimes makes the approximation worse.

Note: If the area is in the tails, then the continuity correction may make the approximation worse.

The number of hits on a website follows a Poisson distribution, with a mean of 27 hits per hour. Find the probability that there will be 90 or more hits in three hours.

$$X \sim \text{Poisson}(27) \Rightarrow X \sim \text{Poisson}(3 \times 27)$$

$$X \sim \text{Poisson}(81)$$

$$X \sim N(81, 9^2)$$

9-SD
81-81

$$P(X \geq 90) = 0.1718$$

Normal without CC

$$P(X \geq 90) = P\left(Z \geq \frac{90 - 81}{9}\right)$$

$$= P(Z \geq 1)$$

$$= 1 - P(Z \leq 1) = 0.1587$$

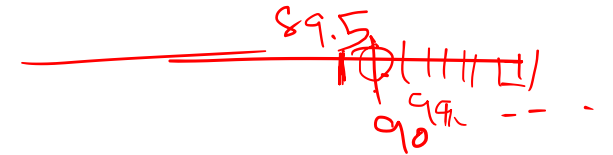
Normal with CC:

$$P(X \geq 90) \approx P(X \geq 89.5)$$

$$= P\left(Z \geq \frac{89.5 - 81}{9}\right)$$

$$= P(Z \geq 1)$$

$$= 1 - P(Z < 1) = 0.1587$$





THANK YOU

D. Uma

Computer Science and Engineering

umaprabha@pes.edu

+91 99 7251 5335