

Assignment Project Exam Help

<https://powcoder.com>



Add WeChat powcoder

Statistics (MAST20005) &
Elements of Statistics
(MAST90058)

School of Mathematics and Statistics
University of Melbourne

Semester 2, 2022

Aims of this module

- Brief information about this subject
- Brief revision of some prerequisite knowledge (probability)
- Introduce some basic elements of statistics, data analysis and visualisation

<https://powcoder.com>

Add WeChat powcoder

Outline

Assignment Project Exam Help

Subject information

<https://powcoder.com>

Descriptive statistics

Add WeChat powcoder

Basic data visualisations

What is statistics?

Assignment Project Exam Help

<https://powcoder.com>

Let's see some examples...

Add WeChat powcoder



Assignment Project Exam Help

Bureau Home > Australia > Victoria > Forecasts > Melbourne Forecast

Melbourne Forecast



[View the current warnings for Victoria](#)

Forecast issued at 10:00 am EST on Tuesday 17 July 2018

Forecast for the rest of Tuesday



Max 16

Showers developing. Windy.

Possible rainfall: 1 to 3 mm

Chance of any rain: 80% 

Melbourne area

Sunny morning. High (80%) chance of showers during the afternoon and evening and the threat of hail and thunder. Winds northerly 25 to 40 km/h increasing to 40 to 55 km/h this morning then turning west to northwesterly 30 to 45 km/h in the late afternoon.

Add WeChat powcoder

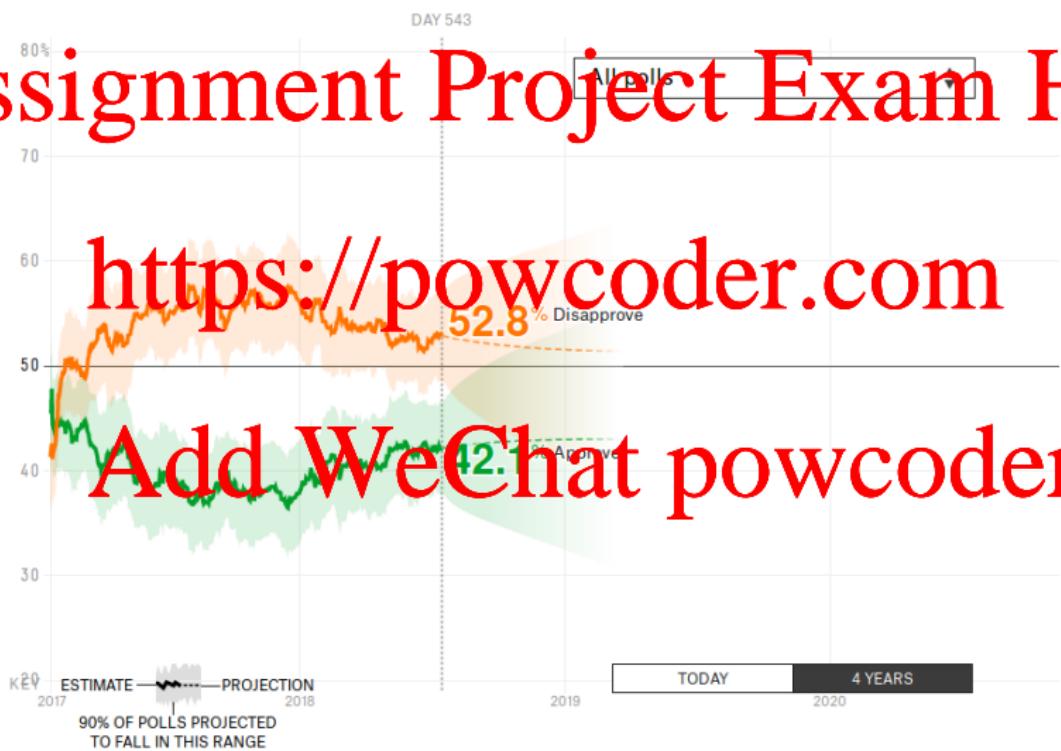
How unpopular is Donald Trump?

An updating calculation of the president's approval rating, accounting for each poll's quality, recency, sample size and partisan lean. [How this works »](#)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Poll of polls: who will win the Australian election?

Assignment Project Exam Help

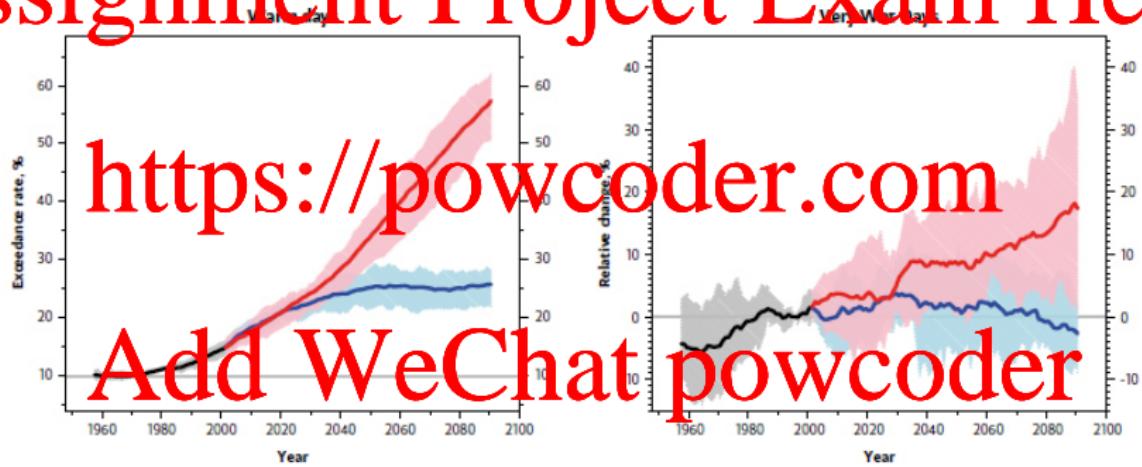


<https://powcoder.com>

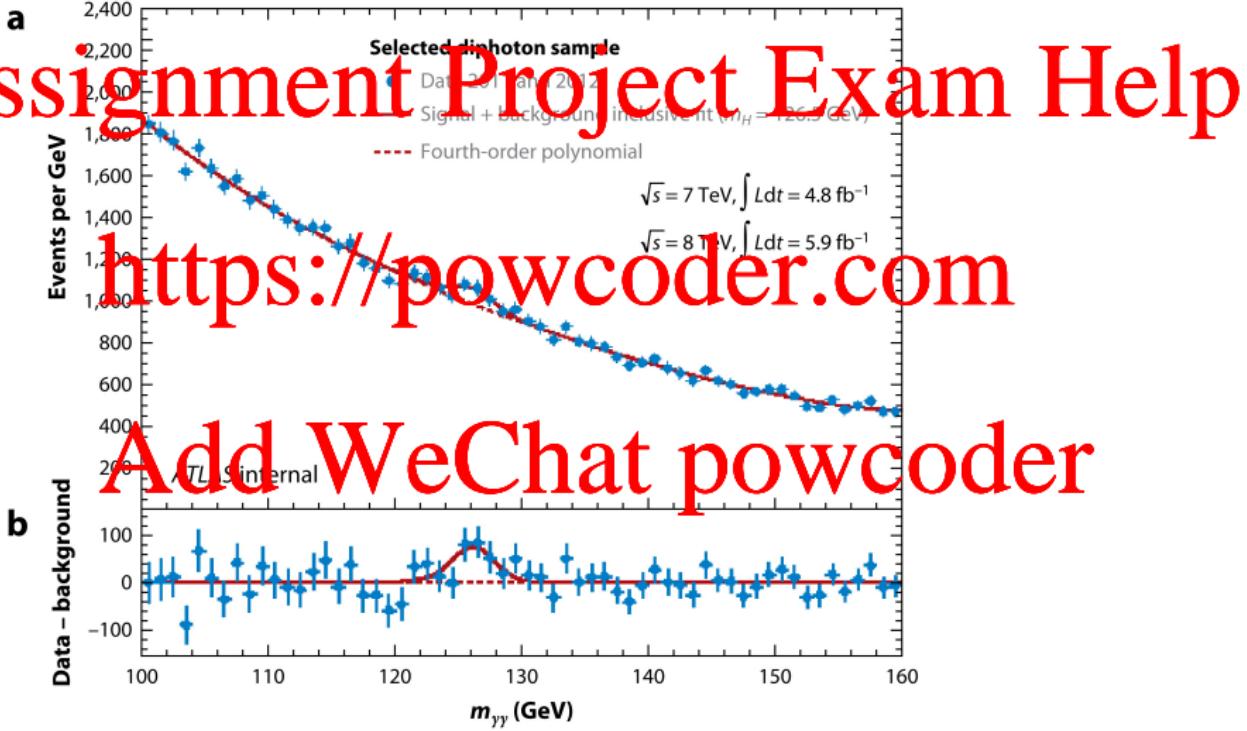
Add WeChat powcoder

Climate change modelling

Assignment Project Exam Help



Discovery of the Higgs Boson (the ‘God Particle’)

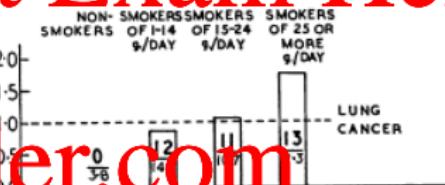


Smoking leads to lung cancer

TABLE IV—Standardized Death Rate Per Annum Per 1,000 Men Aged 35 Years and Above in Relation to the Most Recent Amount of Tobacco Smoked

Cause of Death	No. of Deaths Recorded	Death Rates of Men Smoking a Daily Average of			Death Rate of All Men
		1 g.-	15 g.-	25 g.+	
Lung cancer . . .	36*	0.00	0.48	0.67	1.14
Other cancers . . .	92	2.32	1.41	1.50	1.91
Respiratory disease (other than cancer)	54	0.86	0.85	1.01	0.77
Coronary thrombosis	21	0.85	3.01	4.1	1.15
Other cardiovascular diseases	17	1.2	2.07	1.6	1.78
Other diseases . . .	24	4.27	4.67	3.1	4.52
All causes . . .	789	13.61	13.42	13.38	16.30
					14.00

* 1 case in which lung cancer was recorded as a contributory but not a direct cause of death has been entered in both groups.



Add WeChat powcoder

Richard Doll & A. Bradford Hill

"The Mortality of Doctors in Relation to Their Smoking Habits"
British Medical Journal (1954)

Assignment Project Exam Help
<https://powcoder.com>

A/B testing for websites

Assignment Project Exam Help



https://powcoder.com

Add WeChat powcoder

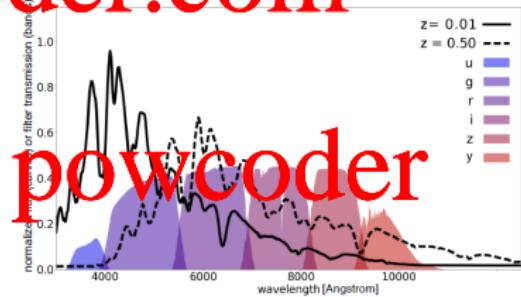
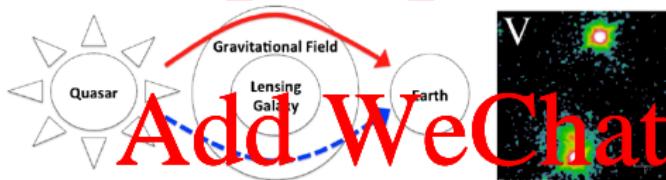


41 shades of blue?!

Assignment Project Exam Help

The Photometric LSST Astronomical
Time-series Classification Challenge

<https://powcoder.com>



Add WeChat powcoder



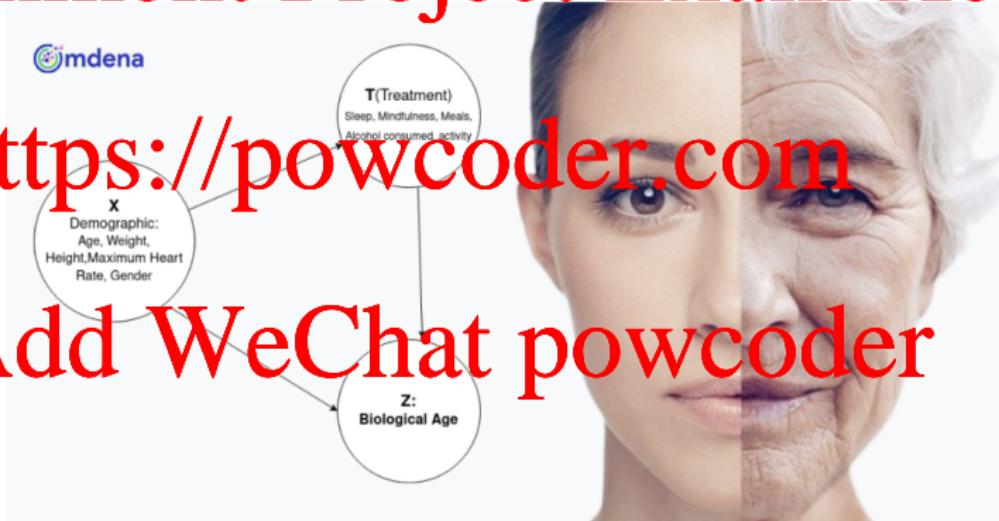
PROJECTS

AI ENGINEERS ▾

ORGANIZATIONS ▾

A^ssignment Project Exam Help

<https://powcoder.com>
Add WeChat powcoder



THE INDEPENDENT



(Ireland, €1) 70p

Thursday 7 June 2007

www.independent.co.uk

PA 14.0895 6461

Assignment Project Exam Help

Hypertension

High blood pressure affects 16 million people in Britain. Can lead to stroke, heart disease and kidney failure

Type 1 diabetes

Diabetic condition in which sufferers have to inject insulin. Affects 350,000 people in UK

Type 2 diabetes

Almost 2 million Britons are affected by this late-onset disease, which is linked with the growing obesity epidemic

Add WeChat powcoder

Add WeChat powcoder

THE GENETIC REVOLUTION

DISCOVERY OF GENES RESPONSIBLE FOR SEVEN OF THE MOST COMMON ILLNESSES OFFERS HOPE TO MILLIONS OF SUFFERERS

FULL STORY, PAGE 2

Web analytics

Assignment Project Exam Help

<https://powcoder.com>



2011 Black Friday Shoppers by State

Week ending November 26, 2011, compared with the Online Population

State	Visits Share Black Friday 2011 Shoppers	Visits Share Online Population	Index
1 California	10.05%	1.05%	68
2 Texas	7.5%	0.75%	82
3 New York	5.23%	6.08%	99
4 Illinois	4.88%	4.34%	104
5 Florida	4.88%	5.92%	102
6 Ohio	4.49%	3.66%	77
7 Pennsylvania	4.03%	4.01%	103
8 Georgia	3.57%	3.23%	83
9 Michigan	3.51%	3.23%	84
10 North Carolina	3.35%	2.92%	107

2011 Black Friday Shoppers by DMA

Week ending November 26, 2011, compared with the Online Population

DMA	Visits Share Black Friday 2011 Shoppers	Visits Share Online Population	Index
1 New York, NY	4.4%	0.4%	68
2 Los Angeles, CA	4.01%	0.41%	82
3 Chicago, IL	3.01%	0.3%	99
4 Dallas-Ft. Worth, TX	2.29%	0.22%	104
5 Atlanta, GA	2.12%	0.21%	102
6 Philadelphia, PA	1.98%	0.19%	77
7 Houston, TX	1.90%	0.18%	103
8 Boston et al, MA-NH	1.74%	0.17%	83
9 Washington et al, DC-MD	1.72%	0.17%	84
10 Minneapolis-St. Paul, MN	1.61%	0.16%	107

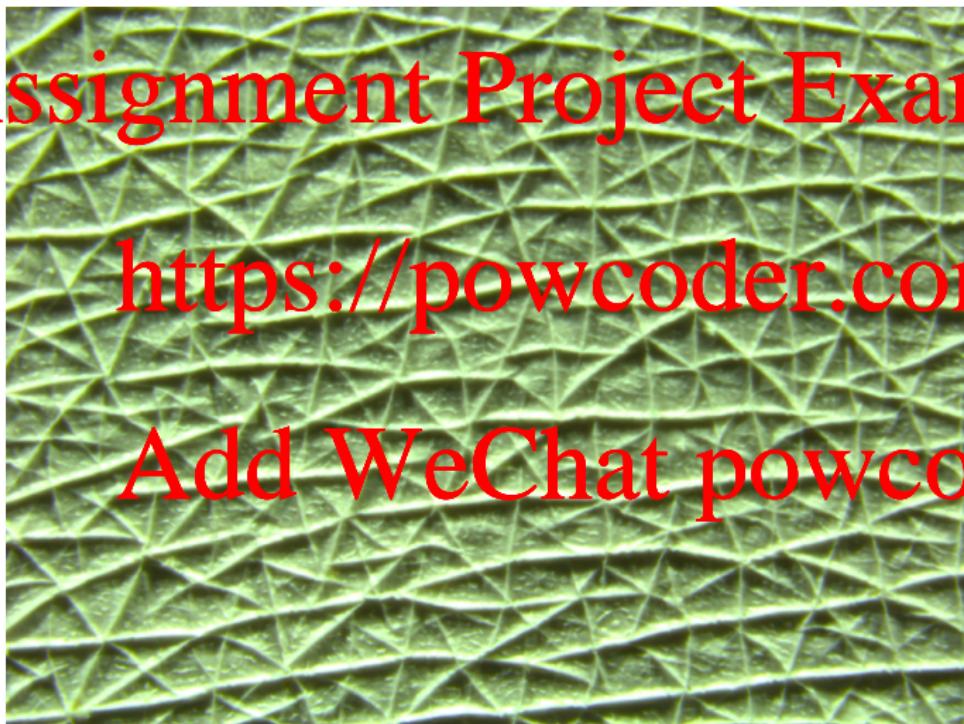
Add WeChat powcoder

Skin texture image analysis

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



NEWS MAGAZINE

[Home](#) | [World](#) | [Asia](#) | [India](#) | [China](#) | [UK](#) | [Business](#) | [Health](#) | [Science/Environment](#) | [Technology](#) | [Entertainment](#)[Magazine](#) | [In Pictures](#) | [Also in the News](#) | [Editors' Blog](#) | [Have Your Say](#) | [World News TV](#) | [World Service](#)

1 February 2014 Last updated at 08:47



Assignment Project Exam Help

A statistically modelled wedding

By Ruth Alexander
BBC News



<https://powcoder.com>

It's the classic dilemma for anyone planning a wedding. The list of friends and family you'd like to invite is seemingly endless. Your budget is not, and neither is the venue's floor space. Could one couple have found a solution via statistical modelling?

Damjan Vukcevic and Joan Ko, planning their wedding in Melbourne, Australia, were struggling to draw up an invitation list of family and friends in places as far-flung as Serbia, Taiwan, the UK and the US.

In today's Magazine

The world's first church-mosque-synagogue?

The face of British

Goals of statistics

- Answer questions using ~~data~~
- Evaluate ~~evidence~~

Assignment Project Exam Help

- Optimise ~~study~~ design
- Make ~~decisions~~

<https://powcoder.com>

And, importantly:

- Clarify ~~assumptions~~
- Quantify ~~uncertainty~~

Add WeChat powcoder

Why study statistics?

Assignment Project Exam Help

"The best thing about being a statistician is that you get to play in everyone's backyard."

—John W. Tukey (1915–2000)

<https://powcoder.com>

"I keep saying the sexy job in the next ten years will be statisticians... The ability to take data – to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it's going to be a hugely important skill in the next decades..."

—Hal Varian, Google's Chief Economist, Jan 2009

Add WeChat powcoder

The best job

U.S. News Best Business Jobs in 2022:

1. Medical and Health Services Manager
2. Financial Manager
3. Statistician

<https://powcoder.com>

CareerCast (recruitment website) Best Jobs of 2021:

1. Data Scientist
 2. Genetic Counselor
 3. Statistician
- Add WeChat powcoder

Subject overview

Statistics (MAST20005) Elements of Statistics (MAST90058)

Assignment Project Exam Help

These subjects introduce the basic elements of **statistical modelling**, **statistical computation** and **data analysis**. They demonstrate that many commonly used statistical procedures arise as applications of a common theory. They are an entry point to further study of both **mathematical** and **applied** statistics, as well as broader data science.

Students will develop the ability to fit statistical models to data, estimate parameters of interest and test hypotheses. Both classical and Bayesian approaches will be covered. The importance of the underlying **mathematical theory of statistics** and the use of **modern statistical software** will be emphasised.

Add WeChat powcoder

Joint teaching

MAST20005 and MAST90058 share the same lectures but have separate tutorials and lab classes. The teaching and assessment material for both subjects will overlap significantly.

<https://powcoder.com>

Add WeChat powcoder

Subject website (LMS)

- Full information is on the subject website, available through the Learning Management System (LMS).
- Only a brief overview is covered in these notes. Please read all of the info on the LMS as well.
- New material (e.g. problem sets, assignments, solutions) and announcements will appear regularly on the LMS.

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

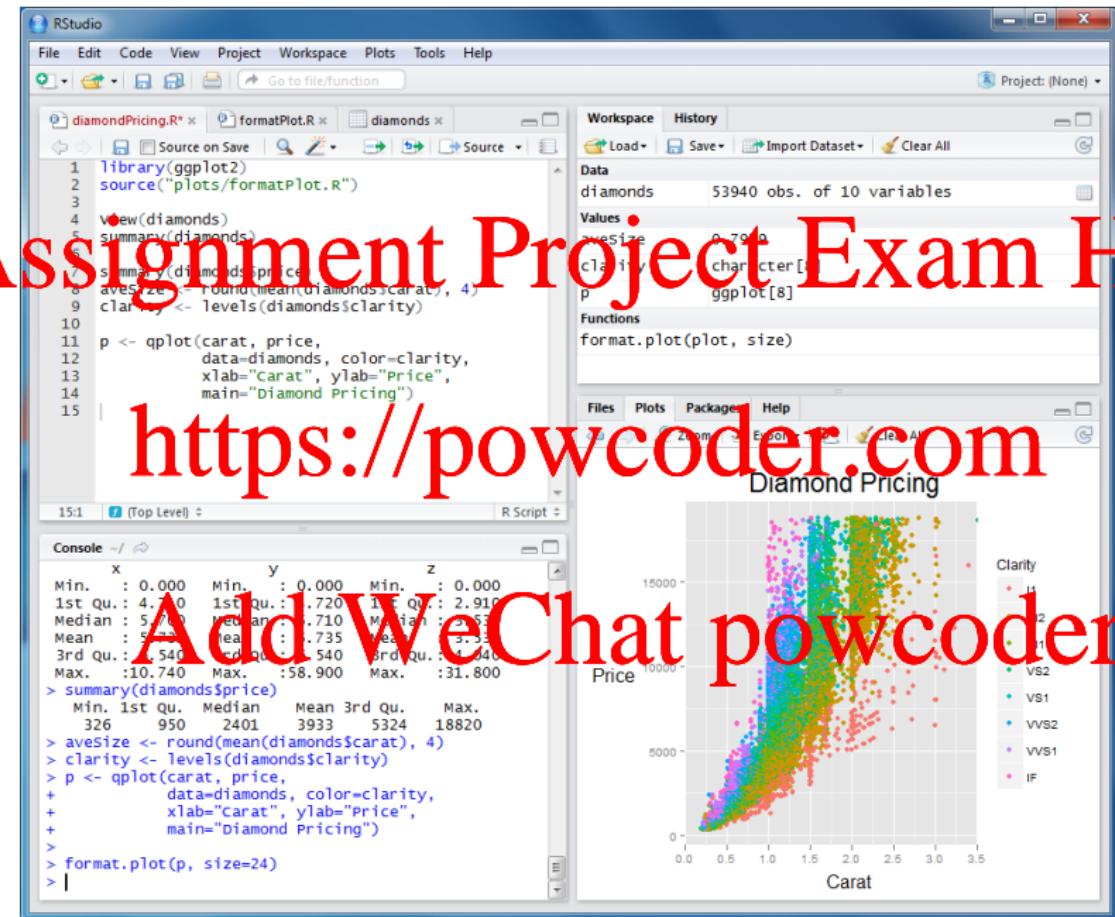
- This subject introduces basic statistical computing and programming skills.
- We make extensive use of the R statistical software environment.
- Knowledge of R will be essential for some of the tutorial problems, assignment questions and will also be examined.
- We will use the RStudio program as a convenient interface with R.

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>



Add WeChat powcoder

Staff contacts

Subject coordinator / Lecturer

Dr Damjan Vukcevic <damjan.vukcevic@unimelb.edu.au>

Lecturer

Dr Wei Huang <wei.huang@unimelb.edu.au>

Tutorial coordinator
<https://powcoder.com>

Dr Rekha Mathur <rekha.mathur@unimelb.edu.au>

See the LMS for details of consultation times

Add WeChat powcoder

Discussion forum

- Access via the LMS
- Post any general questions on the forum

- Do **not** send them by email to staff
- You can answer each others' questions
- Staff will also help to answer questions

<https://powcoder.com>

Add WeChat powcoder

Student representatives

Student representatives assist the teaching staff to ensure good communication and feedback from students.

Assignment Project Exam Help

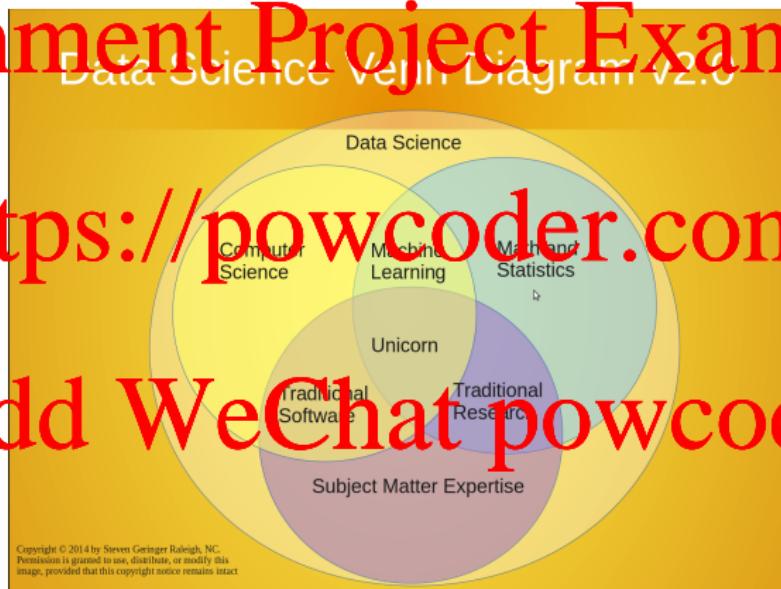
See the LMS to find the contact details of your representatives.

<https://powcoder.com>

Add WeChat powcoder

What is Data Science?

Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder

Data science is a 'team sport'

How to succeed in statistics / data science?

- Get experience with ~~real data~~
- Develop your computational skills, ~~Learn R~~

- Understand the ~~mathematical theory~~
- Collaborate with others, ~~use the discussion forum~~

<https://powcoder.com>

Add WeChat powcoder

This subject is challenging

- It is mathematical

- Manipulating equations

- Calculus

- Probability

- Proofs

- But the 'real' world also matters

- Context can 'trump' mathematics

- More than one correct answer

- Often uncertain about the answer

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

In 2017: 341 students

60% Bachelor of Commerce

24% Bachelor of Science

6% Master of Science (Bioinformatics)

10% 8 other degrees/categories

Add WeChat powcoder

What are your strengths and weaknesses?

Get extra help

- Your classmates
- Discussion forum

- Consultation times
- Textbooks
- Oassis

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Homework

1. Log in to the discussion forum

2. Install RStudio on your computer

3. Start reading the [R introduction and reference guide](#) before week 2

<https://powcoder.com>

Add WeChat powcoder

Tips

The best way to learn statistics is by **solving problems** and 'getting your hands dirty' with data.

Assignment Project Exam Help

We encourage you to attend all lectures, tutorial and computer labs to get as much practice and feedback as possible.

<https://powcoder.com>

Good luck!

Add WeChat powcoder

Outline

Assignment Project Exam Help

Subject information

<https://powcoder.com>

Review of probability

Descriptive statistics

Add WeChat powcoder

Basic data visualisations

Why probability?

- It forms the mathematical foundation for statistical models and procedures
- Let's review what we know already...

<https://powcoder.com>

Add WeChat powcoder

Random variables (notation)

- Random variables (rvs) are denoted by uppercase letters: X , Y , Z , etc.
- Outcomes, or realisations, of random variables are denoted by corresponding lowercase letters: x , y , z , etc.

<https://powcoder.com>

Add WeChat powcoder

Distribution functions

- The cumulative distribution function (cdf) of X is

Assignment Project Exam Help

$$F(x) = \Pr(X \leq x), \quad -\infty < x < \infty$$

- If X is a continuous rv then it has a probability density function (pdf), $f(x)$, that satisfies

<https://powcoder.com>

$$f(x) = F'(x) = \frac{d}{dx} F(x)$$

Add WeChat powcoder

$$F(x) = \int_{-\infty}^x f(t) dt$$

- If X is a discrete rv then it has a probability mass function (pmf),

Assignment Project Exam Help

$$p(x) = \Pr(X=x), \quad x \in \Omega$$

where Ω is a discrete set, e.g. $\Omega = \{1, 2, \dots\}$.

- $\Pr(X \geq x) = 1 - F(x)$ is called a tail probability of X
- $F(x)$ increases to 1 as $x \rightarrow \infty$ and decreases to 0 as $x \rightarrow -\infty$
- If the rv has a certain distribution with pdf f (or pmf p), we write

Add WeChat powcoder
 $X \sim f$ (or $X \sim p$)

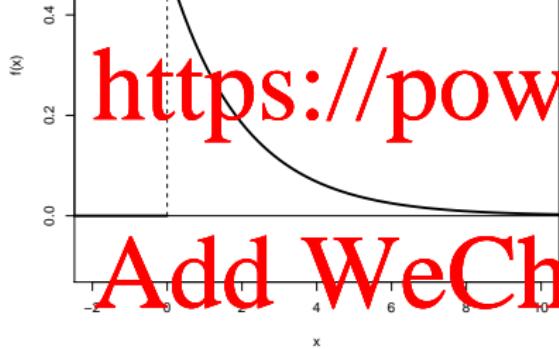
Example: Unemployment duration

A large group of individuals have recently lost their jobs. Let X denote the length of time (in months) that any particular individual will stay unemployed. It was found that this was well-described by the following pdf:

<https://powcoder.com>

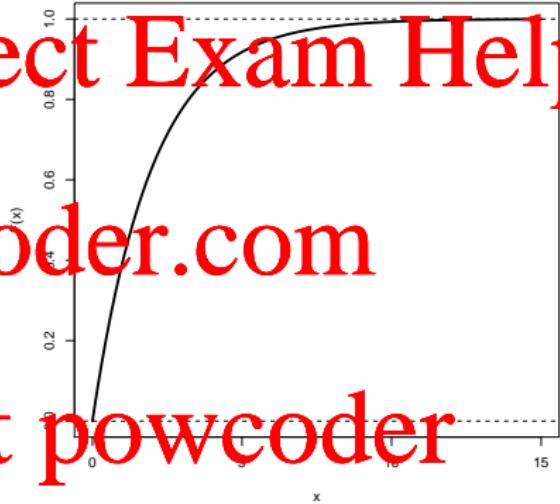
Add WeChat powcoder

Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder



Clearly, $f(x) \geq 0$ for any x and the total area under the pdf is:

$$\Pr(-\infty < X < \infty) = \int_0^{\infty} \frac{1}{2}e^{-x/2} dx = \frac{1}{2} \left[-2e^{-x/2} \right]_0^{\infty} = 1.$$

The probability that a person in the population finds a new job within 3 months is:

$$\Pr(0 \leq X \leq 3) = \int_0^3 \frac{1}{2}e^{-x/2} dx = \frac{1}{2} \left[-2e^{-x/2} \right]_0^3 = 0.7769.$$

Add WeChat powcoder

Example: Received calls

The number of calls received by an office in a given day, X , is well represented by a pmf with the following expression:

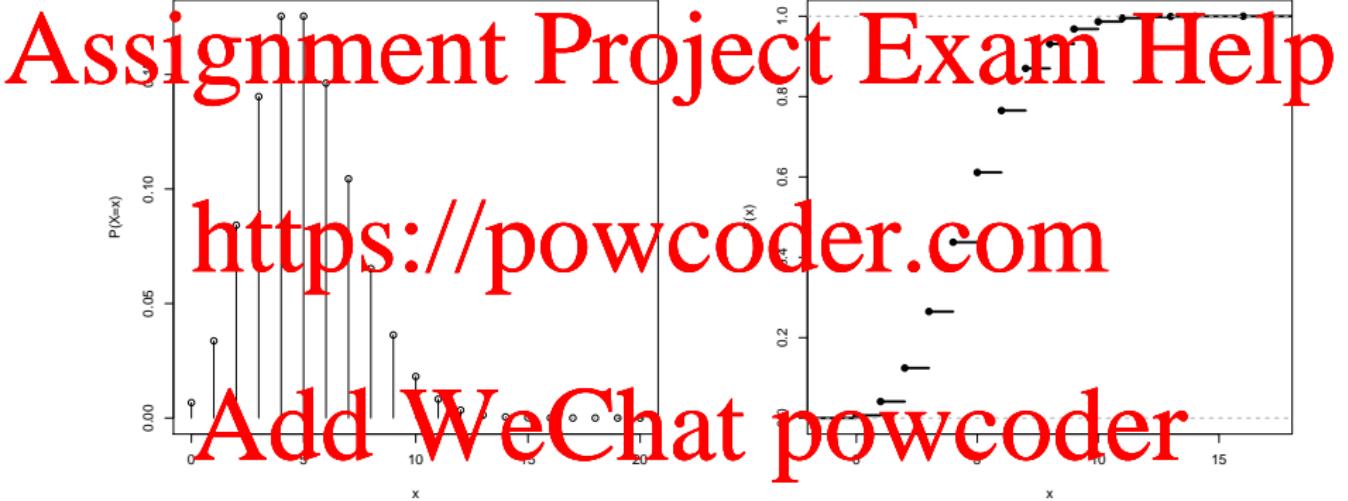
$$p(x) = \frac{e^{-5} 5^x}{x!}, \quad x \in \{0, 1, 2, \dots\},$$

where $x! = 1 \cdot 2 \cdots (x-1) \cdot x$ and $0! = 1$. For example,

$$\Pr(X = 1) = e^{-5} 5 = 0.03368$$

Add WeChat powcoder

$$\Pr(X = 3) = \frac{e^{-5} 5^3}{3 \cdot 2 \cdot 1} = 0.1403$$



To show that $p(x)$ is a pmf we need to show

Assignment Project Exam Help

$$\sum_{x=0}^{\infty} p(x) = p(0) + p(1) + p(2) + \cdots = 1.$$

Since the Taylor Series expansion of e^z is $\sum_{i=0}^{\infty} z^i / i!$, we can write

$$\sum_{x=0}^{\infty} \frac{e^{-5} 5^x}{x!} = e^{-5} \sum_{x=0}^{\infty} \frac{5^x}{x!} = e^5 e^{-5} = 1.$$

Add WeChat powcoder

Moments and variance

- The expected value (or the first moment) of a rv is denoted by

$E(X)$ and

Assignment Project Exam Help

$$E(X) = \sum_{x=-\infty}^{\infty} x p(x) \quad (\text{discrete rv})$$

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx \quad (\text{continuous rv})$$

- Higher moments, $E(X^k) = \mathbb{E}(X^k)$, for $k \geq 1$, can be obtained by

Add WeChat powcoder

$$E(X^k) = \sum_{x=-\infty}^{\infty} x^k p(x) \quad (\text{discrete rv})$$

$$E(X^k) = \int_{-\infty}^{\infty} x^k f(x) dx \quad (\text{continuous rv})$$

- More generally for a function $g(x)$ we can compute

$$\mathbb{E}(g(X)) = \sum_{x=-\infty}^{\infty} g(x)p(x) \quad (\text{discrete rv})$$

$$\mathbb{E}(g(X)) = \int_{-\infty}^{\infty} g(x)f(x) dx \quad (\text{continuous rv})$$

<https://powcoder.com>

Letting $g(x) = x^k$ gives the moments.

- The **variance** of X is defined by

Add WeChat powcoder

and the **standard deviation** of X is $\text{sd}(X) = \sqrt{\text{var}(X)}$

- “Computational” formula: $\text{var}(X) = \mathbb{E}(X^2) - \{\mathbb{E}(X)\}^2$

Basic properties of expectation and variance

- For any rv X and constant c ,

$$\mathbb{E}(cX) = c\mathbb{E}(X), \quad \text{var}(cX) = c^2 \text{var}(X)$$

- For any two rvs X and Y ,

<https://powcoder.com>

- For any two **independent** rvs X and Y ,

Add WeChat powcoder

- More generally,

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2 \text{cov}(X, Y)$$

where $\text{cov}(X, Y)$ is the **covariance** between X and Y

Covariance

- Definition of covariance:

Assignment Project Exam Help

$$\text{cov}(X, Y) = \mathbb{E}\{(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))\}$$

- Specifically, for the continuous case

<https://powcoder.com>

$$\text{cov}(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mathbb{E}(X))(y - \mathbb{E}(Y))f(x, y) dx dy$$

where $f(x, y)$ is the bivariate pdf for pair (X, Y) .

- “Computational” formula:

$$\begin{aligned}\text{cov}(X, Y) &= \mathbb{E}\{(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))\} \\ &= \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)\end{aligned}$$

Correlation

- If $\text{cov}(X, Y) > 0$ then X and Y are positively correlated
- If $\text{cov}(X, Y) < 0$ then X and Y are negatively correlated

- If $\text{cov}(X, Y) = 0$ then X and Y are uncorrelated
- The correlation between X and Y is defined as:

$$\rho = \text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\text{sd}(X) \text{sd}(Y)}, \quad -1 \leq \rho \leq 1$$

- When $\rho = +1$ then X and Y are perfectly correlated.

Add WeChat powcoder

Moment generating functions

- A moment generating function (mgf) of a rv X is

$$M_X(t) = \mathbb{E}(e^{tX}), \quad t \in (-\infty, \infty)$$

- It enables us to generate moments of X by differentiating at $t = 0$

<https://powcoder.com>

$$M_X(0) = \mathbb{E}(X)$$

$$M_X^{(k)}(0) = \mathbb{E}(X^k), \quad k \geq 1$$

- The mgf uniquely determines a distribution. Hence knowing the mgf is the same as knowing the distribution.
- If X and Y are independent rvs,

$$M_{X+Y}(t) = \mathbb{E}\{e^{t(X+Y)}\} = \mathbb{E}\{e^{tX}\} \mathbb{E}\{e^{tY}\} = M_X(t)M_Y(t)$$

i.e. the mgf of the sum is the product of individual mgfs.

Bernoulli distribution

Assignment Project Exam Help

- X takes on the values 1 (success) or 0 (failure)

- $X \sim \text{Be}(p)$ with pmf

$$p(x) = p^x(1-p)^{1-x} \quad x \in \{0, 1\}$$

- Properties:

Add WeChat powcoder

$$\mathbb{E}(X) = p$$
$$\text{var}(X) = p(1 - p)$$

$$M_X(t) = pe^t + 1 - p$$

Binomial distribution

Assignment Project Exam Help

- $X \sim \text{Bi}(n, p)$ with pmf

$$p(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x \in \{0, 1, \dots, n\}$$

- Properties:

Add WeChat powcoder

$$\mathbb{E}(X) = np$$
$$\text{var}(X) = np(1-p)$$

$$M_X(t) = (pe^t + 1 - p)^n$$

Poisson distribution

- $X \sim \text{Pn}(\lambda)$ with pmf

$$p(x) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x \in \{0, 1, \dots\}$$

- Properties:

<https://powcoder.com>

$$\mathbb{E}(X) = \text{var}(X) = \lambda$$

$$M_X(t) = e^{\lambda(e^t - 1)}$$

- Add WeChat powcoder
It arises as an approximation to $\text{Bi}(n, p)$. Letting $\lambda = np$ gives

$$p(x) = \binom{n}{x} p^x (1-p)^{n-x} \approx e^{-\lambda} \frac{\lambda^x}{x!}$$

as $n \rightarrow \infty$ and $p \rightarrow 0$.

Uniform distribution

- $X \sim \text{Unif}(a, b)$ with pdf

$$f(x) = \frac{1}{b-a}, \quad x \in (a, b)$$

- Properties:

<https://powcoder.com>

$$\mathbb{E}(X) = \frac{(a+b)}{2}$$

Add WeChat  powcoder

$$\text{Var}(X) = \frac{(b-a)^2}{12}$$

$$M_X(t) = \frac{e^{tb} - e^{ta}}{t(b-a)}$$

- If $b = 1$ and $a = 0$, this is known as the uniform distribution over the unit interval.

Exponential distribution

- $X \sim \text{Exp}(\lambda)$ with pdf

Assignment Project Exam Help

- It approximates “time until first success” for independent $\text{Be}(p)$ trials every Δt units of time with $p = \lambda \Delta t$ and $\Delta t \rightarrow 0$
- Properties:

<https://powcoder.com>

$$\mathbb{E}(X) = 1/\lambda$$

$$\text{var}(X) = 1/\lambda^2$$

Add WeChat powcoder

$$M_X(t) = \frac{\lambda}{\lambda - t}$$

- It is famous for being the only continuous distribution with the **memoryless property**:

$$\Pr(X > y + x \mid X > y) = \Pr(X > x), \quad x \geq 0, \quad y \geq 0.$$

Normal distribution

- $X \sim N(\mu, \sigma^2)$ with pdf

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in (-\infty, \infty), \quad \mu \in (-\infty, \infty), \quad \sigma > 0$$

- It is important in applications because of the Central Limit Theorem (CLT)

<https://powcoder.com>

- Properties:

Assignment Project Exam Help

$$\mathbb{E}(X) = \mu$$

$$\text{var}(X) = \sigma^2$$

<https://powcoder.com>

- When $\mu = 0$ and $\sigma = 1$ we have the standard normal distribution.
- If $X \sim N(\mu, \sigma^2)$

Add WeChat powcoder

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

Quantiles

Let X be a continuous rv. The p th quantile of its distribution is a number π_p such that $P(X \leq \pi_p) = F(\pi_p)$.

In other words, the area under $f(x)$ to the left of π_p is p :

$$p = \int_{-\infty}^{\pi_p} f(x) dx = F(\pi_p)$$

- π_p is also called the $(100p)$ th percentile
- The 50th percentile (0.5 quantile) is the median, denoted by $m = \pi_{0.5}$
- The 25th and 75th percentiles are the first and third quartiles, denoted by $q_1 = \pi_{0.25}$ and $q_3 = \pi_{0.75}$

Example: Weibull distribution

The time X until failure of a certain product has the pdf

Assignment Project Exam Help

$$f(x) = \frac{3x^2}{4} e^{-(x/4)^3}, \quad x \in (0, \infty).$$

The cdf is

$$F(x) = 1 - e^{-(x/4)^3}, \quad x \in (0, \infty)$$

Then $\pi_{0.3}$ satisfies $0.3 = F(\pi_{0.3})$. Therefore,

$$1 - e^{-(\pi_{0.3}/4)^3} = 0.3$$

$$\Rightarrow \ln(0.7) = -(\pi_{0.3}/4)^3$$

$$\Rightarrow \pi_{0.3} = -4(\ln 0.7)^{1/3} = 2.84.$$

Law of Large Numbers (LLN)

Consider a collection X_1, \dots, X_n of independent and identically distributed (i.i.d.) random variables with $\mathbb{E}(X) = \mu < \infty$, then with probability 1 we have:

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu \quad \text{as } n \rightarrow \infty$$

<https://powcoder.com>

The LLN ‘guarantees’ that long-run averages behave as we expect them to.

Add WeChat powcoder

$$\mathbb{E}(X) \approx \frac{1}{n} \sum_{i=1}^n X_i.$$

Central Limit Theorem (CLT)

Consider a collection X_1, \dots, X_n of iid rvs with $\mathbb{E}(X) = \mu < \infty$ and $\text{Var}(X) = \sigma^2 < \infty$. Let

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

<https://powcoder.com>

Then

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

Add WeChat powcoder

follows a $N(0, 1)$ distribution as $n \rightarrow \infty$.

This is an extremely important theorem!

It provides the 'magic' that will make statistical analysis work.

Example

Let X_1, \dots, X_{25} be iid rvs where $X_i \sim \text{Exp}(\lambda = 1/5)$.

Recall that $\mathbb{E}[X] = \lambda$.

Thus, the LLN implies

$\bar{X} \xrightarrow{\text{LLN}} \mathbb{E}(X) = 5.$

Moreover, since $\text{var}(X) = 1/\lambda^2 = 25$, we have

Add WeChat $\bar{X} \xrightarrow{\text{LLN}} N\left(\frac{1}{\lambda}, \frac{1}{n\lambda^2}\right) = N\left(5, \frac{5^2}{25}\right)$ powcoder

Is $n = 25$ large enough?

A simulation exercise

Generate $B = 1000$ samples of size n . For each sample compute \bar{x} .

The continuous variable is the normal $N(0.5^2/n)$ distribution prescribed by the CLT.

<https://powcoder.com>

Sample 1: $x_1^{(1)}, \dots, x_n^{(1)} \rightarrow \bar{x}^{(1)}$

Sample 2: $x_1^{(2)}, \dots, x_n^{(2)} \rightarrow \bar{x}^{(2)}$

Add WeChat powcoder

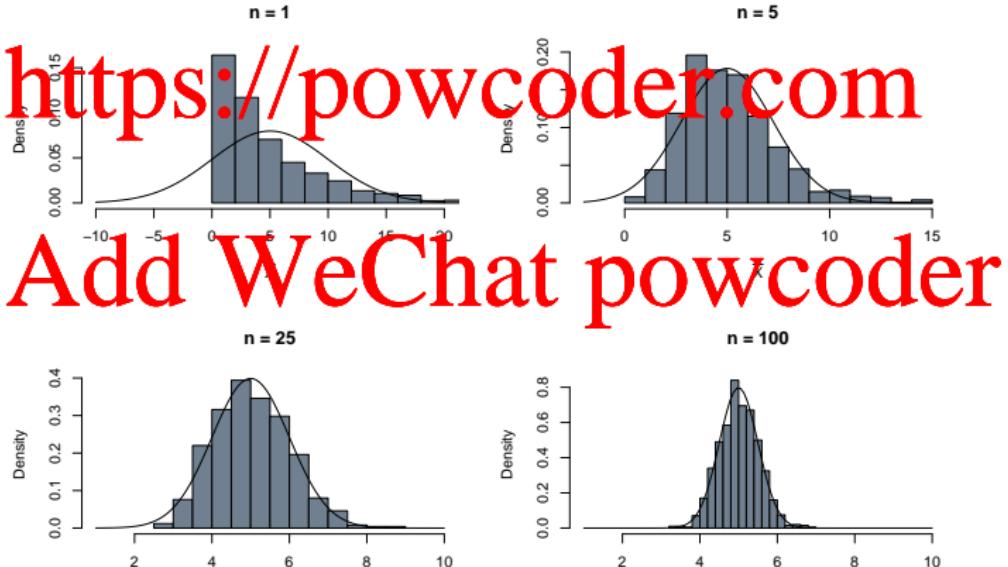
Sample B: $x_1^{(B)}, \dots, x_n^{(B)} \rightarrow \bar{x}^{(B)}$

Then represent the distribution of $\{\bar{x}^{(b)}, b = 1, \dots, B\}$ by a histogram.

A simulation exercise

The distribution of \bar{X} approaches the theoretical distribution (CLT). Moreover it will be more and more concentrated around $\mu_{\bar{X}}$ (LN). To see this, note that $\text{var}(\bar{X}) = \frac{\sigma^2}{n}, n \rightarrow \infty$ as $n \rightarrow \infty$.

Assignment Project Exam Help



Challenge problem

Let X_1, X_2, \dots, X_{25} be iid rvs with pdf $f(x) = ax^3$ where $0 < x < 2$.

Assignment Project Exam Help

1. What is the value of a ?
2. Calculate $\mathbb{E}(X_1)$ and $\text{var}(X_1)$.
3. What is an approximate value of $\Pr(X < 1.5)$?

<https://powcoder.com>

Add WeChat powcoder

Outline

Assignment Project Exam Help

Subject information

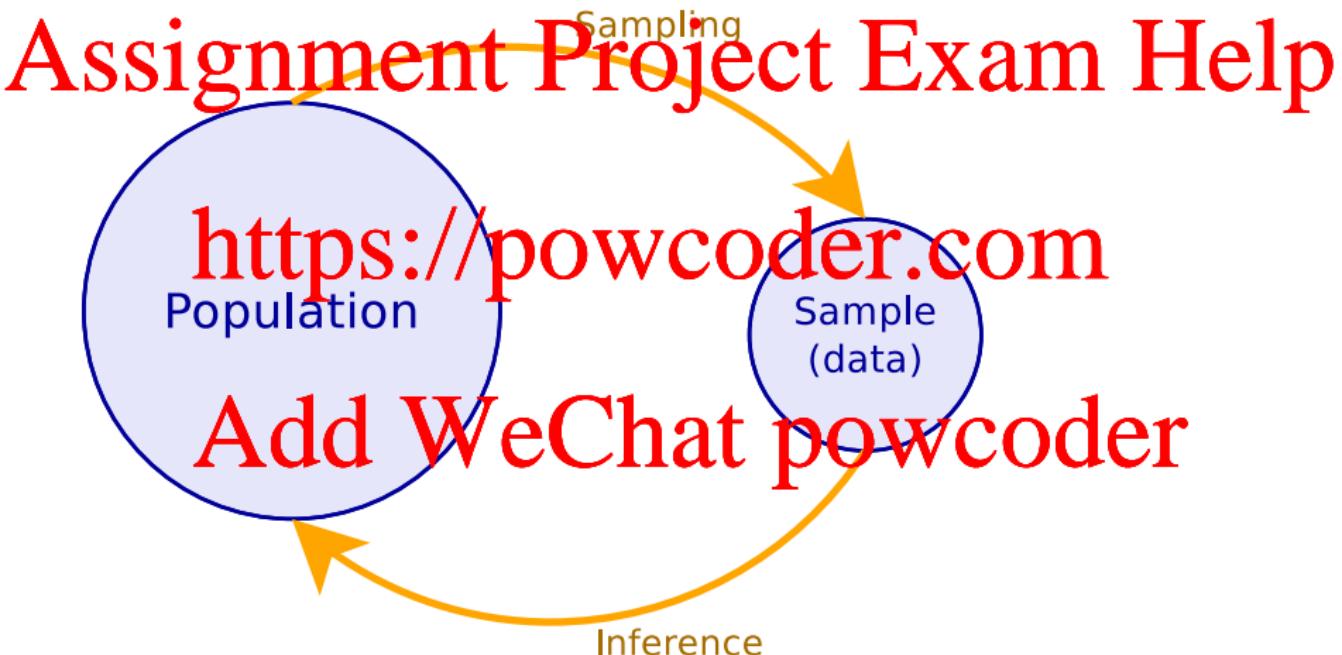
<https://powcoder.com>

Review of probability

Descriptive statistics

Add WeChat powcoder

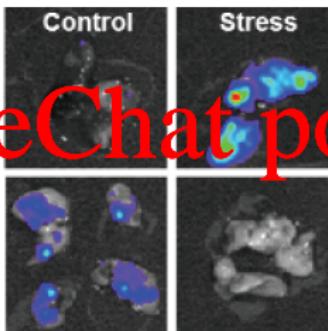
Basic data visualisations



Example: Stress and cancer

- An experiment gives independent measurements on 10 mice
- Mice are divided in control and stress groups
- The biologist considers two different proteins:
 - Vascular endothelial growth factor C (VEGFC)
 - Prostaglandin-endoperoxide synthase 2 (COX2)

<https://powcoder.com>



Add WeChat powcoder

Mouse	Group	VEGFC	COX2
1	Control	0.96718	14.05901
2	Control	0.51940	6.92926
3	Control	0.73270	11.02799
4	Control	0.96008	6.16924
5	Control	1.25964	7.32697
6	Stress	4.05745	6.45443
7	Stress	2.41335	12.95872
8	Stress	1.52595	13.26786
9	Stress	6.07073	55.03024
10	Stress	5.07591	29.92790

<https://powcoder.com>

Add WeChat powcoder

Data & sampling

- The **data** are numbers:

x_1, \dots, x_n

Assignment Project Exam Help

- The model for the data is a **random sample**, that is a sequence of iid rvs:

<https://powcoder.com>

This model is equivalent to random selection from a hypothetical infinite **population**.

- The goal is to use the data to learn about the distribution of the random variables (and, therefore, the population).

Add WeChat powcoder

Statistic

- A **statistic** $T = \phi(X_1, \dots, X_n)$ is a function of the sample and its realisation is denoted by $t = \phi(x_1, \dots, x_n)$.
- Note: the word “statistic” can also be used to refer to both the realisation, t , as well as the random variable, T . Sometime need to be more specific about which one is meant.
- A statistic has two purposes:
 - Describe or summarise the sample — **descriptive statistics**
 - Estimate the distribution generating the sample — **inferential statistics**
- A statistic can be both descriptive and inferential. It depends on how you wish to use/interpret it (see later)
- We now introduce some commonly used descriptive statistics...

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

$$\text{Sample mean} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{23.59}{10} = 2.359$$

$$\text{Sample variance} = s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = 3.98761$$

$$\text{Sample standard deviation} = s = \sqrt{3.98761} = 1.9969$$

Add WeChat powcoder

These are 'sample' or empirical versions of moments of a random variable

Empirical means 'derived from the data'

Order statistics

Arrange the sample x_1, \dots, x_n in order of increasing magnitude and define:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

Then $x_{(k)}$ is the k th order statistic.

Special cases:

- $x_{(1)}$ is the sample minimum
- $x_{(n)}$ is the sample maximum

For the example data,

$$x_{(1)} = 0.52, \quad x_{(2)} = 0.73, \quad \dots, \quad x_{(10)} = 6.07$$

What is $x_{(3.25)}$?

Let it be 0.25 of the way from $x_{(3)}$ to $x_{(4)}$,

Assignment Project Exam Help

$$\begin{aligned}x_{(3.25)} &= x_{(3)} + 0.25 \cdot (x_{(4)} - x_{(3)}) \\&= 0.96 + 0.25 \cdot (0.97 - 0.96) \\&= 0.9625\end{aligned}$$

<https://powcoder.com>

In other words, define it via **linear interpolation**.

Exercise: verify that $A_{(7.5)} = 3.6480$

Add WeChat powcoder

Why do this? It allows us to define...

Sample quantiles

General definition ('Type 7' quantiles):

Assignment Project Exam Help

$$\hat{\pi}_p = x_{(k)}, \text{ where } k = 1 + (n - 1)p$$

Special cases:

https://powcoder.com

$$\text{Sample median} = \hat{\pi}_{0.5} = x_{(5.5)} = \frac{1.26 + 1.53}{2} = 1.395$$

$$\text{Sample 1st quartile} = \hat{\pi}_{0.25} = x_{(3.25)} = 0.9625$$

$$\text{Sample 3rd quartile} = \hat{\pi}_{0.75} = x_{(7.7)} = 3.5480$$

Also:

$$\text{Interquartile range} = \hat{\pi}_{0.75} - \hat{\pi}_{0.25} = 2.685$$

$\hat{\pi}_{0.25}$ and $\hat{\pi}_{0.75}$ contain about 50% of the sample between them

Some descriptive statistics in R

Assignment Project Exam Help

```
> x  
[1] 0.97 0.52 0.73 0.96 1.26 4.06 2.41 1.53 6.07 5.08  
https://powcoder.com  
> sort(x)      # order statistics  
[1] 0.52 0.73 0.96 0.97 1.26 1.53 2.41 4.06 5.08 6.07  
Add WeChat powcoder  
> summary(x) # sample mean & sample quantiles  
   Min. 1st Qu. Median     Mean 3rd Qu.    Max.  
0.5200  0.9625  1.3950  2.3590  3.6475  6.0700
```

```
> var(x) # sample variance  
[1] 3.98761
```

Assignment Project Exam Help

```
> sd(x) # sample standard deviation  
[1] 1.9969
```

```
> IQR(x) # interquartile range  
[1] 2.685
```

<https://powcoder.com>

Add WeChat powcoder

Frequency statistics

Can also define empirical versions of pdf, pmf, cdf

Assignment Project Exam Help

Will see in the next section...

<https://powcoder.com>

Add WeChat powcoder

Outline

Assignment Project Exam Help

Subject information

<https://powcoder.com>

Descriptive statistics

Add WeChat powcoder

Basic data visualisations

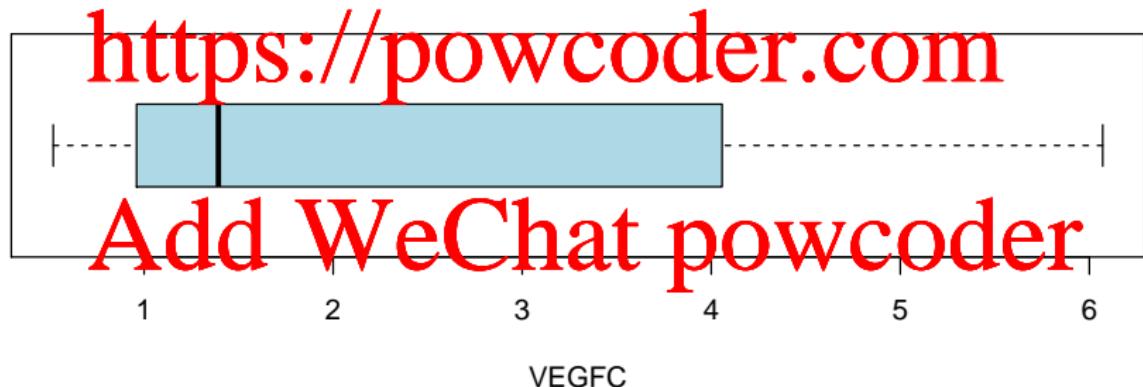
Box plot

Graphical summary of data from a single variable

Main box: $[Q_1, Q_3]$, $Q_1 = 0.7$

'Whiskers': $x_{(1)}, x_{(n)}$

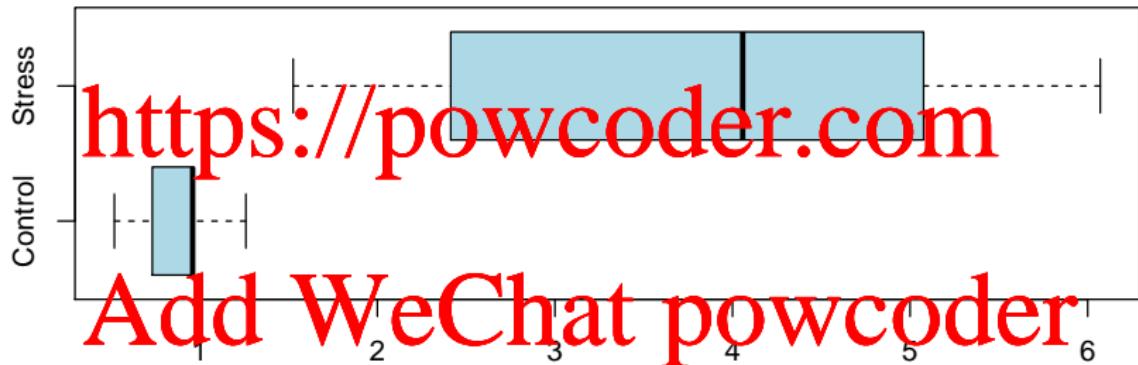
(but R does something more complicated, see tutorial problems)



Convenient way of comparing data from different groups

Example: VEGFC (Stress vs Control)

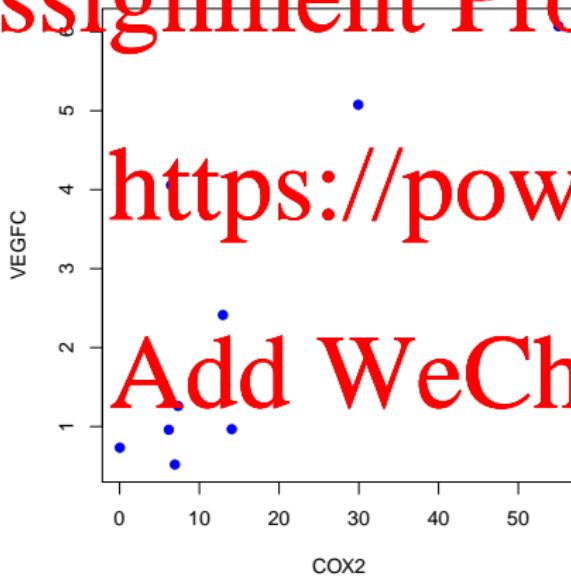
Assignment Project Exam Help



Scatter plot

For comparing data from two variables (usually continuous)

Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder

Empirical cdf

The sample cdf, or empirical cdf, is defined as

Assignment Project Exam Help

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

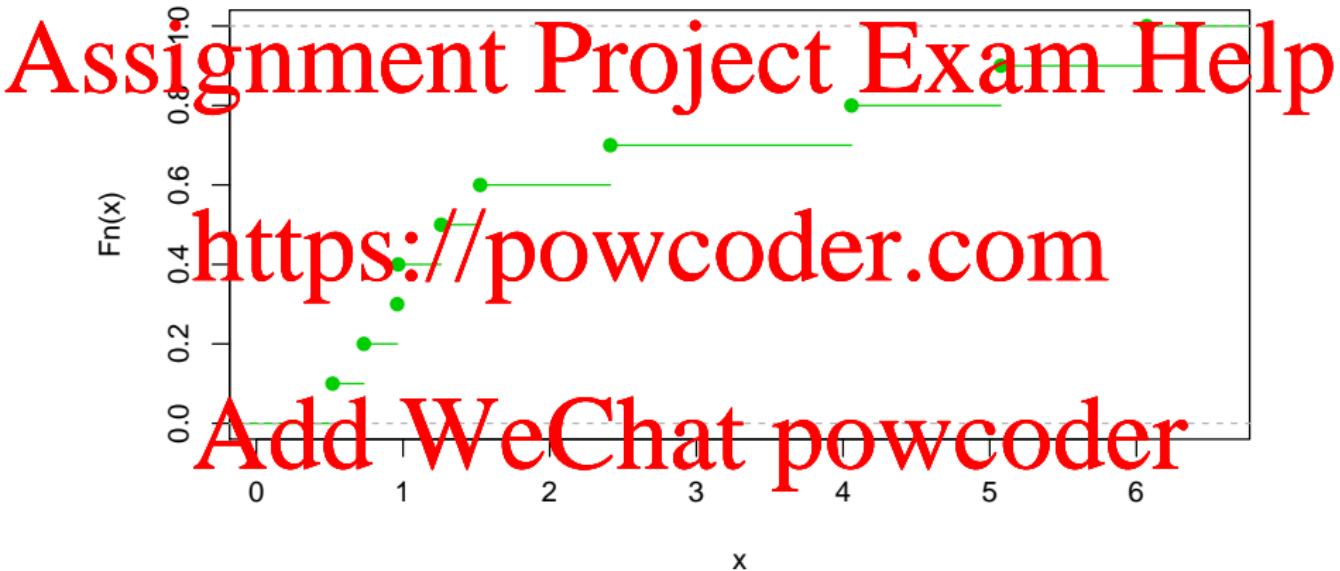
where $I(\cdot)$ is the indicator function ($I(x_i \leq x)$ has value 1 if $x_i \leq x$ and value 0 if $x_i > x$).

For example, for the previous data,

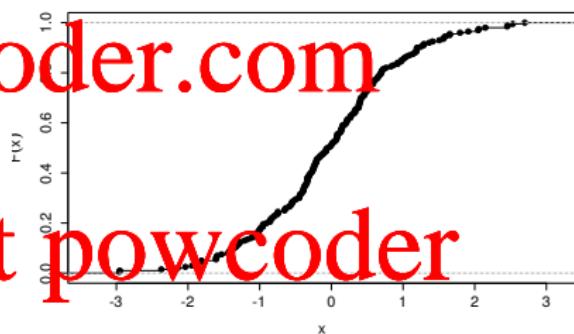
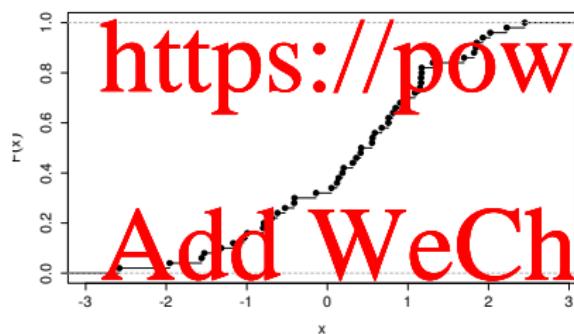
Add WeChat powcoder

$$\hat{F}(2) = \frac{1}{10} \sum_{i=1}^{10} I(x_i \leq 2) = \frac{6}{10} = 0.6$$

ecdf(VEGFC)



It has the form of a discrete cdf. However, it will approximate the cdf of a continuous variable if the sample size is large. The following diagram shows cdfs based on $n = 50$ and $n = 200$ observations sampled from a standard normal distribution, $N(0, 1)$.

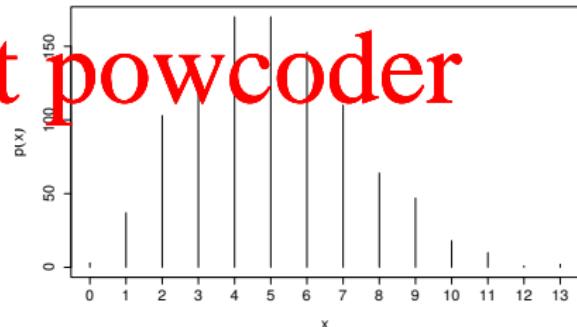
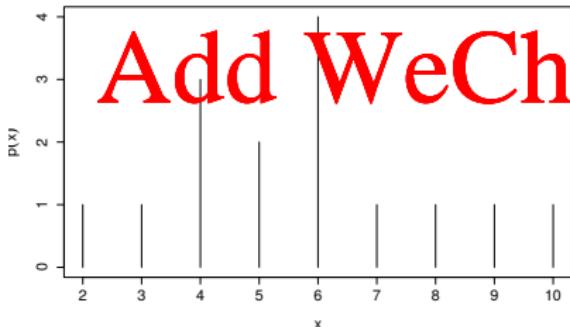


Empirical pmf

If the underlying variable is discrete we use the pmf corresponding to the sample cdf \hat{F}

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i = x)$$

For example, the following shows $\hat{p}(x)$ of size $n = 15$ from $P_n(1)$ (left) and the true pmf $p(x)$ of $P_n(5)$ (right)



Histograms and smoothed pdfs

If the underlying variable is continuous we would prefer to obtain an approximation of the pdf. There are several approaches that can be used:

<https://powcoder.com>

Add WeChat powcoder

1. **Histogram**, \hat{f}_h (h is the bin length). First divide the entire range of values into a series of small intervals (bins) and then count how many values fall into each interval. For interval $[a, b)$, where $b - a = h$, draw a rectangle with height

Assignment Project Exam Help

$$\hat{f}_h(x) = \frac{1}{hn} \sum_{i=1}^n I(a \leq x_i < b)$$

<https://powcoder.com>

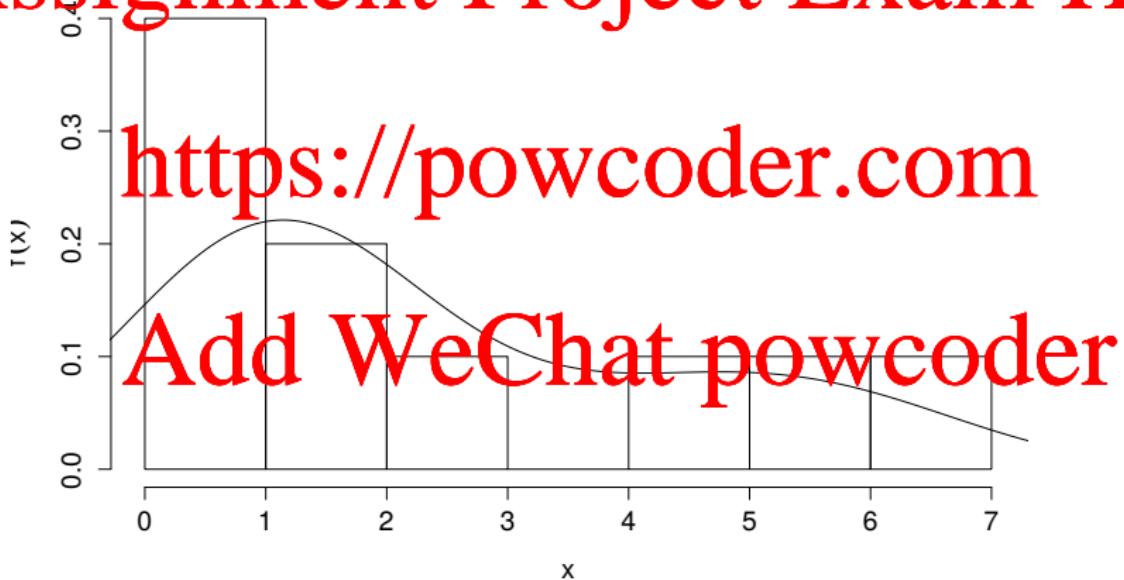
2. **Smoothed pdf**, f_h (h is the 'bandwidth parameter'),

Add WeChat powcoder

where $K(\cdot)$ is the **kernel** (a non-negative function that integrates to 1 and with mean zero) and h is a parameter that controls the level of smoothing.

Example: VEGFC

Assignment Project Exam Help



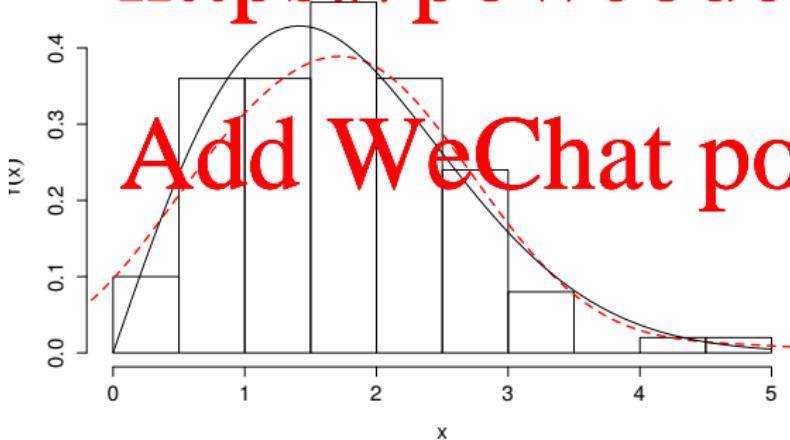
Simulated data

Consider $n = 100$ observations from the Weibull distribution with pdf

$$f(x) = \frac{1}{2}xe^{-x/2}, x > 0$$

True density (solid black curve), smoothed pdf (red dashed curve)

<https://powcoder.com>



Add WeChat powcoder

Quantile-quantile (QQ) plots

- For comparing the similarity of two probability distributions
- We plot their quantiles against each other (as a scatter plot)
- Typically, we compare data against a theoretical distribution
- The points in the plot are $(\hat{\pi}_p, \pi_p)$
- Note: some people plot these the other way around: $(\pi_p, \hat{\pi}_p)$

Add WeChat powcoder

- One axis shows the data, written here as sample quantiles:

$$\hat{\pi}_p = x_{(k)}, \quad \text{where } p = \frac{k}{n+1} \quad (\text{'Type 6' quantiles})$$

Assignment Project Exam Help

- Other axis shows corresponding quantiles for a theoretical distribution:

$\pi_p = F^{-1}(p) = F^{-1}\left(\frac{k}{n+1}\right)$

- The points in the plot therefore are,

Add WeChat powcoder
 $\left\{x_{(k)}, F^{-1}\left(\frac{k}{n+1}\right)\right\}$.

Example: VEGFC

Is the sample from an exponential distribution with cdf

$$F(x) = 1 - e^{-\lambda x}$$

Sample quantiles (data):

$$x_{(1)} = 0.52, x_{(2)} = 0.73, \dots, x_{(10)} = 6.07$$

<https://powcoder.com>

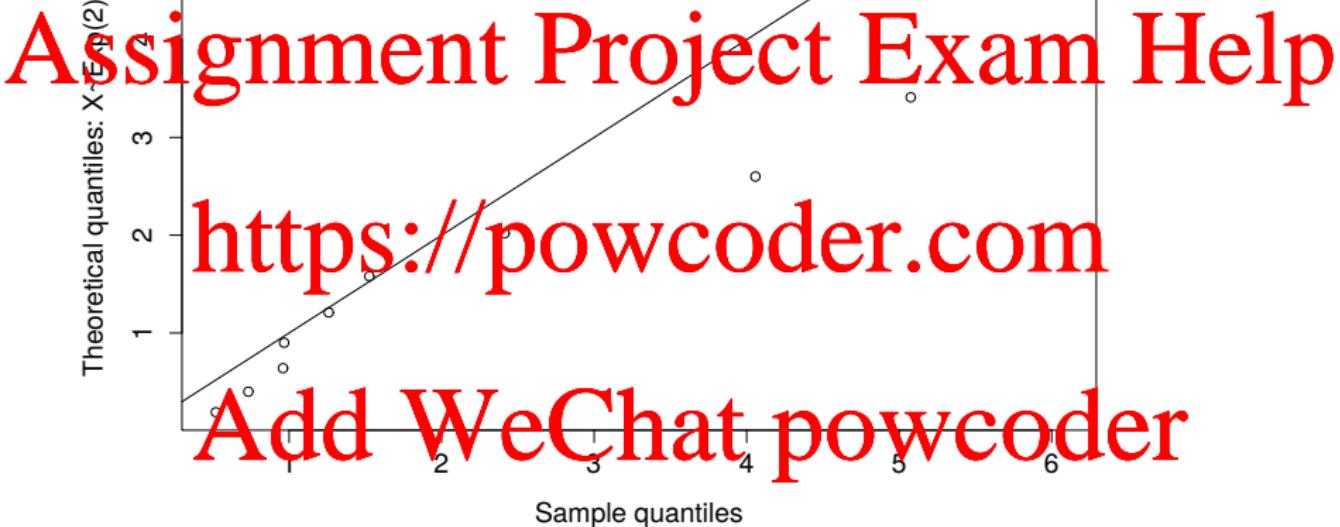
Theoretical quantiles:

$$F^{-1}(p) = -\ln(1-p)/\lambda \quad (\text{e.g. set } \lambda = 0.5)$$

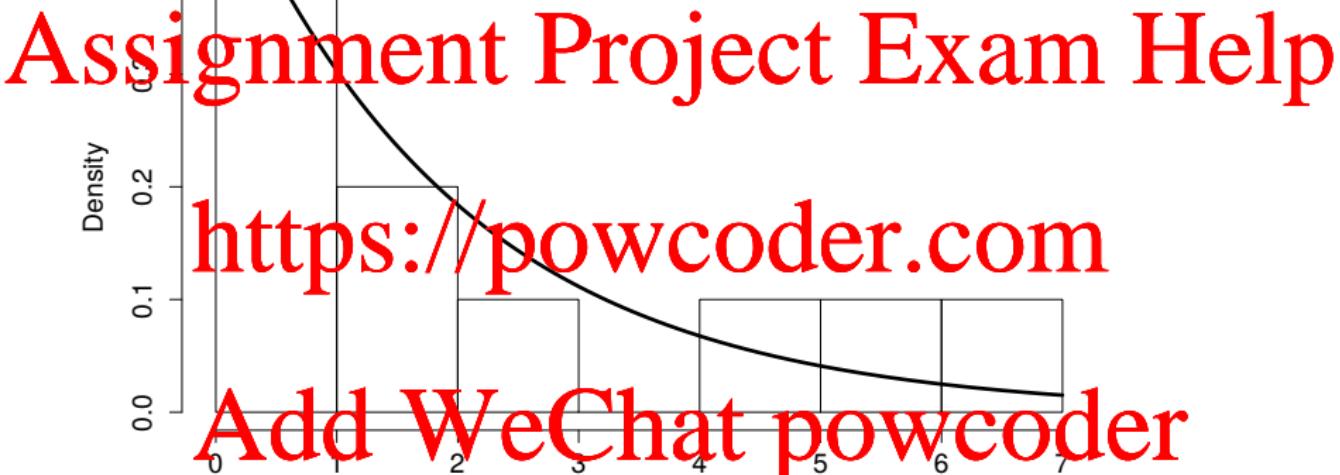
Add WeChat powcoder

$$1/(10+1) = 0.09, 2/(10+1) = 0.18, \dots, 10/(10+1) = 0.91$$

$$F^{-1}(0.09) = 0.19, F^{-1}(0.18) = 0.40, \dots, F^{-1}(0.91) = 4.80$$



The right tail of the sample does not quite match the theoretical model (tail of the sample distribution is heavier).



The right tail of the sample does not quite match the theoretical model (tail of the sample distribution is heavier).

Normal QQ plots

If $X \sim N(\mu, \sigma^2)$, then $X = \mu + \sigma Z$, where $Z \sim N(0, 1)$. Therefore, if the normal model is correct

$$x_{(k)} \approx \mu + \sigma \Phi^{-1} \left(\frac{k}{n+1} \right)$$

<https://powcoder.com>

where $\Phi(z) = P(Z \leq z)$ is the standard normal cdf.

So, if we plot the points

$$\left(\Phi^{-1} \left(\frac{k}{n+1} \right), x_{(k)} \right), \quad k = 1, \dots, n$$

the result should be a straight line with intercept μ and slope σ .
The values $\Phi^{-1}(k/(n+1))$ are called **normal scores**.

Example: simulated data

Consider 25 observations from $X \sim N(10, 2)$.

The histogram is not very helpful.

</

But the QQ plot is much clearer:

