

Lecture Book

B.Sc Mobile Robotics

D. T. McGuiness, PhD

Version: 2025.SS

(2025, D. T. McGuines, Ph.D)

Current version is 2025.SS.

This document includes the contents of Drive Systems, official name being *Mobile Robotics*, taught at MCI in the Mechatronik Design Innovation. This document is the part of the module MECH-B-4-MRV-MRO-ILV taught in the B.Sc degree.

All relevant code of the document is done using *SageMath* where stated using v10.3 and Python v3.13.1.

This document was compiled with *LuaT_EX* v1.18.0, and all editing were done using GNU Emacs v29.4 using AUCT_EX and org-mode package.

This document is based on the books and resources:

The current maintainer of this work along with the primary lecturer
is D. T. McGuines, Ph.D. (dtm@mci4me.at).

Contents

I Mechanics of Mobile Robotics	3
1 Locomotion	5
1.1 Introduction	5
1.1.1 Key Issues for Locomotion	8
1.2 Legged Mobile Robots	10
1.2.1 Examples of Legged Robot Locomotion	13
1.3 Wheeled Mobile Robots	17
1.3.1 Design	17
1.3.2 Stability	19
1.3.3 Manoeuvrability	20
1.3.4 Controllability	21
1.3.5 Case Studies for Wheeled Motion	22
1.3.6 Walking Wheels	24
2 Perception	27
2.1 Introduction	27
2.1.1 Sensors for Mobile Robotics	27
2.1.2 Sensor Classification	28
2.1.3 Characterising Sensor Performance	28
2.1.4 Wheel and Motor Sensors	33
2.2 Active Ranging	40
2.2.1 The Ultrasonic Sensor	41
2.2.2 Motion and Speed Sensors	48
2.3 Vision Based Sensors	51
2.3.1 CMOS Technology	55
2.3.2 Visual Ranging Sensors	56
2.3.3 Depth from Focus	58
2.4 Feature Extraction	61
2.4.1 Defining Feature	61
2.4.2 Using Range Data	63
3 Theory of Probability	65
3.1 Introduction	65
3.2 Experiments & Outcomes	69
3.3 Probability	70

3.4	Permutations & Combinations	75
3.4.1	Permutations	75
3.4.2	Combinations	76
3.4.3	Factorial Function	77
3.4.4	Binomial Coefficients	78
3.5	Random Variables and Probability Distributions	79
3.5.1	Discrete Random Variables and Distributions	80
3.5.2	Continuous Random Variables and Distributions	81
3.6	Mean and Variance of a Distribution	83
3.7	Binomial, Poisson, and Hyper-geometric Distributions	86
3.7.1	Sampling with Replacement	89
3.7.2	Sampling without Replacement: Hyper-geometric Distribution	89
3.8	Normal Distribution	90
3.8.1	Distribution Function	91
3.8.2	Numeric Values	91
3.8.3	Normal Approximation of the Binomial Distribution	92
3.9	Distribution of Several Random Variables	93
3.9.1	Discrete Two-Dimensional Distribution	93
3.9.2	Continuous Two-Dimensional Distribution	94
3.9.3	Marginal Distributions of a Discrete Distribution	94
3.9.4	Marginal Distributions of a Continuous Distribution	95
3.9.5	Independence of Random Variables	96
3.9.6	Functions of Random Variables	97
3.9.7	Addition of Means	97
3.9.8	Addition of Variances	98
4	Statistical Methods	101
4.1	Introduction	101
4.2	Point Estimation of Parameters	104
4.2.1	Maximum Likelihood Method	105
4.3	Confidence Intervals	108
4.4	Testing of Hypotheses and Making Decisions	115
4.4.1	Errors in Tests	117
4.5	Goodness of Fit	121
II	Localisation and Mapping	125
III	GNU/Linux Operating System	127
5	Welcome to Linux	129
5.1	Learning the Linux Command Line	129
5.1.1	A History of Command Line Interface	130

5.1.2	Linux is a Nutshell	131
5.1.3	Linux Distributions	133
6	Command Line Fundamentals	137
6.1	Introduction	137
6.2	The Structure of Commands	140
6.2.1	Some Rules Regarding the Syntax	141
6.3	Helpful Keyboard Shortcuts for the Terminal	143
6.4	When you need help with Commands	145
6.5	Additional Information	149
6.5.1	Use Tab completion on the Shell	149
6.5.2	The info command	149
6.5.3	The whatis command	150
IV	Robot Operating System	153
7	Installation	155
7.1	Introduction	155
7.2	Installing ROS Humble Hawksbill	156
7.2.1	Set locale	156
7.2.2	Setup Sources	156
7.2.3	Install ROS 2 packages	157
7.2.4	Setting up the Environment	157
V	Appendix	161
A	Tables	163
A.1	Introduction	163
A.2	Student-t Distribution	164
A.3	Chi-Square Distribution	166
Bibliography		169

List of Figures

1.1	Types of locomotion mechanisms used in biological systems [1].	6
1.2	Bipedal motion is not unique to only humans as a wide variety of animals show bipedal motion [4].	7
1.3	Specific power versus attainable speed of various locomotion mechanisms (Adapted from [1]).	7
1.4	RoboTrac, a hybrid wheel-leg vehicle for rough terrain.	8
1.6	Types of motions used by different animals.	10
1.5	Dug-beetle are a great example for legged mobile robotics [7] as not only they can manoeuvre in their environment using legged motion, they can also manipulate their environment and generate rotational motion [8].	10
1.7	Main locomotory gaits in <i>Pleurotya</i> caterpillar [15].	11
1.8	The Degrees of Freedom (DoF) a human leg has [16].	11
1.9	An example of a leg possessing three (3) DoF [17].	11
1.10	The Raibert hopper [22].	13
1.11	The 2D single Bow Leg Hopper.	13
1.12	The New ASIMO introduced in 2005 [23].	13
1.13	Spring Flamingo is a planar bipedal walking robot [28].	15
1.14	Genghis, one of the most famous walking robots from MIT uses hobby servomotors as its actuators.	16
1.15	Genghis, one of the most famous walking robots from MIT uses hobby servomotors as its actuators.	17
1.16	The four basic wheel types a)Standard wheel: Two degrees of freedom; rotation around the (motorized) wheel axle and the contact point b)castor wheel: Two degrees of freedom; rotation around an offset steering joint c)Swedish wheel: Three degrees of freedom; rotation around the (motorized) wheel axle, around the rollers and around the contact point	18
1.17	NAVLAB I, the first autonomous highway vehicle that steers and controls the throttle using vision and radar sensors [30].	20
1.18	Example of an Ackerman drive used mostly in automotive industry [31].	21
2.1	An example of a rotary encoder. [32]	34
2.2	An example of an electronic compass [33].	35
2.3	Optical Gyroscopes have no moving parts, (unlike mechanical gyroscopes) making them extremely reliable [34].	36
2.4	37
2.5	Signals of an ultrasonic sensor.	41

2.6	An example of an ultrasonic sensor used in Raspberry Pi applications [35].	42
2.8	Schematic of laser rangefinding by phase-shift measurement.	44
2.7	A laser range finder used in robotics applications	44
2.9	Range estimation by measuring the phase shift between transmitted and received signals.	45
2.10	Principle of 1D laser triangulation.	46
2.11	Structured light sources on display at the 2014 Machine Vision Show in Boston [36].	47
2.12	a) Principle of active two dimensional triangulation b) Other possible light structures c) One-dimensional schematic of the principle	47
2.13	Doppler effect between two moving objects (a) or a moving and a stationary object(b)	49
2.14	Sony ICX493AQ-A 10.14-megapixel APS-C (23.4 × 15.6 mm) Charge Coupled Device (CCD) from digital camera Sony DSLR-A200 or DSLR-A300, sensor side [37].	51
2.15	Normalized Spectral Response of a Typical Monochrome CCD.	52
2.16	Types of colour filter used in commercial and industrial applications	52
2.17	Example of white balance. Here the same scene is emulated to be shot under different light conditions [40].	54
2.18	A close-up view of a Complimentary MOS (CMOS) sensor and its circuitry [44]. . .	55
2.19	Photon noise simulation. Number of photons per pixel increases from left to right and from upper row to bottom row [47].	56
2.20	Depiction of the camera optics and its impact on the image. To get a sharp image, the image plane must coincide with the focal plane. Otherwise the image of the point (x, y, z) will be blurred in the image as can be seen in the drawing above. .	57
2.21	Three images of the same scene taken with a camera at three different focusing positions. Note the significant change in texture sharpness between the near surface and far surface [48].	58
3.1	The histogram of the data given in Exercise 1	67
3.2	A visual comparison of the Stirling formula and the actual values of the factorial function.	77
3.3	A visual representation of the Eq. (3.42).	82
3.4	The Poisson distribution with different mean (μ) values.	88
3.5	The poster child of probability and statistics, the normal distribution.	90
3.6	A visual representation between the relationship of PDF and CDF.	91
3.7	Many samples from a bivariate normal distribution. The marginal distributions are shown on the z-axis. The marginal distribution of X is also approximated by creating a histogram of the X coordinates without consideration of the Y coordinates. . .	94
4.1	The student-t distribution with different degrees of freedom m	111
4.2	Chi-square distribution with different degrees of freedom.	113
4.3	The t -distribution used in Example 4.10.	116
4.4	Illustration of Type I and II errors in testing a hypothesis $\theta = \theta_0$ against an alternative $\theta = \theta_0$	118
5.1	Bourne shell interaction on Version 7 Unix (Original).	131

5.2	The kernel mapping of the Linux operating system.	132
5.3	A Family tree of the debian branch of linux [53].	135
6.1	A graphical interface from the late 1980s, which features a TUI window for a man page, a shaped window (oclock) as well as several iconified windows. In the lower right we can see a terminal emulator running a Unix shell, in which the user can type commands as if they were sitting at a terminal. - <i>From Wikipedia</i>	140

List of Tables

4.1	Useful c values based on a given confidence (γ) value.	109
4.2	Frequency table of the sample given in the question.	123
5.1	Most popular distributions used according to distrowatch.com	133
6.1	Types of shells used in industry and academia. For reference, the authors computer uses zsh.	138
A.1	Values of z for given values of the distribution function $F(z)$ with $m = 1 - 10$	164
A.2	Values of z for given values of the distribution function $F(z)$ with $m = 11 - 20$	164
A.3	Values of z for given values of the distribution function $F(z)$ with $m = 21 - 30$	165
A.4	Values of z for given values of the distribution function $F(z)$ with $m = 1 - 10$	166
A.5	Values of z for given values of the distribution function $F(z)$ with $m = 11 - 20$	166
A.6	Values of z for given values of the distribution function $F(z)$ with $m = 21 - 30$	166

List of Examples

3.1 Recording Data	66
3.2 Leaf Plots	66
3.3 Histogram	66
3.4 Empirical Rule Outliers and z-Score	68
3.5 Sample Spaces of Random Experiments & Events	69
3.6 Fair Die	70
3.7 Coin Tossing	71
3.8 Mutually Exclusive Events	72
3.9 Union of Arbitrary Events	72
3.10 Multiplication Rule	73
3.11 Sampling w/o Replacement	74
3.12 An Encrypted Message	76
3.13 Sampling Light-bulbs	77
3.14 Waiting Time Problem	81
3.15 Continuous Distribution	82
3.16 Mean and Variance	83
3.17 Binomial Distribution	87
3.18 Poisson Distribution	88
3.19 The Parking Problem	88
3.20 Marginal Distributions of a Discrete Two-Dimensional Random Variable	95
3.21 Independence and Dependence	96
4.1 Generating Random Numbers	102
4.2 Maximum Likelihood of Gaussian Distribution	106
4.3 Maximum Likelihood of Poisson Distribution	106
4.4 For Science	107
4.5 Sampling the Population	107
4.6 Confidence Interval for mean with known variance in Normal Distribution	109
4.7 Sample Size Needed for a Confidence Interval of Prescribed Length	110
4.8 Confidence Interval for Mean of Normal Distribution with Unknown Variance	111
4.9 Confidence Interval for the Variance of the Normal Distribution	112
4.10 Test of a Hypothesis	115
4.11 Test for the Mean of the Normal Distribution with Known Variance	119
4.12 Comparison of the Means of Two Normal Distributions	120

4.13 Test of Normality	122
----------------------------------	-----

List of Theorems

3.1	First Definition of Probability	70
3.2	General Definition of Probability	71
3.3	Complementation Rule	71
3.4	Addition Rule for Mutually Exclusive Events	72
3.5	Addition Rule for Arbitrary Events	72
3.6	Multiplication Rule	73
3.7	Permutations	75
3.8	Permutations	76
3.9	Combinations	77
3.10	Random Variable	79
3.11	Mean of a Symmetric Distribution	84
3.12	Transformation of Mean and Variance	84
3.13	Relationship between PDF and CDF	91
3.14	Normal Probabilities for Intervals	91
3.15	Limit Theorem of De Moivre and Laplace	92
3.16	Addition of Means	98
3.17	Multiplication of Means	98
3.18	Addition of Variances	99
4.1	Sum of Independent Normal Random Variables	109
4.2	Student's t-Distribution	111
4.3	Chi-Square Distribution	112
4.4	Central Limit Theorem	113

Part I

Mechanics of Mobile Robotics

Chapter 1

Locomotion

Table of Contents

1.1	Introduction	5
1.2	Legged Mobile Robots	10
1.3	Wheeled Mobile Robots	17

1.1 Introduction

A mobile robot needs locomotion mechanisms which enable it to move **unbounded** throughout its environment. However, as with everything in engineering our solution comes with options, and so the selection of a robot's approach to locomotion is an important aspect of mobile robot design. In laboratory settings, there are robots that can walk, jump, run, slide, swim, fly and of course roll.

Most locomotion mechanisms have been inspired by biological counterparts, shown in **Fig. 1.1.**

There is, however, one (1) exception where there is, practically, **NO** natural equivalent:

Actively powered wheel is a human invention achieving high efficiency on flat ground.

This mechanism is **NOT** completely foreign to biological systems¹. Our bi-pedal walking system can be approximated by a rolling polygon, with sides equal in length to the span of the step. As the step size decreases, the polygon approaches a circle or wheel. But nature did not develop a fully rotating, actively powered joint, which is the technology necessary for wheeled locomotion.

Biological systems succeed in moving through a wide variety of harsh environments. Therefore it can be desirable to copy their selection of locomotion mechanisms². Replicating nature in this regard, however, is extremely difficult for several reasons.

¹While this statement is practically true, single cell organism use a similar locomotion to what we call wheel [2].

²Scientifically, this is called **Biomimetics** [3]

Type of motion	Resistance to motion	Basic kinematics of motion
Flow in a Channel	Hydrodynamic forces	Eddies
Crawl	Friction forces	Longitudinal vibration
Sliding	Friction forces	Transverse vibration
Running	Loss of kinetic energy	Oscillatory movement of a multi-link pendulum
Jumping	Loss of kinetic energy	Oscillatory movement of a multi-link pendulum
Walking	Gravitational forces	Rolling of a polygon (see figure 2.2)

Figure 1.1: Types of locomotion mechanisms used in biological systems [1].

- Mechanical complexity is easily achieved in biological systems through structural replication.

Cell division, in combination with specialisation, can readily produce a millipede with several hundred legs and several tens of thousands of individually sensed cilia³. In man-made structures, each part must be fabricated individually, and therefore, no such economies of scale exist.

- Cell is a microscopic building block that enables extreme miniaturisation. With very small size and weight, insects achieve a level of robustness that we have not been able to match with human fabrication techniques.
- The biological energy storage system and the muscular and hydraulic activation systems used in animals and insects achieve torque, response time and conversion efficiencies that far exceed similarly scaled man-made systems.

Based on these aforementioned limitations, mobile robots generally generate motion, either using wheeled mechanisms, a well-known human technology for vehicles, or using a small number of articulated legs, the simplest of the biological approaches to locomotion (shown in Fig. 1.2).

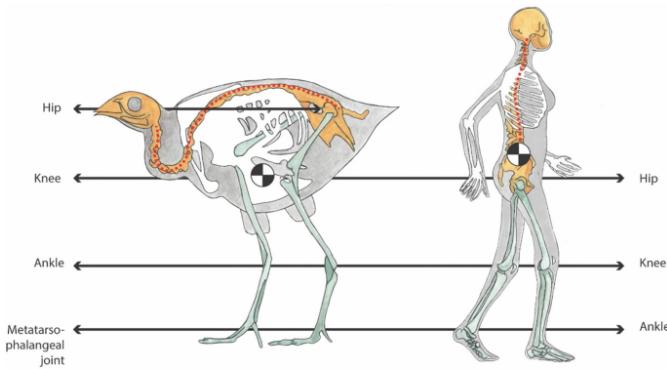


Figure 1.2: Bipedal motion is not unique to only humans as a wide variety of animals show bipedal motion [4].

In general, legged locomotion requires higher DoF and therefore greater mechanical complexity than wheeled locomotion [5]. Wheels, in addition to being simple, are extremely well suited to flat ground. As **Fig. 1.3** depicts, on flat surfaces wheeled locomotion is one to two orders of magnitude more efficient than legged locomotion.

The railway is ideally engineered for wheeled locomotion because rolling friction is minimised using a hard and flat steel surface.

But as the surface becomes soft, wheeled locomotion accumulates inefficiencies due to **rolling friction** while legged locomotion suffers much less because it consists only of **point contacts** with the ground. This is demonstrated in figure 2.3 by the dramatic loss of efficiency in the case of a tire on soft ground.

the efficiency of wheeled locomotion depends greatly on environmental qualities, particularly the flatness and hardness of the ground, while the efficiency of legged locomotion depends on the leg mass and body mass, both of which the robot must support at various points in a legged gait.

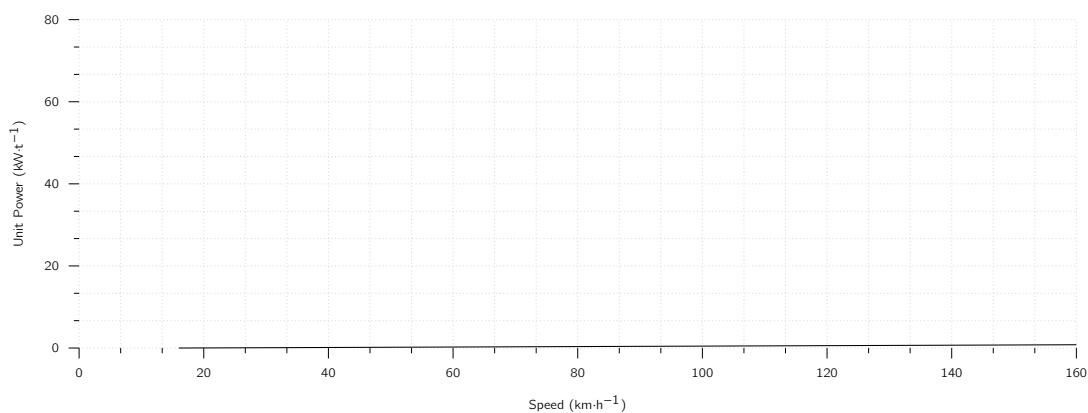


Figure 1.3: Specific power versus attainable speed of various locomotion mechanisms (Adapted from [1]).

It is understandable therefore nature favours legged locomotion, as locomotion systems in nature must operate on rough and unstructured terrain. For example, in the case of insects in a forest the vertical variation in ground height is often an order of magnitude greater than the total height of the insect.

By the same token, the human environment frequently consists of engineered, smooth surfaces both indoors and outdoors. Therefore, it is also understandable that virtually all industrial applications of mobile robotics utilise some form of wheeled locomotion. Recently, for more natural outdoor environments, there has been some progress toward hybrid and legged industrial robots such as the forestry robot [6] shown in **Fig. 1.4**.

In the next section, we present general considerations that concern all forms of mobile robot locomotion. Following this will be overviews of legged locomotion and wheeled locomotion techniques for mobile robots.

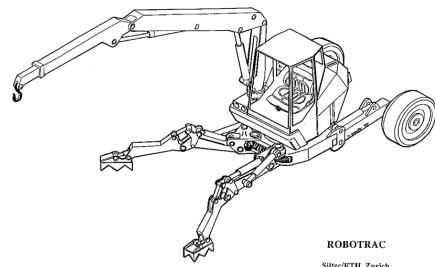


Figure 1.4: RoboTrac, a hybrid wheel-leg vehicle for rough terrain.

1.1.1 Key Issues for Locomotion

Locomotion is the **complement of manipulation**:

- In manipulation, the robot arm is fixed but moves objects in the workspace by imparting force to them.
- In locomotion, the environment is fixed and the robot moves by imparting force to the environment.

For both cases, the scientific basis is the **study of actuators** which generate interaction forces, and mechanisms that implement desired kinematic and dynamic properties. Locomotion and manipulation therefore share the same core issues of stability, contact characteristics and environmental type:

- | | |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"> ■ Stability <ul style="list-style-type: none"> – number and geometry of contact points – centre of gravity – static/dynamic stability – inclination of terrain
 ■ characteristics of contact | <ul style="list-style-type: none"> – contact point/path size and shape – angle of contact – friction
 ■ type of environment <ul style="list-style-type: none"> – structure – medium (e.g., water, air, soft or hard ground) |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

A theoretical analysis of locomotion begins with mechanics and physics. From this starting point, we can formally define and analyse all manner of mobile robot locomotion systems. However, this Lecture Book puts more emphasis on the mobile robot navigation problem, particularly on the topics of perception, localisation and cognition. Therefore, we will not delve deeply into the physical basis of locomotion. Nevertheless, two remaining sections in this chapter present overviews of issues in legged locomotion and wheeled locomotion.



(a) Bipedal motion [9].



(b) Quadpedal motion [10].



(c) Hexapedal motion [11]

Figure 1.6: Types of motions used by different animals.

1.2 Legged Mobile Robots

Legged locomotion is characterised by a **series of point contacts between the robot and the ground**. The primary advantages include adaptability and manoeuvrability in rough terrain. Because only a set of point contacts is required, the quality of the ground between those points does not matter, so long as the robot can maintain adequate ground clearance. In addition, a walking robot is capable of crossing a hole or chasm so long as its reach exceeds the width of the hole. A final advantage of legged locomotion is the potential to manipulate objects in the environment with great skill.

The dung beetle, is capable of rolling a ball while locomotion as a result of its dexterous front legs shown in **Fig. 1.5**.

The main disadvantages of legged locomotion include **power and mechanical complexity**. The leg, which may include several DoF, must be capable of sustaining part of the robot's total weight, and in many robots must be capable of lifting and lowering the robot. Additionally, high manoeuvrability will only be achieved if the legs have a sufficient number of DoF to impart forces in a number of different directions.



Figure 1.5: Dung-beetle are a great example for legged mobile robotics [7] as not only they can manoeuvre in their environment using legged motion, they can also manipulate their environment and generate rotational motion [8].

Leg Configurations and Stability

Because legged robots are biologically inspired, it is instructive to examine biologically successful legged systems. A number of different leg configurations have been successful in a variety of organisms seen in **Fig. 1.6**.

Large animals such as mammals and reptiles have four (4) legs whereas insects have six (6) or more legs. In some mammals, the ability to walk on only two (2) legs has been perfected. Especially in the case of humans, balance has progressed to the point that we can even jump with one leg⁴. This

⁴In child development, one of the tests used to determine if the child is acquiring advanced locomotion skills is the ability to jump on one leg [12].

exceptional manoeuvrability comes at a price:

Bipedal motion is much more complex active control to maintain balance.

In contrast, a creature with three (3) legs can exhibit a static, stable pose provided that it can ensure that its centre of gravity is within the tripod of ground contact. Static stability, demonstrated by a three-legged stool, means that balance is maintained with no need for motion. A small deviation from stability⁵ is passively corrected towards the stable pose when the upsetting force stops. But a robot must be able to lift its legs in order to walk. To achieve static walking, a robot **must** have at least six (6) legs [13]. In such a configuration, it is possible to design a gait⁶ in which a statically stable tripod of legs is in contact with the ground at all times.

Insects⁷ are immediately able to walk when born. For them, the problem of balance during walking is relatively simple. Mammals, with four (4) legs, cannot achieve static walking, but are able to stand easily on four (4) legs. Fauns⁸, for example, spend several minutes attempting to stand before they are able to do so, then spend several more minutes learning to walk without falling [14]. Humans, with two (2) legs, cannot even stand in one place with static stability. Infants require months to stand and walk, and even longer to learn to jump, run and stand on one leg.

There is also the potential for great variety in the complexity of each individual leg. Once again, the biological world provides ample examples at both extremes. For instance, in the case of the caterpillar, each leg is extended using hydraulic pressure by constricting the body cavity and forcing an increase in pressure, and each leg is retracted longitudinally by relaxing the hydraulic pressure, then activating a single tensile muscle that pulls the leg in towards the body, seen in **Fig. 1.7**. Each leg has only a single DoF, which is oriented longitudinally along the leg.

Forward locomotion depends on the hydraulic pressure in the body, which extends the distance between pairs of legs. The caterpillar leg is therefore mechanically very simple, using a minimal number of extrinsic muscles to achieve complex overall locomotion.

At the other extreme, the human leg has more than six (6) major degrees of freedom, combined with further actuation at the toes, shown in **Fig. 1.8**. There are more than 50 muscles in each lower limb and at least half of them participate actively in the control of leg

⁵e.g., such as gently pushing the stool

⁶the pattern of steps of an animal at a particular speed.

⁷such as spiders, ant, beetles, ...

⁸a baby deer

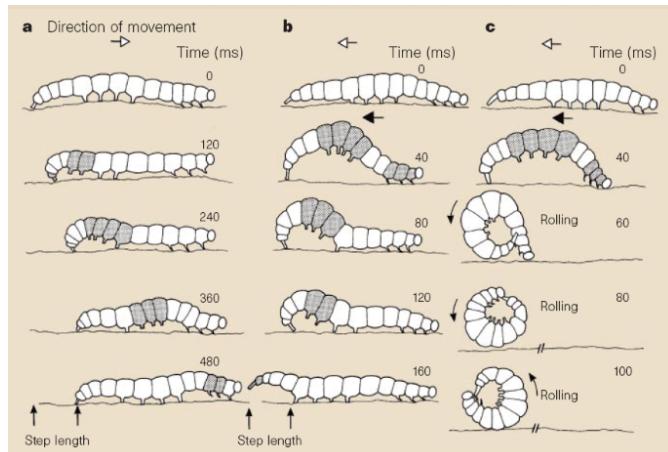
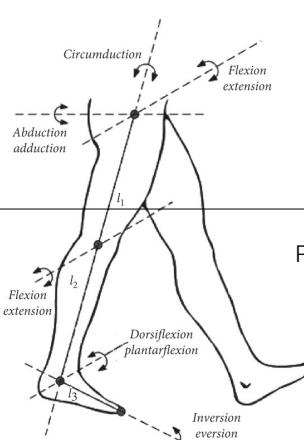


Figure 1.7: Main locomotory gaits in *Pleurotya* caterpillar [15].



motion [18, 19].

In the case of legged mobile robots, a minimum of two (2) DoF is generally required to move a leg forward by lifting the left and swinging it forward. More common is the addition of a 3rd DoF for more complex manoeuvres, resulting in legs such as those shown in **Fig. 1.9**. Recent successes in the creation of bipedal walking robots have added a fourth DOF at the ankle joint [20]. The ankle enables more consistent ground contact by actuating the pose of the sole of the foot. In general, adding DoF to a robot leg increases the manoeuvrability of the robot, both augmenting the range of terrains on which it can travel and the ability of the robot to travel with a variety of gaits. The primary disadvantages of additional joints and actuators is, of course, energy, control and mass. Additional actuators require energy and control, and they also add to leg mass, further increasing power and load requirements on existing actuators.

In the case of a multi-legged mobile robot, there is the issue of leg coordination for locomotion, or gait control.

The number of possible gaits depends on the number of legs [21].

The gait is a sequence of lift and release events for the individual legs. For a mobile robot with k legs, the total number of possible events N for a walking machine is:

$$N = (2k - 1)! \quad (1.1)$$

For a bipedal walker ($k=2$) legs the number of possible events N is:

$$N = (2k - 1)! = 3! = 3 \cdot 2 \cdot 1 = 6 \quad (1.2)$$

The six (6) different events are:

- lift right leg
- lift left leg
- release right leg
- release left leg
- lift both legs together
- release both legs together

As can we see, this list of possible events quickly grows quite large. For example, a robot with six

(6) legs has far more gaits theoretically:

$$N = 11! = 39\,916\,800 \quad (1.3)$$

1.2.1 Examples of Legged Robot Locomotion

Although there are no high-volume industrial applications to date, legged locomotion is an important area of long-term research. Several interesting designs are presented below, beginning with the one-legged robot and finishing with six-legged robots.

Single Leg

The minimum number of legs a legged robot can have is, of course, one. Minimising the number of legs is beneficial for several reasons.

- Body mass is particularly important to walking machines, and the single leg minimises cumulative leg mass.
- Leg coordination is required when a robot has several legs, but with one leg no such coordination is needed.
- The one-legged robot maximises the basic advantage of legged locomotion: legs have single points of contact with the ground in lieu of an entire track as with wheels.

A single legged robot requires only a sequence of single contacts, making it useful in rough terrain.

Perhaps most importantly, a hopping robot can dynamically cross a gap that is larger than its stride by taking a running start, whereas a multi-legged walking robot that cannot run is limited to crossing gaps that are as large as its reach.

The major challenge of creating a single-leg robot is **balance**. For a robot with one leg, static walking is not only impossible, but static stability when stationary is also impossible. The robot must actively balance itself by either changing its centre of gravity or by imparting corrective forces. Thus, the successful single-leg robot **must be dynamically stable**.

Fig. 1.10 shows the Raibert Hopper [24, 25], one of the most well-known single-leg hopping robots created. This robot makes continuous corrections to body attitude and to robot velocity by adjusting the leg angle with respect

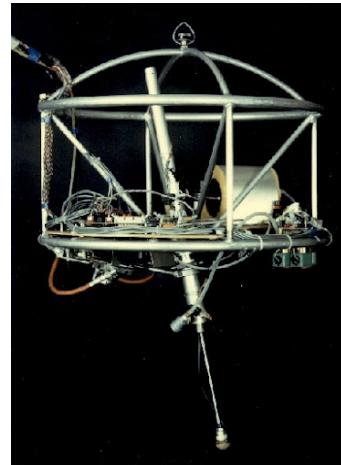


Figure 1.10: The Raibert hopper [22].

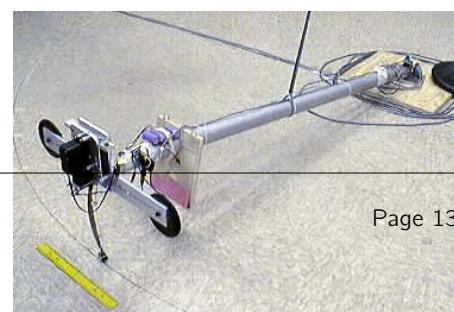


Figure 1.11: The 2D single Bow Leg Hopper.

to the body. The actuation is hydraulic, including high-power longitudinal extension of the leg during stance to hop back into the air. Although powerful, these actuators require a large, off-board hydraulic pump to be connected to the robot at all times. **Fig.** 1.11 shows a more energy efficient design developed [26]. Instead of supplying power by means of an off-board hydraulic pump, the Bow Leg Hopper is designed to capture the kinetic energy of the robot as it lands using an efficient bow spring leg. This spring returns approximately 85% of the energy, meaning that stable hopping requires only the addition of 15% of the required energy on each hop. This robot, which is constrained along one axis by a boom, has demonstrated continuous hopping for 20 minutes using a single set of batteries carried on board the robot. As with the Raibert Hopper, the Bow Leg Hopper controls velocity by changing the angle of the leg to the body at the hip joint. The paper of Ringrose [27] demonstrates the very important duality of mechanics and controls as applied to a single leg hopping machine. Often clever mechanical design can perform the same operations as complex active control circuitry. In this robot, the physical shape of the foot is exactly the right curve so that when the robot lands without being perfectly vertical, the proper corrective force is provided from the impact, making the robot vertical by the next landing. This robot is dynamically stable, and is furthermore passive.

The correction is provided by physical interactions between the robot and its environment, with no computer nor any active control in the loop.

Two Legs (Bipedal)

A variety of successful bipedal robots have been demonstrated. Two-legged robots have been shown to run, jump, travel up and down stairs and even do aerial tricks such as somersaults. **Fig.** 1.12 shows the Honda P2 bipedal robot, which is the product of tens of millions of research dollars and more than a decade of work. This biped can walk on slopes, climb and descend stairs, and push shopping carts. The crucial technology that enables this robot is Honda's research into the fabrication of extremely high torque, low mass motors that serve as the robot's joints. In the case of P2, the most significant obstacle that remains is energy capacity, efficiency and autonomous navigation. This robot can operate for only about 20 minutes with on-board power. An important feature of bipedal robots is their anthropomorphic shape. They can be built to have the same approximate dimensions

as humans, and this makes them excellent vehicles for research in human-robot interaction.

Bipedal robots can only be statically stable within some limits, and so robots such as P2 and Wabian generally must perform continuous balance-correcting servoing even when standing still. Furthermore, each leg must have sufficient capacity to support the full weight of the robot. In the case of four-legged robots, the balance problem is facilitated along with the load requirements of each leg. An elegant design of a biped robot is the Spring Flamingo of MIT seen in **Fig. 1.13**. This robot inserts springs in series with the leg actuators to achieve a more elastic gait. Combined with "kneecaps" that limit knee joint angles, the Flamingo achieves surprisingly biomimetic motion.

Four Legs (Quadruped)

Although standing still on four legs is passively stable, walking remains challenging because to remain stable the robot's center of gravity must be actively shifted during the gait. Sony recently invested several million dollars to develop a four-legged robot (figure 2.14). To create this robot, Sony created both a new robot operating system that is near real-time and invented new geared servomotors that are sufficiently high torque to support the robot, yet backdriveable for safety. In addition to developing custom motors and software, Sony incorporated a color vision system that enables Aibo to chase a brightly colored ball. The robot is able to function for at most one hour before requiring recharging. Early sales of the robot have been very strong, with more than 60,000 units sold in the first year. Nevertheless, the number of motors and the technology investment behind this robot dog have resulted in a very high price of approximately 1500. Four legged robots have the potential to serve as effective artifacts for research in human-robot interaction (fig. 2.15). Humans can treat the Sony robot, for example, as a pet and might develop an emotional relationship similar to that between man and dog. Furthermore, Sony has designed Aibo's walking style and general behavior to emulate learning and maturation, resulting in dynamic behavior over time that is more interesting for the owner who can track the changing behavior. As the challenges of high energy storage and motor technology are solved, it is likely that quadruped robots much more capable than Aibo will become common throughout the human environment.

Six Legs (Hexapod)

Six legged configurations have been extremely popular in mobile robotics because of their static stability during walking, thus reducing the control complexity (figure 2.16 and 2.17). In most cases,

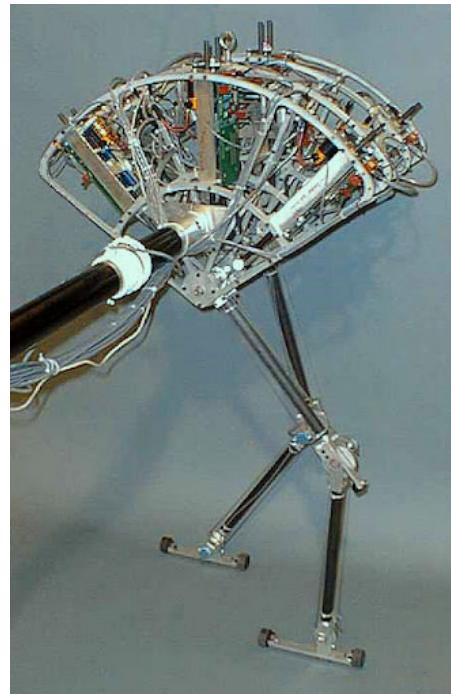


Figure 1.13: Spring Flamingo is a planar bipedal walking robot [28].

each leg has 3 DOF, including hip flexion, knee flexion and hip abduction (figure 2.6).

Genghis is a commercially available hobby robot that has six legs, each of which has 2 DOF provided by hobby servos (figure 2.18). Such a robot, which consists only of hip flexion and hip abduction, has less maneuverability in rough terrain but performs quite well on flat ground. Because it consists of a straightforward arrangement of servo motors and straight legs, such robots can be readily built by a robot hobbyist. Insects, which are arguably the most successful locomoting creatures on earth, excel at traversing all forms of terrain with six legs, even upside down. Currently, the gap between the capabilities of six-legged insects and artificial six-legged robots is still quite large. Interestingly, this is not due to a lack of sufficient numbers of degrees of freedom on the robots. Rather, insects combine a small number of active degrees of freedom with passive structures, such as microscopic barbs and textured pads, that increase the gripping strength of each leg significantly. Robotic research into such passive tip structures has only recently begun. For example, a research group is attempting to recreate the complete mechanical function of the cockroach leg (Roland, reference in notes (Espenschied et al.)). It is clear from all of the above examples that legged robots have much progress to make before they are competitive with the 24 Autonomous Mobile Robots have been realised recently, primarily due to advances in motor design. Creating actuation systems that approach the efficiency of animal muscle remains far from the reach of robotics, as does energy storage with the energy densities found in organic life forms



Figure 1.14: Genghis, one of the most famous walking robots from MIT uses hobby servomotors as its actuators.

1.3 Wheeled Mobile Robots

The wheel has been by far the most popular locomotion mechanism in mobile robotics and in man-made vehicles in general.⁹ It can achieve high efficiencies, as demonstrated in figure 2.3, and does so with a relatively simple mechanical implementation. In addition, balance is not usually a research problem in wheeled robot designs, because wheeled robots are almost always designed so that all wheels are in ground contact at all times.

⁹This should be clear as wheel motion is one of the most efficient method of converting energy to motion.

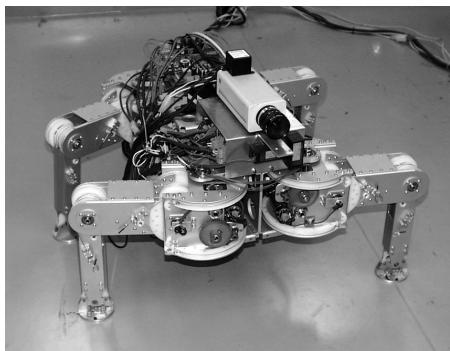


Figure 1.15: Genghis, one of the most famous walking robots from MIT uses hobby servomotors as its actuators.

Therefore, **three wheels are sufficient to guarantee stable balance**, although as we will see below, two-wheeled robots can also be stable.¹⁰ When more than three wheels are used, a suspension system is required in order to allow all wheels to maintain ground contact when the robot encounters uneven terrain. Instead of worrying about balance, researchers in wheeled robots tends to focus on the problems of **traction** and **stability**, maneuverability and control:

can the robot wheels provide sufficient traction and stability for the robot to cover all of the desired terrain, and does the robot's wheeled configuration enable sufficient control over the velocity of the robot?

¹⁰Of course, with clever implementation.

1.3.1 Design

As we will see, there is a very large space of possible wheel configurations when we consider possible techniques for mobile robot locomotion. We will begin by discussing the wheel in detail, as there are a number of different wheel types with specific strengths and weaknesses. Then, we will examine complete wheel configurations that deliver particular forms of locomotion for a mobile robot.

Wheel Design

There are four (4) major wheel classes, as shown in **Fig. 1.16**. They differ widely in their kinematics, and therefore the choice of wheel type has a large effect on the overall kinematics of the mobile robot.

The standard wheel and the castor wheel have a **primary axis of rotation** and therefore are highly directional. To move in a different direction, the wheel must be steered first along a vertical axis. The key difference between these two (2) wheels is that the standard wheel can accomplish this steering motion with no side effects, as the centre of rotation passes through the contact patch

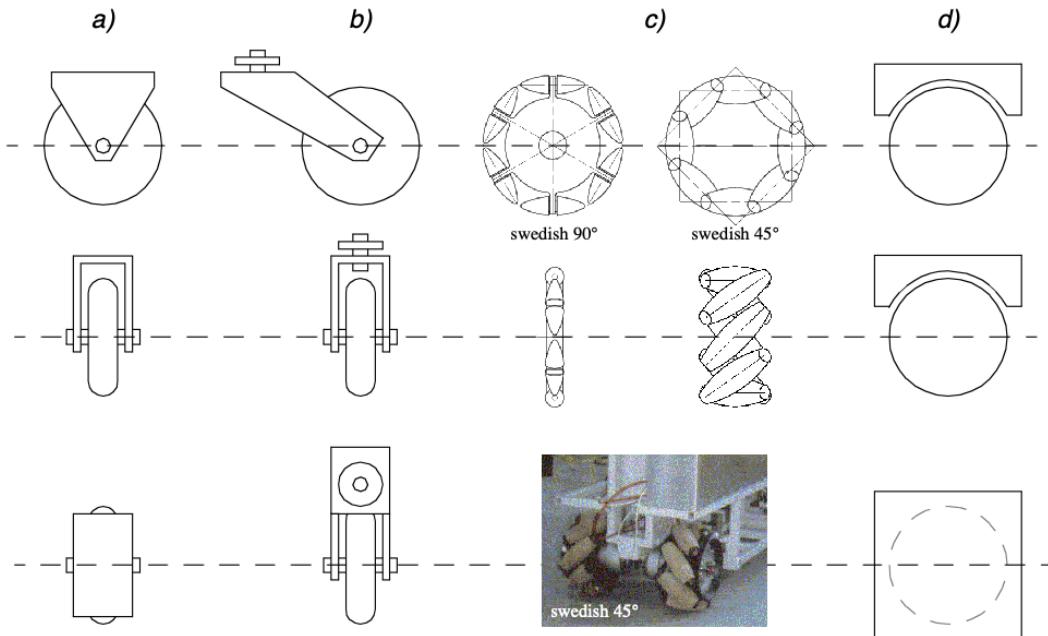


Figure 1.16: The four basic wheel types a)Standard wheel: Two degrees of freedom; rotation around the (motorized) wheel axle and the contact point b)castor wheel: Two degrees of freedom; rotation around an offset steering joint c)Swedish wheel: Three degrees of freedom; rotation around the (motorized) wheel axle, around the rollers and around the contact point

with the ground, while the castor wheel rotates around an offset axis, causing a force to be imparted to the robot chassis during steering.

¹¹Sometimes known as Swedish wheel or Ilon wheel after its inventor Bengt Erland Ilon [29].

The mecanum wheel¹¹ and the spherical wheel are both designs that are less constrained by directionality than the conventional standard wheel. The swedish wheel functions as a normal wheel, but provides low resistance in another direction as well, sometimes perpendicular to the conventional direction as in the swedish 90 and sometimes at an intermediate angle as in the swedish 45. The small rollers attached around the circumference of the wheel are passive and the wheel's primary axis serves as the only actively powered joint. The key advantage of this design is that, although the wheel rotation is powered only along the one principal axis (through the axle), the wheel can kinematically move with very little friction along many possible trajectories, not just forward and backward.

The spherical wheel is a truly omnidirectional wheel, often designed so that it may be actively powered to spin along any direction. One mechanism for implementing this spherical design imitates the computer mouse, providing actively powered rollers that rest against the top surface of the sphere and impart rotational force.

Regardless of what wheel is used, in robots designed for all-terrain environments and in robots with more than three (3) wheels, a suspension system is normally required to maintain wheel contact with the ground. One of the simplest approaches to suspension is to design flexibility into the wheel itself. For instance, in the case of some four-wheeled indoor robots that use castor wheels, manufacturers have applied a deformable tire of soft rubber to the wheel in order to create a

primitive suspension. Of course, this limited solution cannot compete with a sophisticated suspension system in applications where the robot needs a more dynamic suspension for significantly non-flat terrain.

Wheel Geometry

The choice of wheel types for a mobile robot is strongly linked to the choice of wheel arrangement, or wheel geometry. When designing a mobile robot locomotion we must consider these two (2) issues simultaneously. Why does wheel type and wheel geometry matter? Three fundamental characteristics of a robot are governed by these choices:

- maneuverability,
- controllability
- stability.

Unlike automobiles, which are largely designed for a highly standardised environment¹², mobile robots are designed for applications in a wide variety of situations. Automobiles all share similar wheel configurations as there is one region in the design space that maximises maneuverability, controllability and stability for their standard environment:

the paved road.

However, there is no single wheel configuration that maximises these qualities for the variety of environments faced by different mobile robots. So, we will see great variety in the wheel configurations of mobile robots. In fact, few robots use the Ackerman wheel configuration of the automobile because of its poor maneuverability, with the exception of mobile robots designed for the road system (figure 2.20).

Table 2.1 gives an overview of wheel configurations ordered by the number of wheels. This table shows both the selection of particular wheel types and their geometric configuration on the robot chassis. Note that some of the configurations shown are of little use in mobile robot applications. For instance, the 2-wheeled bicycle arrangement has moderate maneuverability and poor controllability. Like a single-leg hopping machine, it can never stand still. Nevertheless, this table provides an indication of the large variety of wheel configurations that are possible in mobile robot design.

1.3.2 Stability

Surprisingly, the minimum number of wheels required for static stability is two (2). As shown above, a two-wheel differential drive robot can achieve static stability if the center of mass is below the wheel axle. Cye is a commercial mobile robot that uses this wheel configuration

¹²such as the road network



Figure 1.17: NAVLAB I, the first autonomous highway vehicle that steers and controls the throttle using vision and radar sensors [30].

However, under ordinary circumstances such a solution requires wheel diameters that are impractically large. Dynamics can also cause a two-wheeled robot to strike the floor with a third point of contact, for instance with sufficiently high motor torques from standstill. Conventionally, static stability requires a minimum of three (3) wheels, with the additional caveat that the center of gravity must be contained within the triangle formed by the ground contact points of the wheels. Stability can be further improved by adding more wheels, although once the number of contact points exceeds three, the hyperstatic nature of the geometry will require some form of flexible suspension on uneven terrain.

1.3.3 Manoeuvrability

Some robots are omnidirectional, meaning that they can move at any time in any direction along the ground plane (X, Y) regardless of the orientation of the robot around its vertical axis. This level of maneuverability requires wheels that can move in more than just a single direction, and so omnidirectional robots usually employ swedish or spherical wheels that are powered. A good example is Uranus, shown in figure 2.24. This robot uses four swedish wheels to rotate and translate independently and without constraints. In general, the ground clearance of robots with swedish and spherical wheels is somewhat limited, due to the mechanical constraints of constructing omnidirectional wheels. An interesting recent solution to the problem of omnidirectional navigation while solving this ground clearance problem is the four castor-wheeled configuration in which each castor wheel is actively steered and actively translated. In this configuration, the robot is truly omnidirectional because, even if the castor wheels are facing a direction perpendicular to the desired direction of travel, the robot can still move in the desired direction by steering these wheels. Because the vertical axis is offset from the ground contact path, the result of this steering motion is robot motion.

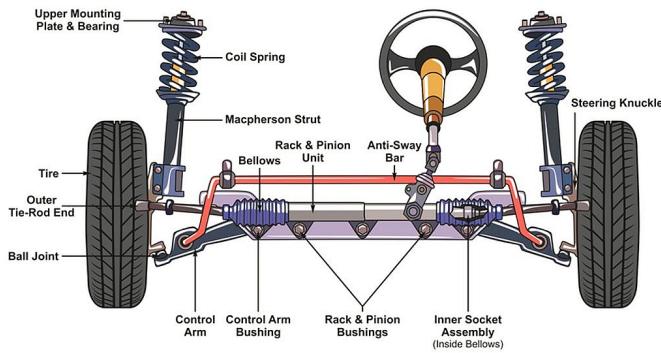


Figure 1.18: Example of an Ackerman drive used mostly in automotive industry [31].

In the research community, another classes of mobile robots are popular which achieve high maneuverability, only slightly inferior to that of the omnidirectional configurations. In such robots, motion in a particular direction may initially require a rotational motion. With a circular chassis and an axis of rotation at the center of the robot, such a robot can spin without changing its ground footprint. The most popular such robot is the two-wheel differential drive robot where the two wheels rotate around the center point of the robot. One or two additional ground contact points may be used for stability, based on the application specifics.

In contrast to the above configurations, consider the Ackerman steering configuration common in automobiles. Such a vehicle typically has a turning diameter that is larger than the car. Furthermore, for such a vehicle to move sideways requires a parking maneuver consisting of repeated changes in direction forward and backward. Nevertheless, Ackerman steering geometries have been especially popular in the hobby robotics market, where a robot can be built by starting with a remote-control race car kit and adding sensing and autonomy to the existing mechanism. In addition, the limited maneuverability of Ackerman steering has an important advantage: its directionality and steering geometry provide it with very good lateral stability in high speed turns.

1.3.4 Controllability

There is generally an **inverse correlation** between controllability and maneuverability. For example, the omni-directional designs such as the four castor-wheeled configuration require significant processing to convert desired rotational and translational velocities to individual wheel commands. Furthermore, such omni-directional designs often have greater degrees of freedom at the wheel. For instance, the swedish wheel has a set of free rollers along the wheel perimeter. These degrees of freedom cause an **accumulation of slippage**, tend to reduce dead-reckoning accuracy and increase the design complexity.

Controlling an omnidirectional robot for a specific direction of travel is also more difficult and often less accurate when compared to less manoeuvrable designs.

For example, an Ackerman steering vehicle can go straight simply by locking the steerable wheels and driving the drive wheels, which can be seen in **Fig.** 1.18.

In a differential drive vehicle, the two (2) motors attached to the two wheels must be driven along exactly the same velocity profile, which can be challenging considering variations between wheels, motors and environmental differences. With four-wheel omnidrive, such as the Uranus robot which has four swedish wheels, the problem is even harder because all four wheels must be driven at exactly the same speed for the robot to travel in a perfectly straight line.

In summary, there is **NO** “ideal” drive configuration that simultaneously maximises stability, manoeuvrability and controllability. Each mobile robot application places unique constraints on the robot design problem, and the designer’s task is to choose the most appropriate drive configuration possible from among this space of compromises.

1.3.5 Case Studies for Wheeled Motion

Now let’s describe four (4) specific wheel configurations, in order to demonstrate concrete applications of the concepts above to mobile robots built for real-world activities.

Synchro Drive

The synchro drive configuration (figure 2.22) is a popular arrangement of wheels in indoor mobile robot applications. It is an interesting configuration because, although there are three driven and steered wheels, only two motors are used in total. The one translation motor sets the speed of all three wheels together, and the one steering motor spins all the wheels together about each of their individual vertical steering axes. But note that the wheels are being steered with respect to the robot chassis, and therefore there is no direct way of re-orienting the robot chassis. In fact, the chassis orientation does drift over time due to uneven tire slippage, causing rotational dead-reckoning error.

Synchro drive is particularly advantageous in cases where omnidirectionality is needed. So long as each vertical steering axis is aligned with the contact path of each tire, the robot can always re-orient its wheels and move along a new trajectory without changing its footprint. Of course, if the robot chassis has directionality and the designers intend to re-orient the chassis purposefully, then synchro drive is only appropriate when combined with an independently rotating turret that attaches to the wheel chassis. Commercial research robots such as the Nomadics 150 or the RWI B21r have been sold with this configuration (figure 1.12). In terms of dead-reckoning, synchro drive systems are generally superior to true omni-directional configurations but inferior to differential drive and Ackerman steering systems. There are two main reasons for this. First and foremost, the translation motor generally drives the three wheels using a single belt. Due to slop and backlash in the drivetrain, whenever the drive motor engages, the closest wheel begins spinning before the furthest wheel, causing a small change in the orientation of the chassis. With additional changes in motor speed, these small angular shifts accumulate to create a large error in orientation during dead-reckoning.

Second, the mobile robot has no direct control over the orientation of the chassis. Depending on the orientation of the chassis, the wheel thrust can be highly asymmetric, with two wheels on one side and the third wheel alone, or symmetric, with one wheel on each side and one wheel straight ahead or behind, as shown in (2.22). The asymmetric cases results in a variety of errors when tire-ground slippage can occur, again causing errors in dead-reckoning of robot orientation.

Omnidirectional Drive

As we will see later in chapter 3.4.2, omnidirectional movement is of great interest for complete maneuverability. Omnidirectional robots that are able to move in any direction () at any time are also holonomic (see chapter 3.4.2). They can be realized by either using spheric, castor or swedish wheels. Three examples of such holonomic robots are presented below.

Omnidirectional locomotion with three spheric wheels

The omnidirectional robot depicted in figure 2.23 is based on three spheric wheels, each actuated by one motor. In this design, the spheric wheels are suspended by three contact points, two given by spherical bearings and one by the a wheel connected to the motor axle. This concept provides excellent maneuverability and is simple in design. However, it is limited to flat surfaces and small loads, and it is quite difficult to find round wheels with high friction coefficients

Omnidirectional locomotion with four swedish wheels

The omnidirectional arrangement depicted in figure 2.24 has been used successfully on several research robots, including the CMU Uranus. This configuration consists of four swedish 45 degree wheels, each driven by a separate motor. By varying the direction of rotation and relative speeds of the four wheels, the robot can be moved along any trajectory in the plane

and, even more impressively, can simultaneously spin around its vertical axis. For example, when all four wheels spin "forward" or "backward", the robot as a whole moves in a straight line forward and backward, respectively. However, when one diagonal pair of wheels is spun in the same direction and the other diagonal pair is spun in the opposite direction, the robot moves laterally. This four-wheel arrangement of swedish wheels is not minimal in terms of control motors. Because there are only 3 degrees of freedom in the plane, one can build a three-wheeled omnidirectional robot chassis using three swedish 90 degree wheels as shown in Table 2.1. However, existing examples such as Uranus have been designed with four wheels due to capacity and stability considerations. One application for which such omnidirectional designs are particular amenable is mobile manipulation. In this case, it is desirable to reduce the degrees of freedom of the manipulator arm to save arm mass by using the mobile robot chassis motion for gross motion. As with humans, it would be ideal if the base could move omnidirectionally without greatly impact-

Omnidirectional locomotion with four castor wheels and eight motors

Another solution for omnidirectionality is to use castor wheels. This is done for the Nomad XR4000 from Nomadics (fig. 2.25) giving it an excellent maneuverability. Unfortunately Nomadics Technology has ceased the production of mobile robots. The above two examples are drawn from Table 2.1, but this is not an exhaustive list of all wheeled locomotion techniques. Hybrid approaches that combine legged and wheeled locomotion, or tracked and wheeled locomotion, can also offer particular advantages. Below are two unique designs created for specialized applications.

Tracked Slip/Skid Locomotion

In the wheel configurations discussed above, we have made the assumption that wheels are not allowed to skid against the surface. An alternative form of steering, termed slip/skid, may be used to re-orient the robot by spinning wheels that are facing the same direction at different speeds or in opposite directions. The army tank operates this way, and Nanokhod, pictured below (figure 2.26) is an example of a mobile robot based on the same concept. Robots that make use of tread have much larger ground contact patches, and this can significantly improve their maneuverability in loose terrain compared to conventional wheeled designs. However, due to this large ground contact patch, changing the orientation of the robot usually requires a skidding turn, wherein a large portion of the track must slide against the terrain. The disadvantage of such configurations is coupled to the slip/skid steering. Because of the large amount of skidding during a turn, the exact center of rotation of the robot is hard to predict and the exact change in position and orientation is also subject to variations depending on the ground friction. Therefore, dead-reckoning on such robots is highly inaccurate. This is the trade-off that is made in return for extremely good maneuverability and traction over rough and loose terrain. Furthermore, a slip/skid approach on a high-friction surface can quickly overcome the torque capabilities of the motors being used. In terms of power efficiency, this approach is reasonably efficient on loose terrain but extremely inefficient otherwise.

1.3.6 Walking Wheels

Walking robots might offer the best maneuverability in rough terrain. However, they are inefficient on flat ground and need sophisticated control. Hybrid solutions, combining the adaptability of legs with the efficiency of wheels offer an interesting compromise. Solutions that passively adapt to the terrain are of particular interest for field and space robotics. The Sojourner robot of NASA/JPL (fig. 1.2) represents such a hybrid solution, able to overcome objects up to the size of the wheels. A more advanced mobile robot design for similar applications has recently been produced by EPFL (fig. 2.27). This robot, called Shrimp2, has 6 motorized wheels and is capable of climbing objects up to two times its wheel diameter [84,85]. This enables it to climb regular stairs though the robot is even smaller than the Sojourner. Using a rhombus configuration, the Shrimp has a steering wheel in the front and the rear, and two wheels arranged on a bogie on each side. The front wheel has a spring suspension to guarantee optimal ground contact of all wheels at any time. The

steering of the rover is realized by synchronizing the steering of the front and rear wheels and the speed difference of the bogie wheels. This allows for high precision maneuvers and turning on the spot with minimum slip/skid of the four center wheels. The use of parallel articulations for the front wheel and the bogies creates a virtual center of rotation at the level of the wheel axis. This ensures maximum stability and climbing abilities even for very low friction coefficients between the wheel and the ground. As mobile robotics research matures we find ourselves able to design more intricate mechanical systems. At the same time, the control problems of inverse kinematics and dynam2 Locomotion 37 R. Siegwart, EPFL, Illah Nourbakhsh, CMU ics are now so readily conquered that these complex mechanics can in general be controlled. So, in the near future, you should expect to see a great number of unique, hybrid mobile robots that draw together advantages from several of the underlying locomotion mechanisms that we have discussed in this chapter. They will each be technologically impressive, and each will be designed as the expert robot for its particular environmental niche.

Chapter 2

Perception

Table of Contents

2.1	Introduction	27
2.2	Active Ranging	40
2.3	Vision Based Sensors	51
2.4	Feature Extraction	61

2.1 Introduction

One of the most important tasks of an Autonomous Mobile Robotics (AMR) is to acquire knowledge about its environment.¹ This is achieved by taking measurements using various sensors and then extracting meaningful information from those measurements.

In this chapter we present the most common sensors used in AMR and then discuss strategies for extracting information from the sensors.

¹One could even argue it is the definition of life, if you ask a biologist as the ability to feel and act on its environment is the bare necessity.

2.1.1 Sensors for Mobile Robotics

There is a wide variety of sensors used in AMRs (Fig. 4.1). Some are used to measure simple values like the internal temperature of a robot's electronics or the rotational speed of the motors in its wheels or actuators. Other, more sophisticated sensors can be used to acquire information about the robot's environment or even to directly measure a robot's global position. Here, we focus primarily on sensors used to extract information about the robot's environment. Because a AMR moves around, it will frequently encounter unforeseen environmental characteristics, and therefore such sensing is particularly critical. We begin with a functional classification of sensors. Then, after presenting basic tools for describing a sensor's performance, we proceed to describe selected sensors

in detail.

2.1.2 Sensor Classification

We classify sensors using two (2) important functional axes. Let's define these terms for clarity;

Proprioceptive sensors which measure values **internal** to the robot.

e.g., motor speed, wheel load, robot arm joint angles, battery voltage.

Exteroceptive sensors which measure information from the **robot's environment**;

e.g., distance measurements, light intensity, sound amplitude.

exteroceptive sensor measurements are interpreted by the robot to extract meaningful environmental features.

Passive sensors measure ambient environmental energy entering the sensor.

e.g., temperature probes, microphones and CCD or CMOS cameras.

Active sensors emit energy into the environment, then measure the environmental reaction. Because active sensors can manage more controlled interactions with the environment, they often achieve superior performance. However, active sensing introduces several risks: the outbound energy may affect the very characteristics that the sensor is attempting to measure. Furthermore, an active sensor may suffer from interference between its signal and those beyond its control. For example, signals emitted by other nearby robots, or similar sensors on the same robot may influence the resulting measurements. Examples of active sensors include wheel quadrature encoders, ultrasonic sensors and laser rangefinders.

The sensor classes in Table (4.1) are arranged in ascending order of complexity and descending order of technological maturity. Tactile sensors and proprioceptive sensors are critical to virtually all mobile robots, and are well understood and easily implemented. Commercial quadrature encoders, for example, may be purchased as part of a gear-motor assembly used in a AMR. At the other extreme, visual interpretation by means of one or more CCD/CMOS cameras provides a broad array of potential functionalities, from obstacle avoidance and localisation to human face recognition. However, commercially available sensor units that provide visual functionalities are only now beginning to emerge

2.1.3 Characterising Sensor Performance

The sensors we describe in this chapter vary greatly in their performance characteristics. Some sensors provide extreme accuracy in well-controlled laboratory settings, but are overcome with error

when subjected to real-world environmental variations. Other sensors provide narrow, high precision data in a wide variety settings. To quantify such performance characteristics, first we formally define the sensor performance terminology that will be valuable throughout the rest of this chapter.

Basic Sensor Response Ratings

A number of sensor characteristics can be rated **quantitatively** in a laboratory setting. Such performance ratings will necessarily be best-case scenarios when the sensor is placed on a real-world robot, but are nevertheless useful.

Dynamic Range Used to measure the spread between the lower and upper limits of inputs values to the sensor while maintaining normal sensor operation. Formally, the dynamic range is the ratio of the maximum input value to the minimum measurable input value. Because this raw ratio can be unwieldy, it is usually measured in Decibels, which is computed as ten times the common logarithm of the dynamic range. However, there is potential confusion in the calculation of Decibels, which are meant to measure the ratio between powers, such as Watts or Horsepower.

Suppose your sensor measures motor current and can register values from a minimum of 1 mA to 20 A. The dynamic range of this current sensor is defined as:

$$10 \cdot \log \left[\frac{20}{0.001} \right] = 43 \text{ dB} \quad (2.1)$$

Now suppose you have a voltage sensor that measures the voltage of your robot's battery, measuring any value from 1 mV to 20 V. Voltage is **NOT** a unit of power, but the square of voltage is proportional to power. Therefore, we use 20 instead of 10:

$$20 \cdot \log \left[\frac{20}{0.001} \right] = 86 \text{ dB} \quad (2.2)$$

Range An important rating in AMR because often robot sensors operate in environments where they are frequently exposed to input values beyond their working range. In such cases, it is critical to understand how the sensor will respond. For example, an optical rangefinder will have a minimum operating range and can thus provide spurious data when measurements are taken with object closer than that minimum.

Resolution The minimum difference between two (2) values that can be detected by a sensor. Usually, the lower limit of the dynamic range of a sensor is equal to its resolution. However, in the case of digital sensors, this is not necessarily so. For example, suppose that you have a sensor that measures voltage, performs an analogue-to-digital conversion and outputs the converted value as an 8-bit number linearly corresponding to between 0 and 5 Volts. If this sensor is truly linear, then it has $2^8 - 1$ total output values or a resolution of:

$$\frac{5}{255} = 20 \text{ mV}$$

Linearity is an important measure governing the behaviour of the sensor's output signal as the input signal varies. A linear response indicates that if two (2) inputs, say x and y result in the two outputs $f(x)$ and $f(y)$, then for any values a and b , the following relation can be derived:

$$f(x + y) = f(x) + f(y).$$

This means that a plot of the sensor's input/output response is simply a straight line.

Bandwidth or Frequency is used to measure the speed with which a sensor can provide a stream of readings. Formally, the number of measurements per second is defined as the sensor's frequency in Hz. Because of the dynamics of moving through their environment, mobile robots often are limited in maximum speed by the bandwidth of their obstacle detection sensors. Thus increasing the bandwidth of ranging and vision-based sensors has been a high-priority goal in the robotics community.

In Situ Sensor Performance

The above sensor characteristics can be reasonably measured in a laboratory environment, with confident extrapolation to performance in real-world deployment. However, a number of important measures cannot be reliably acquired without deep understanding of the complex interaction between all environmental characteristics and the sensors in question. This is most relevant to the most sophisticated sensors, including active ranging sensors and visual interpretation sensors.

Sensitivity A measure of the degree to which an incremental change in the target input signal changes the output signal. Formally, sensitivity is the ratio of output change to input change. Unfortunately, however, the sensitivity of exteroceptive sensors is often confounded by undesirable sensitivity and performance coupling to other environmental parameters.

Cross-Sensitivity is the technical term for sensitivity to environmental parameters that are orthogonal to the target parameters for the sensor. For example, a flux-gate compass can demonstrate high sensitivity to magnetic north and is therefore of use for AMR navigation. However, the compass will also demonstrate high sensitivity to ferrous building materials, so much so that its cross-sensitivity often makes the sensor useless in some indoor environments. High cross-sensitivity of a sensor is generally undesirable, especially so when it cannot be modelled.

Error of a sensor is defined as the difference between the sensor's output measurements and the true values being measured, within some specific operating context.

As an example, given a true value v and a measured value m , we can define error as:

$$\text{Error} = m - v.$$

Accuracy defined as the degree of conformity between the sensor's measurement and the true value, and is often expressed as a proportion of the true value (e.g. 97.5% accuracy):

$$\text{Accuracy} = 1 - \frac{|m - v|}{v}.$$

Of course, obtaining the ground truth (v), can be difficult or impossible, and so establishing a confident characterisation of sensor accuracy can be problematic. Further, it is important to distinguish between two different sources of error:

- Systematic errors are caused by factors or processes that can in theory be modelled. These errors are, therefore, deterministic.²

Poor calibration of a laser rangefinder, un-modelled slope of a hallway floor and a bent stereo camera head due to an earlier collision are all possible causes of systematic sensor errors

²Meaning, its value is not determined by a random process and therefore should, in theory, be predictable.

- Random errors cannot be predicted using a sophisticated model nor can they be mitigated with more precise sensor machinery. These errors can only be described in probabilistic terms (i.e. stochastic). Hue instability in a colour camera, spurious range-finding errors and black level noise in a camera are all examples of random errors.

Precision is often confused with accuracy, and now we have the tools to clearly distinguish these two terms. Intuitively, high precision relates to reproducibility of the sensor results. For example, one sensor taking multiple readings of the same environmental state has high precision if it produces the same output. In another example, multiple copies of this sensor taking readings of the same environmental state have high precision if their outputs agree. Precision does not, however, have any bearing on the accuracy of the sensor's output with respect to the true value being measured. Suppose that the random error of a sensor is characterised by some mean value (μ) and a standard deviation (σ). The formal definition of precision is the ratio of the sensor's output range to the standard deviation:

$$\text{Precision} = \frac{\text{Range}}{\sigma}.$$

Only σ and **NOT** μ has impact on precision. In contrast mean error is directly proportional to overall sensor error and inversely proportional to sensor accuracy.

Characterising Error

Mobile robots depend heavily on **exteroceptive** sensors. Many of these sensors concentrate on a central task for the robot:

acquiring information on objects in the robot's immediate vicinity so that it may interpret the state of its surroundings.

Of course, these "objects" surrounding the robot are all detected from the viewpoint of its local reference frame.³ Since the systems we study are **mobile**, their ever-changing position and their motion has a significant impact on overall sensor behaviour.

³In this case we are referring to the robot reference frame.

Now that we have the necessary knowledge on the fundamental concepts and terminology, we can

now describe how dramatically the sensor error of an AMR **disagrees** with the ideal picture drawn in the previous section.

Blurring of Systematical and Random Errors

Active ranging sensors tend to have failure modes which are triggered largely by specific relative positions of the sensor and environment targets.

⁴The incident light is reflected into a single outgoing direction.

For example, a sonar sensor will produce specular reflections,⁴ producing grossly inaccurate measurements of range, at specific angles to a smooth sheet-rock wall.

During motion of the robot, such relative angles occur at stochastic intervals. This is especially true in a AMR outfitted with a ring of multiple sonars. The chances of one sonar entering this error mode during robot motion is high. From the perspective of the moving robot, the sonar measurement error is a **random error** in this case. However, if the robot were to stop, becoming motionless, then a very different error modality is possible.

If the robot's static position causes a particular sonar to fail in this manner, the sonar will fail consistently and will tend to return precisely the same (and incorrect!) reading time after time. Once the robot is motionless, the error appears to be systematic and high precision.

The fundamental mechanism at work here is the cross-sensitivity of AMR sensors to robot pose and robot-environment dynamics.

The models for such cross-sensitivity are **NOT**, in an underlying sense, truly random. However, these physical interrelationships are rarely modelled and therefore, from the point of view of an incomplete model, the errors appear random during motion and systematic when the robot is at rest. Sonar is not the only sensor subject to this blurring of systematic and random error modality. Visual interpretation through the use of a CCD camera is also highly susceptible to robot motion and position because of camera dependency on lighting.⁵

⁵such as glare and reflections.

The important point is to realise that, while systematic error and random error are well-defined in a controlled setting, the AMR can exhibit error characteristics that bridge the gap between deterministic and stochastic error mechanisms.

Multi-Modal Error Distributions

It is common to characterise the behaviour of a sensor's random error in terms of a probability distribution over various output values. In general, one knows very little about the causes of random error and therefore several simplifying assumptions are commonly used. For example, we can assume that the error is zero-mean ($\mu = 0$), in that it symmetrically generates both positive and negative measurement error. We can go even further and assume that the probability density curve is Gaussian.

Although we discuss the mathematics of this in detail later, it is important for now to recognise the fact that one frequently assumes symmetry as well as unimodal distribution. This means that measuring the correct value is most probable, and any measurement that is further away from the correct value is less likely than any measurement that is closer to the correct value. These are strong assumptions that enable powerful mathematical principles to be applied to AMR problems, but it is important to realise how wrong these assumptions usually are.

Consider, for example, the sonar sensor once again. When ranging an object that reflects the sound signal well, the sonar will exhibit high accuracy, and will induce random error based on noise, for example, in the timing circuitry. This portion of its sensor behaviour will exhibit error characteristics that are fairly **symmetric** and **unimodal**. However, when the sonar sensor is moving through an environment and is sometimes faced with materials that cause coherent reflection rather than returning the sound signal to the sonar sensor, then the sonar will grossly overestimate distance to the object. In such cases, the error will be biased toward positive measurement error and will be far from the correct value. The error is not strictly systematic, and so we are left modelling it as a probability distribution of random error. So the sonar sensor has two (2) separate types of operational modes, one in which the signal does return and some random error is possible, and the second in which the signal returns after a multi-path reflection, and gross overestimation error occurs. The probability distribution could easily be at least bimodal in this case, and since overestimation is more common than underestimation it will also be asymmetric.

As a second example, consider ranging via stereo vision. Once again, we can identify two (2) modes of operation. If the stereo vision system correctly correlates two images, then the resulting random error will be caused by camera noise and will limit the measurement accuracy. But the stereo vision system can also correlate two images incorrectly, matching two fence posts for example that are not the same post in the real world. In such a case stereo vision will exhibit gross measurement error, and one can easily imagine such behaviour violating both the unimodal and the symmetric assumptions. The thesis of this section is that sensors in a AMR may be subject to multiple modes of operation and, when the sensor error is characterised, uni modality and symmetry may be grossly violated. Nonetheless, as you will see, many successful AMR systems make use of these simplifying assumptions and the resulting mathematical techniques with great empirical success. The above sections have presented a terminology with which we can characterise the advantages and disadvantages of various mobile robot sensors. In the following sections, we do the same for a sampling of the most commonly used AMR sensors today.

2.1.4 Wheel and Motor Sensors

Wheel/motor sensors are devices used to measure the internal state and dynamics of a mobile robot. These sensors have vast applications outside of AMR and, as a result, AMR has enjoyed the benefits of high-quality, low-cost wheel and motor sensors which offer excellent resolution.

In the next part, we sample just one such sensor, the optical incremental encoder.

Optical Encoders

Optical incremental encoders have become the most popular device for measuring angular speed and position within a motor drive or at the shaft of a wheel or steering mechanism. In mobile robotics, encoders are used to control the position or speed of wheels and other motor-driven joints. Because these sensors are proprioceptive, their estimate of position is best in the reference frame of the robot and, when applied to the problem of robot localisation, significant corrections are required as discussed in Chapter 5.

An optical encoder is basically a mechanical light chopper that produces a certain number of sine or square wave pulses for each shaft revolution. It consists of an illumination source, a fixed grating that masks the light, a rotor disc with a fine optical grid that rotates with the shaft, and fixed optical detectors. As the rotor moves, the amount of light striking the optical detectors varies based on the alignment of the fixed and moving gratings. In robotics, the resulting sine wave is transformed into a discrete square wave using a threshold to choose between light and dark states. Resolution is measured in Cycles Per Revolution (CPR).

The minimum angular resolution can be readily computed from an encoder's CPR rating. A typical encoder in AMR may have 2,000 CPR while the optical encoder industry can readily manufacture encoders with 10,000 CPR. In terms of required bandwidth, it is of course critical that the encoder be sufficiently fast to count at the shaft spin speeds that are expected. Industrial optical encoders present no bandwidth limitation to AMR applications. Usually in AMR the quadrature encoder is used. In this case, a second illumination and detector pair is placed 90° shifted with respect to the original in terms of the rotor disc. The resulting twin square waves, shown in Fig. 4.2, provide significantly more information. The ordering of which square wave produces a rising edge first identifies the direction of rotation. Furthermore, the four detectability different states improve the resolution by a factor of four with no change to the rotor disc. Thus, a 2,000 CPR encoder in quadrature yields 8,000 counts. Further improvement is possible by retaining the sinusoidal wave measured by the optical detectors and performing sophisticated interpolation. Such methods, although rare in AMR, can yield 1000-fold improvements in resolution. As with most proprioceptive sensors, encoders are generally in the controlled environment of a AMR's internal structure, and so systematic error and cross-sensitivity can be engineered away. The accuracy of optical encoders is often assumed to be 100% and, although this may not entirely correct, any errors at the level of an optical encoder are dwarfed by errors downstream of the motor shaft.



Figure 2.1: An example of a rotary encoder. [32]

Heading Sensors

Heading sensors can be proprioceptive (gyroscope, inclinometer) or exteroceptive (compass). They are used to determine the robots orientation and inclination. They allow us, together with appropriate velocity information, to integrate the movement to a position estimate. This procedure,

which has its roots in vessel and ship navigation, is called dead reckoning.

Compasses

The two most common modern sensors for measuring the direction of a magnetic field are the Hall Effect and Flux Gate compasses. Each has advantages and disadvantages, as described below. The Hall Effect describes the behaviour of electric potential in a semiconductor when in the presence of a magnetic field. When a constant current is applied across the length of a semi-conductor, there will be a voltage difference in the perpendicular direction, across the semi-conductor's width, based on the relative orientation of the semiconductor to magnetic flux

lines. In addition, the sign of the voltage potential identifies the direction of the magnetic field. Thus, a single semiconductor provides a measurement of flux and direction along one dimension. Hall Effect digital compasses are popular in AMR, and contain two such semiconductors at right angles, providing two axes of magnetic field (thresholded) direction, thereby yielding one of 8 possible compass directions. The instruments are inexpensive but also suffer from a range of disadvantages. Resolution of a digital hall effect compass is poor. Internal sources of error include the nonlinearity of the basic sensor and systematic bias errors at the semiconductor level. The resulting circuitry must perform significant filtering, and this lowers the bandwidth of hall effect compasses to values that are slow in AMR terms. For example the hall effect compasses pictured in figure 4.3 needs 2.5 seconds to settle after a 90° spin. The Flux Gate compass operates on a different principle. Two small coils are wound on fer- rite cores and are fixed perpendicular to one-another. When alternating current is activated in both coils, the magnetic field causes shifts in the phase depending upon its relative alignment with each coil. By measuring both phase shifts, the direction of the magnetic field in two dimensions can be computed. The flux-gate compass can accurately measure the strength of a magnetic field and has improved resolution and accuracy; however it is both larger and more expensive than a Hall Effect compass. Regardless of the type of compass used, a major drawback concerning the use of the Earth's magnetic field for AMR applications involves disturbance of that magnetic field by other magnetic objects and man-made structures, as well as the bandwidth limitations of electronic compasses and their susceptibility to vibration. Particularly in indoor environments AMR applications have often avoided the use of compasses, although a compass can conceivably provide useful local orientation information indoors, even in the presence of steel structures.



Figure 2.2: An example of an electronic compass [33].

Gyroscope

Gyroscopes are heading sensors which preserve their orientation in relation to a fixed reference frame. Thus they provide an absolute measure for the heading of a mobile system. Gyroscopes

can be classified in two categories, mechanical gyroscopes and optical gyroscopes.

Mechanical Gyroscopes

The concept of a mechanical gyroscope relies on the inertial properties of a fast spinning rotor. The property of interest is known as the gyroscopic precession. If you try to rotate a fast spinning wheel around its vertical axis, you will feel a harsh reaction in the horizontal axis. This is due to the angular momentum associated with a spinning wheel and will keep the axis of the gyroscope inertially stable. The reactive torque τ and thus the tracking stability with the inertial frame are proportional to the spinning speed ω , the precession speed Ω and the wheel's inertia I .

$$\tau = I\omega\Omega$$

By arranging a spinning wheel as seen in Figure 4.4, no torque can be transmitted from the outer pivot to the wheel axis. The spinning axis will therefore be space-stable (i.e. fixed in an inertial reference frame). Nevertheless, the remaining friction in the bearings of the gyro- axis introduce small torques, thus limiting the long term space stability and introducing small errors over time. A high quality mechanical gyroscope can cost up to \$100,000 and has an angular drift of about 0.1̄ in 6 hours. For navigation, the spinning axis has to be initially selected. If the spinning axis is aligned with the north-south meridian, the earth's rotation has no effect on the gyro's horizontal axis. If it points east-west, the horizontal axis reads the earth rotation. Rate gyros have the same basic arrangement as shown in Figure 4.4 but with a slight modification. The gimbals are restrained by a torsional spring with additional viscous damping. This enables the sensor to measure angular speeds instead of absolute orientation.

Optical Gyroscopes

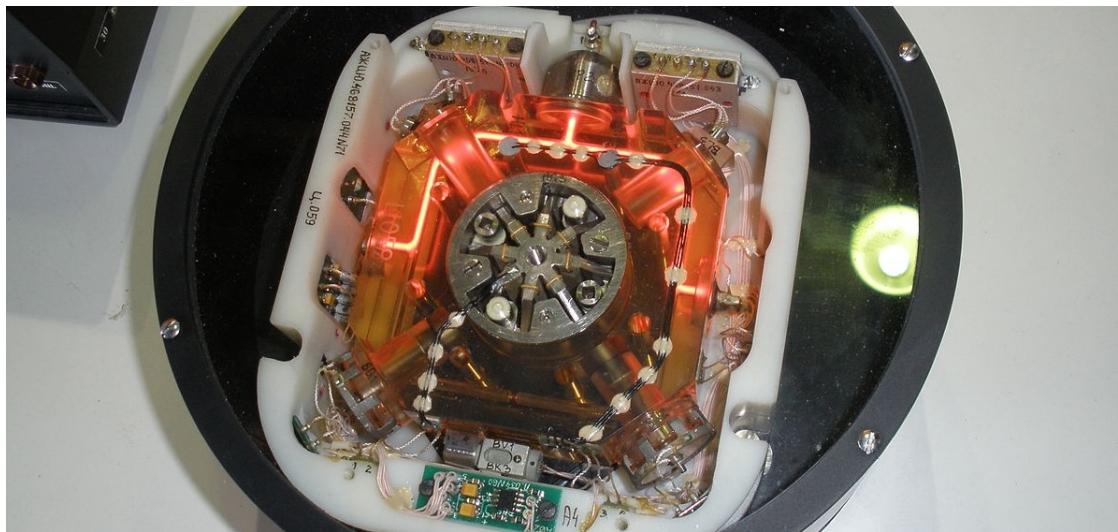


Figure 2.3: Optical Gyroscopes have no moving parts, (unlike mechanical gyroscopes) making them extremely reliable [34].

Optical gyroscopes are a relatively new innovation. Commercial use began in the early 1980's when they were first installed in aircraft. Optical gyroscopes are angular speed sensors that use two monochromatic light beams, or lasers, emitted from the same source instead of moving, mechanical parts. They work on the principle that the speed of light remains unchanged and, therefore, geometric change can cause light to take a varying amount of time to reach its destination. One laser beam is sent traveling clockwise through a fiber while the other travels counterclockwise. Because the laser traveling in the direction of rotation has a slightly shorter path, it will have a higher frequency. The difference in frequency of the two beams is proportional to the angular velocity of the cylinder. New solid-state optical gyroscopes based on the same principle are built using microfabrication technology, thereby providing heading information with resolution and bandwidth far beyond the needs of mobile robotic applications. Bandwidth, for instance, can easily exceed 100KHz while resolution can be smaller than 0.0001°/hr.

Ground Based Beacons



Figure 2.4

One elegant approach to solving the localization problem in AMR is to use active or passive beacons. Using the interaction of on-board sensors and the environmental beacons, the robot can identify its position precisely. Although the general intuition is identical to that of early human navigation beacons, such as stars, mountains and lighthouses, modern technology has enabled sensors to localize an outdoor robot with accuracies of better than 5 cm within areas that are kilometres in size.

In the following subsection, we describe one such beacon system, the Global Positioning System (GPS), which is extremely effective for outdoor ground-based and flying robots. Indoor beacon systems have been generally less successful for a number of reasons. The expense of environmental modification in an indoor setting is not amortized over an extremely large useful area, as it is for example in the case of GPS. Furthermore, indoor environments offer significant challenges not seen outdoors, including multipath and environment dynamics. A laser-based indoor beacon system, for example, must disambiguate the one true laser signal from possibly tens of other powerful signals that have reflected off of walls, smooth floors and doors. Confounding this, humans and other obstacles

may be constantly changing the environment, for example occluding the one true path from the beacon to the robot. In commercial applications such as manufacturing plants, the environment can be carefully controlled to ensure success. In less structured indoor settings, beacons have nonetheless been used, and the problems are mitigated by careful beacon placement and the useful of passive sensing modalities.

Global Positioning System

The Global Positioning System (GPS) was initially developed for military use but is now freely available for civilian navigation. There are at least 24 operational GPS satellites at all times. The satellites orbit every 12 hours at a height of 20.190km. There are four (4) satellites which located in each of six planes inclined 55° with respect to the plane of the earth's equator (figure 4.5).

Each satellite continuously transmits data which indicates its location and the current time. Therefore, GPS receivers are **completely passive** but **exteroceptive** sensors. The GPS satellites synchronise their transmissions to allow their signals to be sent at the same time. When a GPS receiver reads the transmission of two (2) or more satellites, the arrival time differences inform the receiver as to its relative distance to each satellite.

By combining information regarding the arrival time and instantaneous location of four (4) satellites, the receiver can infer its own position.

In theory, such triangulation requires only three (3) data points. However, timing is extremely critical in the GPS application because the time intervals being measured are in ns.

It is, of course, mandatory the satellites to be well synchronised. To this end, they are updated by ground stations regularly and each satellite carries on-board atomic clocks⁶ for timing. The GPS receiver clock is also important so that the travel time of each satellite's transmission can be accurately measured. But GPS receivers have a simple quartz clock. So, although 3 satellites would ideally provide position in three axes, the GPS receiver requires 4 satellites, using the additional information to solve for 4 variables: three position axes plus a time correction. The fact that the GPS receiver must read the transmission of 4 satellites simultaneously is a significant limitation. GPS satellite transmissions are extremely low-power, and reading them successfully requires direct line-of-sight communication with the satellite. Thus, in confined spaces such as city blocks with tall buildings or dense forests, one is unlikely to receive 4 satellites reliably. Of course, most indoor spaces will also fail to provide sufficient visibility of the sky for a GPS receiver to function. For these reasons, GPS has been a popular sensor in AMR, but has been relegated to projects involving AMR traversal of wide-open spaces and autonomous flying machines. A number of factors affect the performance of a localization sensor that makes use of GPS. First, it is important to understand that, because of the specific orbital paths of the GPS satellites, coverage is not geometrically identical in different portions of the Earth and therefore resolution is not uniform. Specifically, at the North and South poles, the satellites are very close to the horizon and, thus, while resolution in the latitude and longitude directions is good, resolution of altitude is relatively poor as compared to



⁶An example of a cesium clock for use in GPS.

more equatorial locations.

The second point is that GPS satellites are merely an information source. They can be employed with various strategies in order to achieve dramatically different levels of localisation resolution. The basic strategy for GPS use, called pseudorange and described above, generally performs at a resolution of 15m. An extension of this method is differential GPS, which makes use of a second receiver that is static and at a known exact position. A number of errors can be corrected using this reference, and so resolution improves to the order of 1m or less. A disadvantage of this technique is that the stationary receiver must be installed, its location must be measured very carefully and of course the moving robot must be within kilometers of this static unit in order to benefit from the DGPS technique. A further improved strategy is to take into account the phase of the carrier signals of each received satellite transmission. There are two carriers, at 19cm and 24cm, therefore significant improvements in precision are possible when the phase difference between multiple satellites is measured successfully. Such receivers can achieve 1cm resolution for point positions and, with the use of multiple receivers as in DGPS, sub-1cm resolution. A final consideration for AMR applications is bandwidth. GPS will generally offer no better than 200 - 300ms latency, and so one can expect no better than 5Hz GPS updates. On a fast-moving AMR or flying robot, this can mean that local motion integration will be required for proper control due to GPS latency limitations.

2.2 Active Ranging

Active range sensors continue to be the most popular sensors used in AMR. Many ranging sensors have a low price point, and most importantly all ranging sensors provide easily interpreted outputs:

Direct measurements of distance from the robot to objects in its vicinity.

For obstacle detection and avoidance, most AMR rely heavily on active ranging sensors. But the local free-space information provided by range sensors can also be accumulated into representations beyond the robot's current local reference frame. Therefore, active range sensors are also commonly found as part of the localisation and environmental modelling processes of AMRs.

It is only with the slow advent of successful visual interpretation competency that we can expect the class of active ranging sensors to gradually lose their primacy as the sensor class of choice among AMR engineers.

Below, we present two (2) Time-of-Flight (ToF) active range sensors:

- the ultrasonic sensor,
- the laser rangefinder.

Continuing onwards, we then present two (2) geometric active range sensors:

- the optical triangulation sensor,
- the structured light sensor.

Time-of-Flight Active Ranging

ToF ranging makes use of the [propagation speed of sound](#) or an [electromagnetic wave](#). In general, the travel distance of a sound or electromagnetic wave is given by:

$$d = ct,$$

where d is the distance travelled usually round-trip (m), c the speed of wave propagation (ms^{-1}), and t is the time it takes to travel (s).

It is important to point out the propagation speed v of sound is approximately 0.3 m ms^{-1} whereas the speed of an electromagnetic signal is 0.3 m ns^{-1} , which is one million times faster. The ToF for a typical distance, say 3 m, is 10 ms for an ultrasonic system but only 10 ns for a laser rangefinder. It is therefore obvious that measuring the time of flight t with electromagnetic signals is more technologically challenging.⁷

The quality of ToF range sensors depends mainly on the following:

⁷This explains why laser range sensors have only recently become affordable and robust for use on mobile robots.

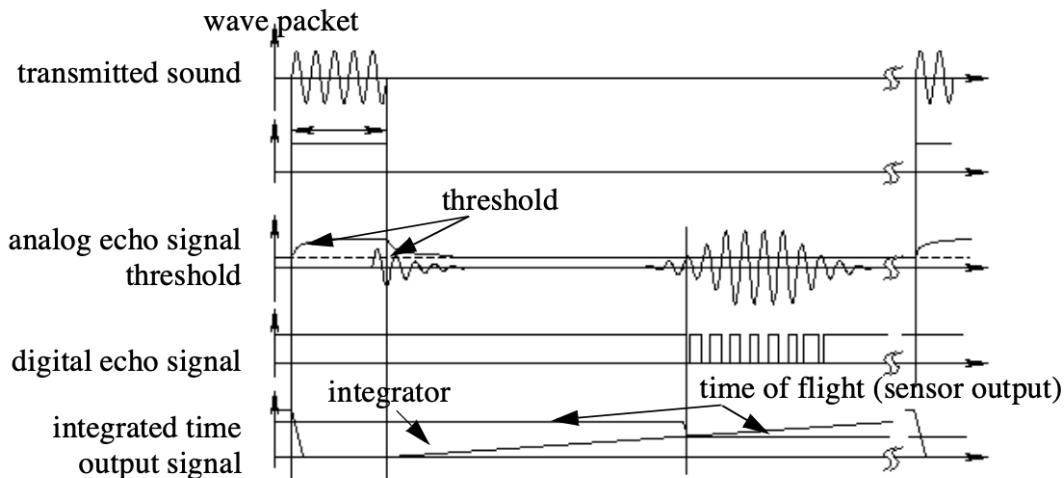


Figure 2.5: Signals of an ultrasonic sensor.

- Uncertainties in determining the exact time of arrival of the reflected signal,
- Inaccuracies in the time of flight measurement, particularly with laser range sensors,
- The dispersal cone of the transmitted beam mainly with ultrasonic range sensors
- Interaction with the target (e.g., surface absorption, specular reflections)
- Variation of propagation speed, and
- The speed of the AMR and target (in the case of a dynamic target).

As discussed below, each type of ToF sensor is sensitive to a particular subset of the above list of factors.

2.2.1 The Ultrasonic Sensor

The main ethos of an ultrasonic⁸ sensor is to transmit a packet of ultrasonic pressure waves and to measure the time it takes for this wave to reflect and return to the receiver. The distance d of the object causing the reflection can be calculated based on the propagation speed of sound⁹ c and the time of flight t .

$$d = \frac{c \times t}{2}$$

The speed of sound (v) in air is given by the following relation:

$$v = \sqrt{\gamma RT}$$

where γ is the ratio of specific heat, R is the gas constant ($\text{J mol}^{-1} \text{K}^{-1}$), and T is the temperature

⁸Ultrasound is sound with frequencies greater than 20 kHz.

⁹Of course in this regard careful consideration needs to be made if the medium is significantly different than that of air (i.e., water).

in Kelvin (K). In air, at standard pressure, and 20 °C the speed of sound is approximately:

$$v = 343 \text{ m s}^{-1}.$$

We can see the different signal output and input of an ultrasonic sensor in **Fig. 2.5**.

First, a series of sound pulses are emitted, which creates the wave packet. An integrator also begins to **linearly climb** in value, measuring the time from the transmission of these sound waves to detection of an echo. A threshold value is set for triggering an incoming sound wave as a valid echo.

This threshold is often decreasing in time, because the amplitude of the expected echo decreases over time based on dispersal as it travels longer.

But during transmission of the initial sound pulses and just afterwards, the threshold is set very high to suppress triggering the echo detector with the outgoing sound pulses. A transducer will continue to ring for up to several ms after the initial transmission, and this governs the blanking time of the sensor.

If, during the blanking time, the transmitted sound were to reflect off of an extremely close object and return to the ultrasonic sensor, it may fail to be detected.

However, once the blanking interval has passed, the system will detect any above-threshold reflected sound, triggering a digital signal and producing the distance measurement using the integrator value.

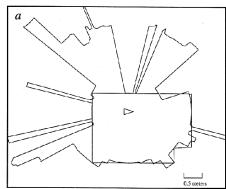
The ultrasonic wave typically has a frequency between 40 and 180 kHz and is usually generated by a piezo or electrostatic transducer. Often the same unit is used to measure the reflected signal, although the required blanking interval can be reduced through the use of separate output and input devices. Frequency can be used to select a useful range when choosing the appropriate ultrasonic sensor for a AMR. Lower frequencies correspond to a longer range, but with the disadvantage of longer post-transmission ringing and, therefore, the need for longer blanking intervals.

Most ultrasonic sensors used by AMRs have an effective range of roughly 12 cm to 5 metres. The published accuracy of commercial ultrasonic sensors varies between 98% and 99.1%. In AMR applications, specific implementations generally achieve a resolution of approximately 2 cm.

In most cases one may want a narrow opening angle for the sound beam in order to also obtain precise directional information about objects that are encountered. This is a major limitation since sound propagates in a cone-like manner with opening angles around 20° and 40°. Consequently, when using ultrasonic ranging one does not acquire depth data points but, rather, entire regions of constant depth. This means that the sensor tells us only that there is an object at a certain distance in within the area of the measurement cone. The sensor readings must be plotted as segments of an arc (sphere for 3D) and not as point measurements.¹⁰ However, recent research developments



Figure 2.6: An example of an ultrasonic sensor used in Raspberry Pi applications [35].



¹⁰The results of a 360° scan of a room.

show significant improvement of the measurement quality in using sophisticated echo processing. Ultrasonic sensors suffer from several additional drawbacks, namely in the areas of **error**, **bandwidth** and **cross-sensitivity**. The published accuracy values for ultrasonic sensors are nominal values based on successful, perpendicular reflections of the sound wave off an acoustically reflective material.

This does not capture the effective error modality seen on a AMR moving through its environment. As the ultrasonic transducer's angle to the object being ranged varies away from perpendicular, the chances become good that the sound waves will coherently reflect away from the sensor, just as light at a shallow angle reflects off of a mirror. Therefore, the true error behavior of ultrasonic sensors is compound, with a well-understood error distribution near the true value in the case of a successful retro-reflection, and a more poorly-understood set of range values that are grossly larger than the true value in the case of coherent reflection.

Of course the acoustic properties of the material being ranged have direct impact on the sensor's performance. Again, the impact is discrete, with one material possibly failing to produce a reflection that is sufficiently strong to be sensed by the unit. For example, foam, fur and cloth can, in various circumstances, acoustically absorb the sound waves. A final limitation for ultrasonic ranging relates to bandwidth. Particularly in moderately open spaces, a single ultrasonic sensor has a relatively slow cycle time.

For example, measuring the distance to an object that is 3 m away will take such a sensor 20ms, limiting its operating speed to 50 Hz. But if the robot has a ring of 20 ultrasonic sensors, each firing sequentially and measuring to minimize interference between the sensors, then the ring's cycle time becomes 0.4s and the overall update frequency of any one sensor is just 2.5 Hz. For a robot conducting moderate speed motion while avoiding obstacles using ultrasonic sensor, this update rate can have a measurable impact on the maximum speed possible while still sensing and avoiding obstacles safely.

Ultrasonic measurements may be limited through barrier layers with large salinity, temperature or vortex differentials.

Laser Rangefinder

The laser rangefinder is a ToF sensor which achieves significant improvements over the ultrasonic range sensor due to the **use of laser light instead of sound**. This type of sensor consists of a transmitter which illuminates a target with a collimated¹¹ beam (e.g. laser), and a receiver capable of detecting the component of light which is essentially coaxial with the transmitted beam. Often referred to as optical radar or Light Detection and Ranging (LIDAR), these devices produce a range estimate based on the time needed for the light to reach the target and return.

¹¹meaning all the rays in questions are made accurately parallel.

A mechanical mechanism with a mirror sweeps the light beam to cover the required scene in a plane or even in 3 dimensions, using a rotating mirror. One way to measure the ToF for the light beam is to use a pulsed laser and then measured the elapsed time directly, just as in the ultrasonic solution

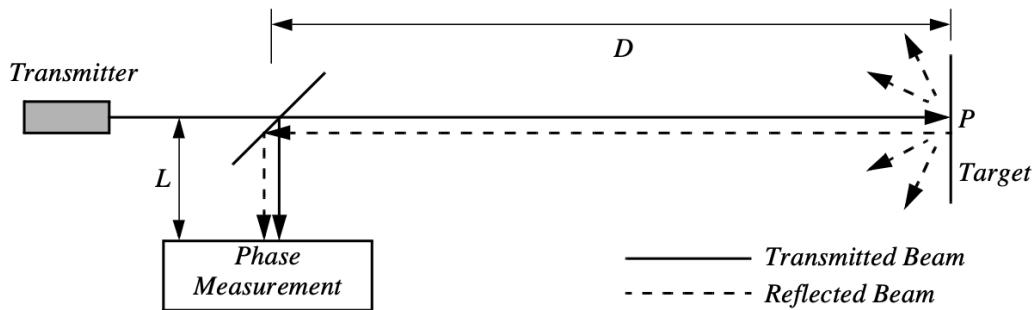


Figure 2.8: Schematic of laser rangefinding by phase-shift measurement.

described in just a little bit. Electronics capable of resolving ps are required in such devices and they are therefore very expensive. A second method is to measure the beat frequency between a frequency modulated continuous wave and its received reflection. Another, even easier method is to measure the phase shift of the reflected light.

Continuous Wave Radar It is a type of radar system where a known stable frequency continuous wave radio energy is transmitted and then received from any reflecting objects. Individual objects can be detected using the Doppler effect, which causes the received signal to have a different frequency from the transmitted signal, allowing it to be detected by filtering out the transmitted frequency.

Doppler-analysis of radar returns can allow the filtering out of slow or non-moving objects, thus offering immunity to interference from large stationary objects and slow-moving clutter. This makes it particularly useful for looking for objects against a background reflector, for instance, allowing a high-flying aircraft to look for aircraft flying at low altitudes against the background of the surface. Because the very strong reflection off the surface can be filtered out, the much smaller reflection from a target can still be seen.



Figure 2.7: A laser range finder used in robotics applications

Phase Shift Measurement Near infrared light, which could be from an Light-Emitting Diode (LED) or a laser, is collimated and transmitted from the transmitter T in Fig. 2.8 and hits a point P in the environment.

For surfaces having a roughness greater than the wavelength of the incident light, diffuse reflection will occur, meaning that the light is reflected almost isotropically¹². The wavelength of the infrared light emitted is 824 nm and so most surfaces with the exception of only highly polished reflecting objects, will be diffuse reflectors. The component of the infrared light which falls within the receiving aperture of the sensor will return almost parallel to the transmitted beam, for distant objects. The sensor transmits 100% amplitude modulated light at a known frequency and measures the phase

¹²Something that is isotropic has the same size or physical properties when it is measured in different directions

shift between the transmitted and reflected signals.

Fig. 2.9 shows how this technique can be used to measure range. The wavelength of the modulating signal obeys the equation $c = f\lambda$ where c is the speed of light and f the modulating frequency.

For example, $f = 5 \text{ MHz}$, the wavelength is $\lambda = 60 \text{ m}$.

The total distance D' covered by the emitted light is:

$$D' = L + 2D = L \frac{\theta}{2\pi} \lambda$$

where D and L are the distances defined in **Fig.** 2.8. The required distance D , between the beam splitter and the target, is therefore given by:

$$D = \frac{\lambda}{4\pi} \theta$$

where θ is the electronically measured phase difference between the transmitted and reflected light beams, and λ the known modulating wavelength. It can be seen that the transmission of a single frequency modulated wave can theoretically result in ambiguous range estimates since

For example if $\lambda = 60\text{m}$, a target at a range of 5 m would give an indistinguishable phase measurement from a target at 65 m , since each phase angle would be 360° apart.

We therefore define an **ambiguity interval** of λ , but in practice we note that the range of the sensor is much lower than λ due to the attenuation of the signal in air. It can be shown that the confidence in the range (phase estimate) is inversely proportional to the square of the received signal amplitude, directly affecting the sensor's accuracy. Hence dark, distant objects will not produce as good range estimates as close, bright objects.

As with ultrasonic ranging sensors, an important error mode involves coherent reflection of the energy. With light, this will only occur when striking a highly polished surface. Practically, a AMR may encounter such surfaces in the form of a polished desktop, file cabinet or of course a mirror. Unlike ultrasonic sensors, laser rangefinders cannot detect the presence of optically transparent materials such as glass, and this can be a significant obstacle in environments, for example museums, where glass is commonly used.

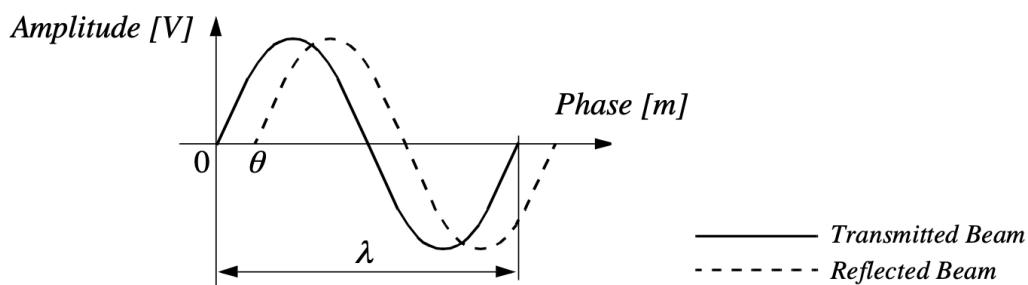


Figure 2.9: Range estimation by measuring the phase shift between transmitted and received signals.

Triangulation-based Active Ranging

Triangulation-based ranging sensors use geometrical properties in their measuring strategy to establish distance readings to objects. The simplest class of triangulation-based rangers are active because they project a known light pattern (e.g., a point, a line or a texture) onto the environment. The reflection of the known pattern is captured by a receiver and, together with known geometric values, the system can use simple triangulation to establish range measurements. If the receiver measures the position of the reflection along a single axis, we call the sensor an optical triangulation sensor in 1D. If the receiver measures the position of the reflection along two orthogonal axes, we call the sensor a structured light sensor.

Optical Triangulation (1D Sensor)

The principle of optical triangulation in 1D is straightforward, as depicted in **Fig. 2.10**. A collimated

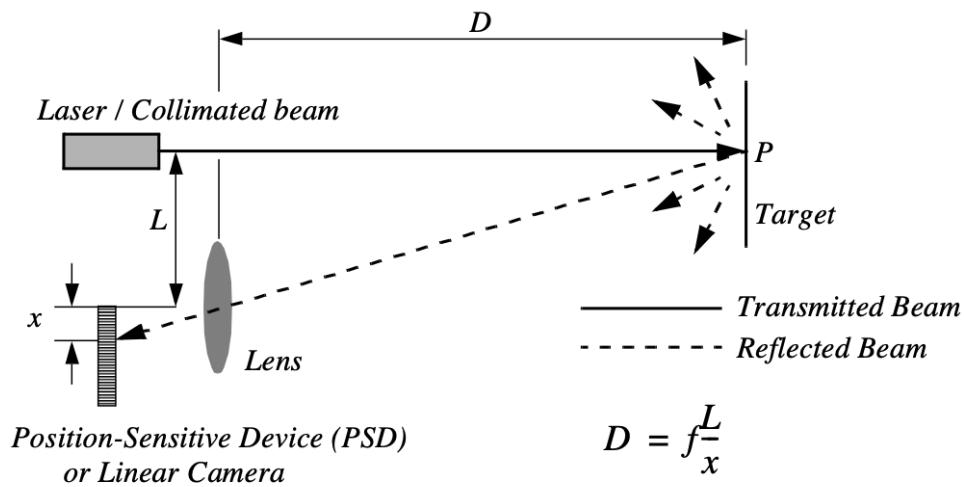


Figure 2.10: Principle of 1D laser triangulation.

beam is transmitted toward the target. The reflected light is collected by a lens and projected onto a position sensitive device¹³ or linear camera. Given the geometry of **Fig. 2.10** the distance D is given by:

$$D = f \frac{L}{x}$$

The distance is proportional to $\frac{1}{x}$, therefore the sensor resolution is best for close objects and becomes worse as distance increases. Sensors based on this principle are used in range sensing up to one or two m, but also in high precision industrial measurements with resolutions far below one μm . Optical triangulation devices can provide relatively high accuracy with very good resolution for close objects. However, the operating range of such a device is normally fairly limited by **geometry**. For



¹³A position sensitive device and/or position sensitive detector is an optical position sensor which can measure a position of a light spot in one or two-dimensions on a sensor surface.

example, an off-the-shelf optical triangulation sensor can operate over a distance range of between 8 cm and 80 cm.

It is inexpensive compared to ultrasonic and laser rangefinder sensors.

Although more limited in range than sonar, the optical triangulation sensor has high bandwidth and does not suffer from cross-sensitivities that are more common in the sound domain.

Structured Light (2D Sensor)

If one replaced the linear camera or Position Sensing Device (PSD) of an optical triangulation sensor with a two-dimensional receiver such as a CCD or CMOS camera, then one can recover distance to a large set of points instead of to only one point. The emitter must project a known pattern, or structured light, onto the environment. Many systems exist which either project light textures, which can be seen in **Fig. 2.12**, or emit collimated light by means of a rotating mirror. Yet another popular alternative is to project a laser stripe by turning a laser beam into a plane using a prism. Regardless of how it is created, the projected light has a known structure, and therefore the image taken by the CCD or CMOS receiver can be filtered to identify the pattern's reflection.

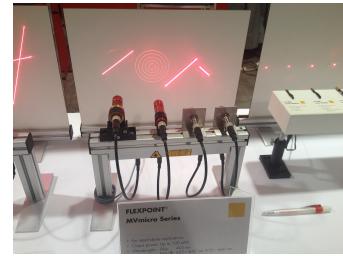


Figure 2.11: Structured light sources on display at the 2014 Machine Vision Show in Boston [36].

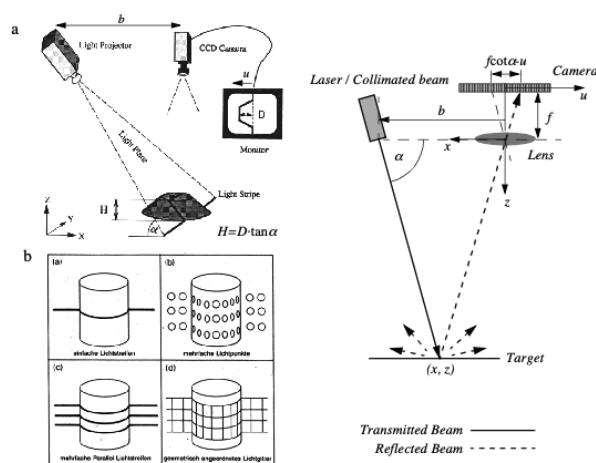


Figure 2.12: a) Principle of active two dimensional triangulation b) Other possible light structures c) One-dimensional schematic of the principle

The problem of recovering depth here far simpler than the problem of passive image analysis.

In passive image analysis, as we discuss later, existing features in the environment must be used to perform correlation, while the present method projects a **known pattern upon the environment** and thereby avoids the standard correlation problem altogether. Furthermore, the structured light sensor

is an active device; so, it will continue to work in dark environments as well as environments in which the objects are featureless¹⁴. In contrast, stereo vision would fail in such texture-free circumstances.

Figure 4.15c shows a one-dimensional active triangulation geometry. We can examine the trade-off in the design of triangulation systems by examining the geometry in figure 4.15c. The measured values in the system are α and u , the distance of the illuminated point from the origin in the imaging sensor.¹⁵ From figure 4.15c, simple geometry shows that:

$$x = \frac{bu}{f \cot \alpha - u} \quad \text{and} \quad z = \frac{bf}{f \cot \alpha - u}.$$

where f is the distance of the lens to the imaging plane. In the limit, the ratio of image resolution to range resolution is defined as the triangulation gain G_p and from equation 4.12 is given by:

$$\frac{\partial u}{\partial z} = G_p = \frac{bf}{z^2}$$

This shows that the ranging accuracy, for a given image resolution, is proportional to source/detector separation b and focal length f , and decreases with the square of the range z . In a scanning ranging system, there is an additional effect on the ranging accuracy, caused by the measurement of the projection angle α . From equation 4.12 we see that:

$$\frac{\partial \alpha}{\partial z} = G_{ff} = \frac{b \sin \alpha^2}{z^2}$$

We can summarise the effects of the parameters on the sensor accuracy as follows:

Baseline Length (b) the smaller b is the more compact the sensor can be. The larger b is the better the range resolution will be. Note also that although these sensors do not suffer from the correspondence problem, the disparity problem still occurs. As the baseline length b is increased, one introduces the chance that, for close objects, the illuminated point(s) may not be in the receiver's field of view.

Detector length and focal length f A larger detector length can provide either a larger field of view or an improved range resolution or partial benefits for both. Increasing the detector length however means a larger sensor head and worse electrical characteristics (increase in random error and reduction of bandwidth). Also, a short focal length gives a large field of view at the expense of accuracy and vice versa.

At one time, laser stripe-based structured light sensors were common on several mobile robot bases as an inexpensive alternative to laser range-finding devices. However, with the increasing quality of laser range-finding sensors in the 1990's the structured light system has become relegated largely to vision research rather than applied mobile robotics.

2.2.2 Motion and Speed Sensors

Some sensors directly measure the relative motion between the robot and its environment. Since such motion sensors detect **relative motion**, so long as an object is moving relative to the robot's

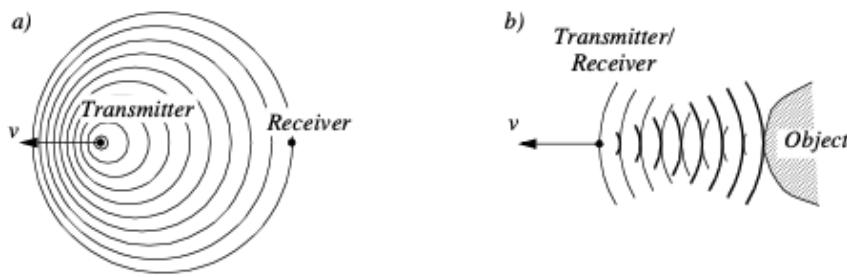


Figure 2.13: Doppler effect between two moving objects (a) or a moving and a stationary object(b)

reference frame, it will be detected and its speed can be estimated. There are a number of sensors that inherently measure some aspect of motion or change.

For example, a pyroelectric¹⁶ sensor detects change in heat.

When someone walks across the sensor's field of view, his motion triggers a change in heat in the sensor's reference frame. In the next subsection, we describe an important type of motion detector based on the **Doppler effect**. These sensors represent a well-known technology with decades of general applications behind them.



¹⁶An example of a pyroelectric sensor.

For fast-moving AMRs such as autonomous highway vehicles and unmanned flying vehicles, Doppler-based motion detectors are the obstacle detection sensor of choice.

Doppler Effect

Anyone who has noticed the change in siren pitch when an ambulance approaches and then passes by is familiar with the Doppler effect.¹⁷

A transmitter emits an electromagnetic or sound wave with a frequency f_t . It is either received by a receiver **Fig. 2.13(a)** or reflected from an object **Fig. 2.13 (b)**. The measured frequency f_r at the receiver is a function of the relative speed v between transmitter and receiver according to

$$f_r = f_t \frac{1}{1 + \frac{v}{c}}$$

if the transmitter is moving and

$$f_r = f_t \left(1 + \frac{v}{c} \right)$$

if the receiver is moving. In the case of a reflected wave **Fig. 2.13 (b)** there is a factor of two introduced, since any change x in relative separation affects the round-trip path length by $2x$.

In such situations it is generally more convenient to consider the change in frequency Δf , known as the Doppler shift, as opposed to the Doppler frequency notation above.

¹⁷For anyone who needs a bit more information, it is the change in the frequency of a wave in relation to an observer who is moving relative to the source of the wave. The Doppler effect is named after the physicist Christian Doppler, who described the phenomenon in 1842. A common example of Doppler shift is the change of pitch heard when a vehicle sounding a horn approaches and recedes from an observer. Compared to the emitted frequency, the received frequency is higher during the approach, identical at the instant of passing by, and lower during the recession.

$$\Delta f = f_t - f_r = \frac{2f_t v \cos \theta}{c} \quad \text{and} \quad v = \frac{\Delta f c}{2f_t \cos \theta}$$

A current application area is both autonomous and manned highway vehicles. Both micro-wave and laser radar systems have been designed for this environment. Both systems have equivalent range, but laser can suffer when visual signals are deteriorated by environmental conditions such as rain, fog, etc. Commercial microwave radar systems are already available for installation on highway trucks. These systems are called VORAD (vehicle on-board radar) and have a total range of approximately 150m. With an accuracy of approximately 97%, these systems report range rate from 0 to 160 km/hr with a resolution of 1 km/ hr. The beam is approximately 4° wide and 5° in elevation. One of the key limitations of radar technology is its bandwidth. Existing systems can provide information on multiple targets at approximately 2 Hz.

2.3 Vision Based Sensors

Vision is our most powerful sense. It provides us with an enormous amount of information about the environment and enables rich, intelligent interaction in dynamic environments. It is therefore not at all surprising that a great deal of effort has been devoted to providing machines with sensors which can at least try to mimic the capabilities of the human vision system.

The first step in this process is the creation of sensing devices that capture the same raw information which is the light the human vision system uses. The main topics which will be described are the two (2) current technologies for creating vision sensors:

1. CCD,
2. CMOS.

Of course, these sensors have specific limitations in performance compared to the human eye, and it is important to understand these limitations. Later sections describe vision-based sensors which are commercially available, similar to the sensors discussed previously, along with their disadvantages and most popular applications.

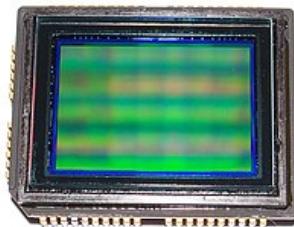


Figure 2.14: Sony ICX493AQ 10.14-megapixel APS-C (23.4 × 15.6 mm) CCD from digital camera Sony DSLR-A200 or DSLR-A300, sensor side [37].

CCD and CMOS Sensors

When it comes to the marketplace, CCD is the most popular fundamental ingredient for robotic vision systems.¹⁸ The CCD chip, which you can see in **Fig. 2.14** is an array of light-sensitive picture elements, or pixels, usually with between 20 000 and 2 million pixels total.

Each pixel can be thought of as a **light-sensitive, discharging capacitor** that is 5 to 25 µm in size. First, the capacitors of all pixels are fully charged, then the integration period begins. As photons of light strike each pixel, the electrons are liberated, which are captured by electric fields and retained at the pixel. Over time, each pixel accumulates a varying level of charge based on the total number of photons that have struck it. After the integration period is complete, the relative charges of all pixels need to be **frozen and read**.

In a CCD, the reading process is performed at one corner of the CCD chip.¹⁹ The bottom row of pixel charges are transported to this corner and read, then the rows above shift down and the process repeats. This means that each charge **must be transported across the chip**, and it is critical the value be preserved.

This requires specialised control circuitry and custom fabrication techniques to ensure the stability of transported charges.

¹⁸Willard Boyle and George E. Smith invented the CCD in 1969 at AT&T Bell Labs. Their original idea was to create a memory device. However, with its publication in 1970, other scientists began experimenting with the technology on a range of applications. Astronomers discovered that they could produce high-resolution images of distant objects, because CCDs offered a photo-sensitivity one hundred times greater than film [38].

¹⁹Because the entire array is read through a single amplifier the output can be highly optimised to give very low noise and extremely high dynamic range. CCDs can have over 100 dB dynamic range with less than 2e of noise [38].

²⁰This also includes CMOS as well.

The photo-diodes used in CCD chips²⁰ are **NOT** equally sensitive to all frequencies of light. They are sensitive to light between 400 nm and 1000 nm wavelength.²¹

²¹This number range is usually given for easier numbers as both CCD and CMOS have sensitivity values at approximately 350 - 1050 nm.

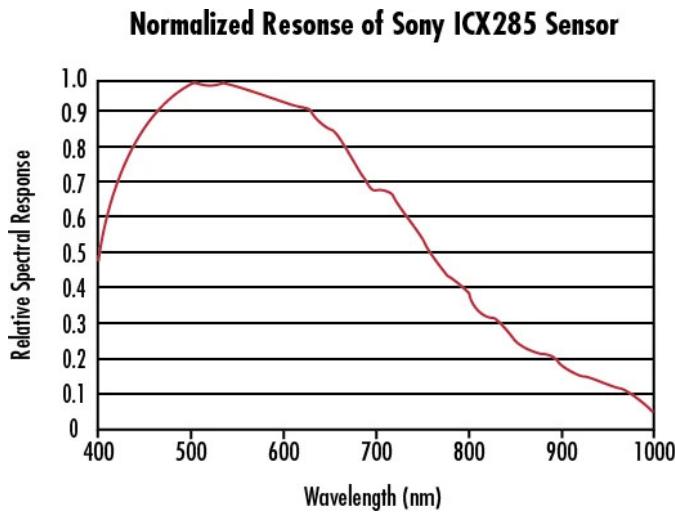


Figure 2.15: Normalized Spectral Response of a Typical Monochrome CCD.

It is important to remember that photodiodes are **less sensitive to the ultraviolet** part of the spectrum and are overly **sensitive to the infrared** portion (e.g. heat) which you can see in Fig. 2.15. You can see that the basic light-measuring process is colourless.²²

There are two (2) common approaches for creating color images. If the pixels on the CCD chip are grouped into 2-by-2 sets of four (4), then red, green and blue dyes can be applied to a colour filter so each individual pixel receives only light of just one color.

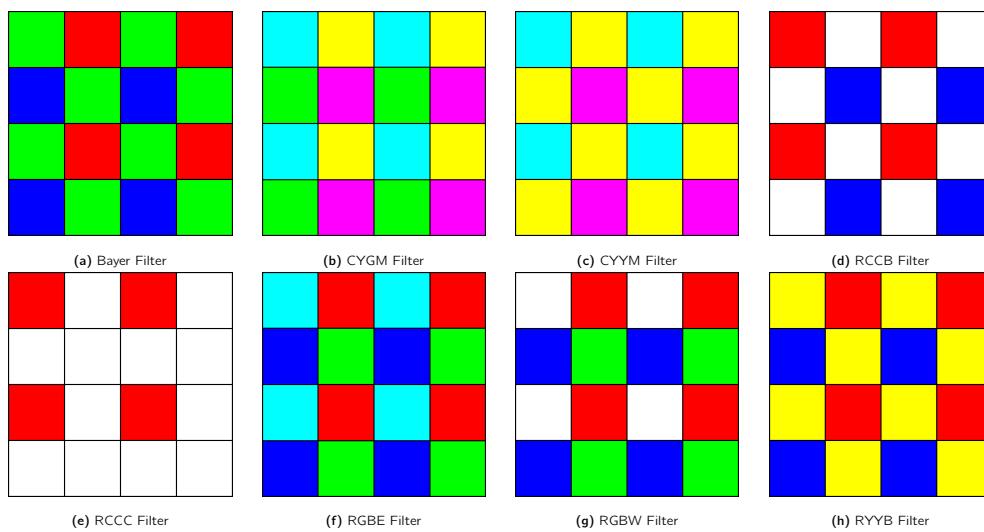


Figure 2.16: Types of colour filter used in commercial and industrial applications

Normally, two (2) pixels measure green while one pixel each measures red and blue light intensity. Of course, this 1-chip color CCD has a geometric resolution disadvantage.

The number of pixels in the system has been effectively cut by a factor of 4, and therefore the image resolution output by the CCD camera will be sacrificed.

The 3-chip color camera avoids these problems by splitting the incoming light into three (3) complete²³ copies. Three separate CCD chips receive the light, with one red, green or blue filter over each entire chip. Thus, in parallel, each chip measures light intensity for just one color, and the camera must combine the CCD chips' outputs to create a joint color image.

²³Albeit, with lower resolution.

Resolution is preserved in this solution, although the 3-chip color cameras are, as one would expect, significantly more expensive and therefore more rarely used in mobile robotics.

Both 3-chip and single chip color CCD cameras suffer from the fact that photo-diodes are much more sensitive to the near-infrared end of the spectrum. This means that the overall system detects blue light much more poorly than red and green. To compensate, the gain must be increased on the blue channel, and this introduces greater absolute noise on blue²⁴ than on red and green. It is not uncommon to assume at least 1 - 2 bits of additional noise on the blue channel.

²⁴This is generally defined as the amplifier noise.

The CCD camera has several camera parameters that affect its behavior. In some cameras, these parameter values are fixed. In others, the values are constantly changing based on built-in feedback loops. In higher-end cameras, the user can modify the values of these parameters via software embedded into the device. The iris position and shutter speed²⁵ regulate the amount of light being measured by the camera. The iris is simply a mechanical aperture that constricts incoming light, just as in standard 35mm cameras. Shutter speed regulates the integration period of the chip. In higher-end cameras, the effective shutter speed can be as brief at 1/30,000s and as long as 2s. Camera gain controls the overall amplification of the analog signal, prior to A/D conversion. However, it is very important to understand that, even though the image may appear brighter after setting high gain, the shutter speed and iris may not have changed at all. Thus gain merely amplifies the signal, and amplifies along with the signal all of the associated noise and error. Although useful in applications where imaging is done for human consumption (e.g. photography, television), gain is of little value to a mobile roboticist.

²⁵It's the speed at which the shutter of the camera closes. A fast shutter speed creates a shorter exposure - the amount of light the camera takes in - and a slow shutter speed gives a longer exposure.

In colour cameras, an additional control exists for white balance. Depending on the source of illumination in a scene²⁶ the relative measurements of red, green and blue light which combine to define pure white light will change dramatically which can be seen in **Fig. 2.17** which can also be adjusted with algorithms [39]. The human eyes compensate for all such effects in ways that are not fully understood, however, the camera can demonstrate glaring inconsistencies in which the same table looks blue in one image, taken during the night, and yellow in another image, taken during the day. White balance controls enable the user to change the relative gain for red, green and blue in order to maintain more consistent color definitions in varying contexts.

²⁶For example this could be fluorescent lamps, incandescent lamps, sunlight, underwater filtered light, etc.

The key disadvantages of CCD cameras are primarily in the areas of inconstancy and **dynamic range**.

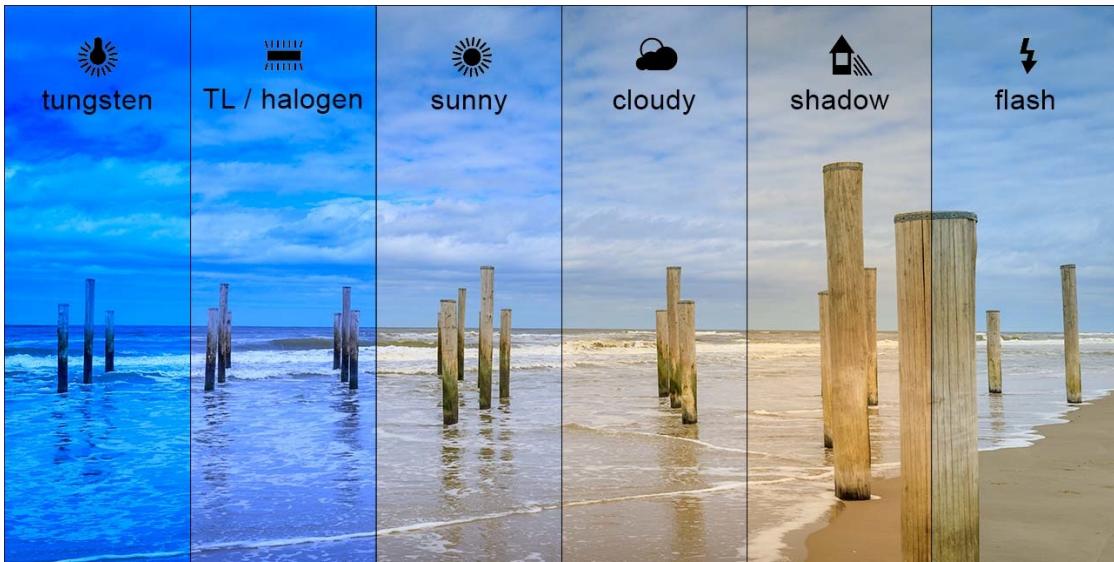


Figure 2.17: Example of white balance. Here the same scene is emulated to be shot under different light conditions [40].

Dynamic Range

Dynamic range in photography describes the ratio between the maximum and minimum measurable light intensities (white and black, respectively). In the real world, one never encounters true white or black - only varying degrees of light source intensity and subject reflectivity. Therefore the concept of dynamic range becomes more complicated, and depends on whether you are describing a capture device (such as a camera or scanner), a display device (such as a print or computer display), or the subject itself.

As mentioned above, a number of parameters can change the brightness and colours with which a camera creates its image.

Manipulating these parameters in a way to provide consistency over time and over environments, for example ensuring a green shirt always looks green, and something dark grey is always dark grey, remains an open problem [41].

The second type of disadvantages relates to the behavior of a CCD chip in environments with **extreme illumination**. In cases of very low illumination, each pixel will receive only a small number of photons. The longest possible shutter speed and camera optics (i.e. pixel size, chip size, lens focal length and diameter) will determine the minimum level of light for which the signal is stronger than random error noise. In cases of very high illumination, a pixel fills its well with free electrons and, as the well reaches its limit, the probability of trapping additional electrons falls and therefore the linearity between incoming light and electrons in the well degrades. This is termed **saturation**²⁷ and can indicate the existence of a further problem related to cross-sensitivity [43]. When a well has reached its limit, then additional light within the remainder of the integration period may cause further charge to leak into neighbouring pixels, causing them to report incorrect values or even reach secondary saturation. This effect, called **blooming**, means that individual pixel values are **NOT** truly



²⁷Example of blooming caused by saturation of a sensor pixel. The sun is so bright in the image that there is blooming on the sun itself, leaking into the surrounding pixels, and a vertical smear across the whole image [42].

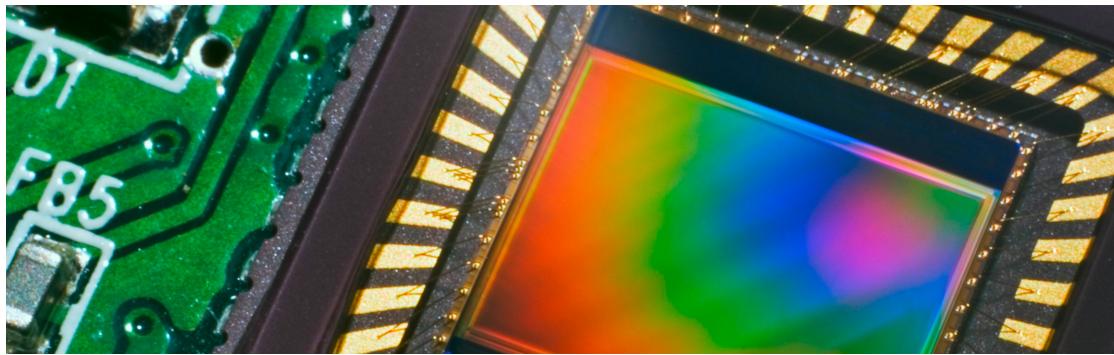


Figure 2.18: A close-up view of a CMOS sensor and its circuitry [44].

independent. The camera parameters may be adjusted for an environment with a particular light level, but the problem remains that the dynamic range of a camera is limited by the well capacity of the individual pixels.

For example, a high quality CCD may have pixels that can hold 40 000 electrons. The noise level for reading the well may be 11 electrons, and therefore the dynamic range will be 40,000:11, or 3,600:1, which is 35 dB.

2.3.1 CMOS Technology

The Complementary Metal Oxide Semiconductor (CMOS) chip is a significant departure from the CCD. Similar to CCD, it too has an array of pixels, but located alongside each pixel are **several transistors specific to that pixel**. Just as in CCD chips, all of the pixels accumulate charge during the integration period. During the data collection step, the CMOS takes a new approach:

The pixel-specific circuitry next to every pixel measures and amplifies the pixel's signal, all in parallel for every pixel in the array.

Using more traditional traces from general semiconductor chips, the resulting pixel values are all carried to their destinations. CMOS has a number of advantages over CCD technologies. First and foremost, there is no need for the specialized clock drivers and circuitry required in the CCD to transfer each pixel's clock down all of the array columns and across all of its rows.²⁸

This also means that specialized semiconductor manufacturing processes are not required to create CMOS chips.

Therefore, the same production lines that create microchips can create inexpensive CMOS chips as well. The CMOS chip is so much simpler that it consumes significantly less power, it operates with a power consumption a tenth the power consumption of a CCD chip [46].

In a AMR, power is a scarce resource and therefore this is an important advantage.



²⁸-CAM80CUNX is an 8MP Ultra-lowlight MIPI CSI-2 camera capable of streaming 4K @ 44 fps. This 8MP camera is based on SONY STARVIS IMX415 CMOS image sensor [45]

On the other hand, the CMOS chip also faces several disadvantages.

- Most importantly, the circuitry next to each pixel consumes valuable real estate on the face of the light-detecting array. Many photons hit the transistors rather than the photodiode, making the CMOS chip significantly less sensitive than an equivalent CCD chip.
- CMOS, compared to CCD is still finding ground in the marketplace, and as a result, the best resolution that one can purchase in CMOS format continues to be far inferior to the best CCD chips available.
- CMOS sensors have a lower dynamic range,
- CMOS sensors have higher levels of noise.

Compared to the human eye, these chips all have worse performance, cross-sensitivity and a limited dynamic range. As a result, vision sensors today continue to be fragile. Only over time, as the underlying performance of imaging chips improves, will significantly more robust vision-based sensors for AMRs be available.

Shot Noise

Shot noise or Poisson noise is a type of noise which can be modeled by a Poisson process. In electronics shot noise originates from the discrete nature of electric charge. Shot noise also occurs in photon counting in optical devices, where shot noise is associated with the particle nature of light.

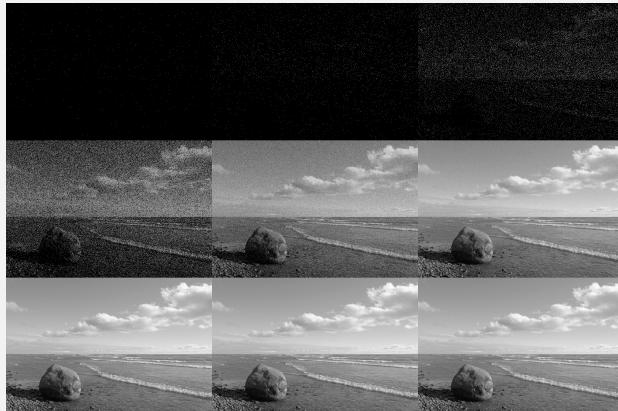


Figure 2.19: Photon noise simulation. Number of photons per pixel increases from left to right and from upper row to bottom row [47].

2.3.2 Visual Ranging Sensors

Range sensing is extremely important in AMR as it is a basic input for successful obstacle avoidance. As we have seen earlier, a number of sensors are popular in robotics specifically for their ability to recover depth estimates:

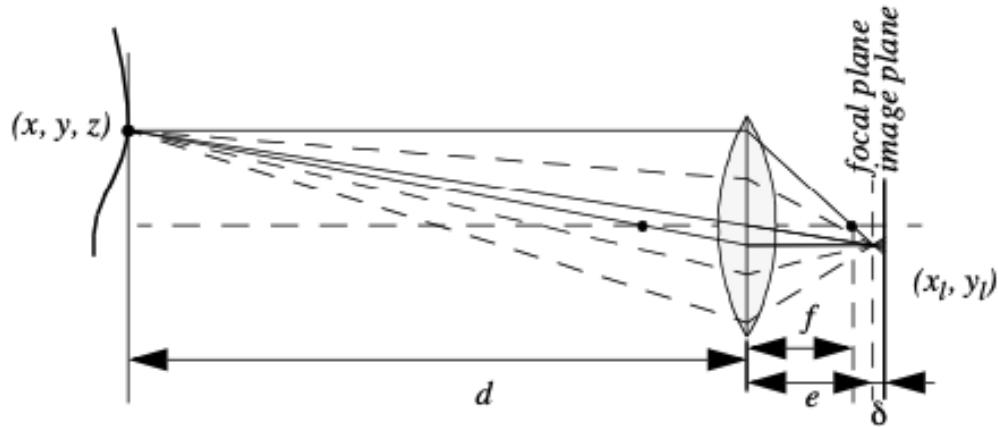


Figure 2.20: Depiction of the camera optics and its impact on the image. To get a sharp image, the image plane must coincide with the focal plane. Otherwise the image of the point (x, y, z) will be blurred in the image as can be seen in the drawing above.

ultrasonic, laser rangefinder, optical rangefinder, etc.

It is natural to attempt to implement ranging functionality using vision chips as well. However, a fundamental problem with visual images makes rangefinding relatively difficult.

Any vision chip collapses the three-dimensional world into a two-dimensional image plane, thereby losing depth information. If one can make strong assumptions regarding the size of objects in the world, or their particular colour and reflectance, then one can directly interpret the appearance of the two-dimensional image to recover depth. But such assumptions are rarely possible in real-world AMR applications.

Without such assumptions, a single picture does not provide enough information to recover spatial information.

The general solution is to recover depth by looking at several images of the scene to gain more information, which will be hopefully enough to at least partially recover depth. The images used **must be different**, so that taken together they provide additional information. They could differ in viewpoint, which would allow the use of stereo or motion algorithms.

An alternative is to create different images, not by changing the viewpoint, but by changing the camera geometry, such as the focus position or lens iris. This is the fundamental idea behind depth from focus and depth from defocus techniques. We will now look into the general approach to the depth from focus techniques as it presents a straightforward and efficient way to create a vision-based range sensor.

2.3.3 Depth from Focus

The depth from focus class of techniques relies on the fact that image properties not only change as a function of the **scene**, but also as a function of the **camera parameters**. The relationship between camera parameters and image properties is depicted in **Fig. 2.20**. The fundamental formula governing image formation relates the distance of the object from the lens, d in **Fig. 2.20**, to the distance e from the lens to the focal point, based on the focal length f of the lens:

$$\frac{1}{f} = \frac{1}{d} + \frac{1}{e}$$

²⁹A three-dimensional counterpart to a pixel. If the image plane is located at distance e from the lens, then for the specific object voxel²⁹ depicted, all light will be focused at a single point on the image plane and the object voxel will be focused. However, when the image plane is **NOT** at e , as is seen in **Fig. 2.20**, then the light from the object voxel will be cast on the image plane as a **blur circle**. To a first approximation, the light is homogeneously distributed throughout this blur circle, and the radius R of the circle can be characterized according to the equation:

$$R = \frac{L\delta}{2e}$$

where L is the diameter of the lens or aperture and δ is the displacement of the image plan from the focal point.

Given these formulae, several basic optical effects are clear.

³⁰The aperture is the opening in the lens that allows light to enter the camera and onto the sensor or film.

For example, if the aperture³⁰ or lens is reduced to a point, as in a pin-hole camera, then the radius of the blur circle approaches zero.

This is consistent with the fact that decreasing the iris aperture opening causes the depth of field to increase until all objects are in focus. Of course, the disadvantage of doing so is that we are allowing less light to form the image on the image plane and so this is practical only in bright circumstances. The second property to be deduced from these optics equations relates to the sensitivity of blurring as a function of the distance from the lens to the object.

Suppose the image plane is at a fixed distance 1.2 from a lens with diameter $L = 0.2$ and focal length $f = 0.5$. We can see from Equation (4.20) that the size of the blur circle R changes proportionally



Figure 2.21: Three images of the same scene taken with a camera at three different focusing positions. Note the significant change in texture sharpness between the near surface and far surface [48].

with the image plane displacement . If the object is at distance $d = 1$, then from Equation (4.19) we can compute $e=1$ and therefore $= 0.2$. Increase the object distance to $d = 2$ and as a result $= 0.533$. Using Equation (4.20) in each case we can compute $R = 0.02$ $R = 0.08$ respectively. This demonstrates high sensitivity for defocusing when the object is close to the lens. In contrast suppose the object is at $d = 10$. In this case we compute $e = 0.526$. But if the object is again moved one unit, to $d = 11$, then we compute $e = 0.524$. Then resulting blur circles are $R = 0.117$ and $R = 0.129$, far less than the quadrupling in R when the obstacle is $1/10$ the distance from the lens. This analysis demonstrates the fundamental limitation of depth from focus techniques: they lose sensitivity as objects move further away (given a fixed focal length). Interestingly, this limitation will turn out to apply to virtually all visual ranging techniques, including depth from stereo and depth from motion. Nevertheless, camera optics can be customised for the depth range of the intended application. For example, a "zoom" lens with a very large focal length f will enable range resolution at significant distances, of course at the expense of field of view. Similarly, a large lens diameter, coupled with a very fast shutter speed, will lead to larger, more detectable blur circles. Given the physical effects summarised by the above equations, one can imagine a visual ranging sensor that makes use of multiple images in which camera optics are varied (e.g. image plane displacement) and the same scene is captured (see Fig. 4.20). In fact this approach is not a new invention. The human visual system uses an abundance of cues and techniques, and one system demonstrated in humans is depth from focus. Humans vary the focal length of their lens continuously at a rate of about 2 Hz. Such approaches, in which the lens optics are actively searched in order to maximise focus, are technically called depth from focus. In contrast, depth from defocus means that depth is recovered using a series of images that have been taken with different camera geometries. Depth from focus methods are one of the simplest visual ranging techniques. To determine the range to an object, the sensor simply moves the image plane (via focusing) until maximizing the sharpness of the object. When the sharpness is maximised, the corresponding position of the image plane directly reports range. Some autofocus cameras and virtually all autofocus video cameras use this technique. Of course, a method is required for measuring the sharpness of an image or an object within the image. The most common techniques are approximate measurements of the sub-image gradient:

$$\text{sharpness}_1 = \sum_{x, y} |I(x, y) - I(x - 1, y)| \quad (2.3)$$

$$\text{sharpness}_2 = \sum_{x, y} (I(x, y) - I(x - 2, y - 2))^2 \quad (2.4)$$

A significant advantage of the horizontal sum of differences technique (Equation (4.21)) is that the calculation can be implemented in analog circuitry using just a rectifier, a low-pass filter and a high-pass filter. This is a common approach in commercial cameras and video recorders. Such systems will be sensitive to contrast along one particular axis, although in practical terms this is rarely an issue. However depth from focus is an active search method and will be slow because it takes time to change the focusing parameters of the camera, using for example a servo-controlled focusing ring. For this reason this method has not been applied to AMRs. A variation of the depth from focus technique has been applied to a AMR, demonstrating obstacle avoidance in a variety of environments as well as avoidance of concave obstacles such as steps and ledges [95]. This robot uses three monochrome cameras placed as close together as possible with different, fixed lens focus positions (Fig. 4.21).

Several times each second, all three frame-synchronised cameras simultaneously capture three images of the same scene. The images are each divided into five columns and three rows, or 15 subregions. The approximate sharpness of each region is computed using a variation of Equation (4.22), leading to a total of 45 sharpness values. Note that Equation 22 calculates sharpness along diagonals but skips one row. This is due to a subtle but important issue. Many cameras produce images in interlaced mode. This means that the odd rows are captured first, then afterwards the even rows are captured. When such a camera is used in dynamic environments, for example on a moving robot, then adjacent rows show the dynamic scene at two different time points, differing by up to 1/30 seconds. The result is an artificial blurring due to motion and not optical defocus. By comparing only even-number rows we avoid this interlacing side effect.

Recall that the three images are each taken with a camera using a different focus position. Based on the focusing position, we call each image close, medium or far. A 5x3 coarse depth map of the scene is constructed quickly by simply comparing the sharpness values of each three corresponding regions. Thus, the depth map assigns only two bits of depth information to each region using the values close, medium and far. The critical step is to adjust the focus positions of all three cameras so that flat ground in front of the obstacle results in medium readings in one row of the depth map. Then, unexpected readings of either close or far will indicate convex and concave obstacles respectively, enabling basic obstacle avoidance in the vicinity of objects on the ground as well as drop-offs into the ground. Although sufficient for obstacle avoidance, the above depth from focus algorithm presents unsatisfyingly coarse range information. The alternative is depth from defocus, the most desirable of the focus-based vision techniques. Depth from defocus methods take as input two or more images of the same scene, taken with different, known camera geometry. Given the images and the camera geometry settings, the goal is to recover the depth information of the three-dimensional scene represented by the images. We begin by deriving the relationship between the actual scene properties (irradiance and depth), camera geometry settings and the image g that is formed at the image plane. The focused image $f(x,y)$ of a scene is defined as follows. Consider a pinhole aperture ($L=0$) in lieu of the lens. For every point p at position (x,y) on the image plane, draw a line through the pinhole aperture to the corresponding, visible point P in the actual scene. We define $f(x,y)$ as the irradiance (or light intensity) at p due to the light from P . Intuitively, $f(x,y)$ represents the intensity image of the scene perfectly in focus

2.4 Feature Extraction

An AMR must be able to determine its relationship to the environment by making measurements with its sensors and then using those measured signals. A wide variety of sensing technologies are available, as we discussed previously. But every sensor we have presented is imperfect:

measurements always have error and, therefore, uncertainty associated with them.

Therefore, sensor inputs must be used in a way that enables the robot to interact with its environment successfully in spite of measurement uncertainty. There are two (2) strategies for using uncertain sensor input to guide the robot's behavior. One strategy is to use each sensor measurement as a raw and individual value. Such raw sensor values could for example be tied directly to robot behavior, whereby the robot's actions are a function of its sensor inputs. Alternatively, the raw sensors values could be used to update an intermediate model, with the robot's actions being triggered as a function of this model rather than the individual sensor measurements.

The second strategy is to extract information from one or more sensor readings first, generating a higher-level percept that can then be used to inform the robot's model and perhaps the robot's actions directly. We call this process feature extraction, and it is this next, optional step in the perceptual interpretation pipeline (Fig. 4.34) that we will now discuss.

In practical terms, mobile robots do not necessarily use feature extraction and scene interpretation for every activity. Instead, robots will interpret sensors to varying degrees depending on each specific functionality. For example, in order to guarantee emergency stops in the face of immediate obstacles, the robot may make direct use of raw forward-facing range readings to stop its drive motors. For local obstacle avoidance, raw ranging sensor strikes may be combined in an occupancy grid model, enabling smooth avoidance of obstacles meters away. For map-building and precise navigation, the range sensor values and even vision sensor measurements may pass through the complete perceptual pipeline, being subjected to feature extraction followed by scene interpretation to minimize the impact of individual sensor uncertainty on the robustness of the robot's map-making and navigation skills. The pattern that thus emerges is that, as one moves into more sophisticated, long-term perceptual tasks, the feature extraction and scene interpretation aspects of the perceptual pipeline become essential.

2.4.1 Defining Feature

Features are recognizable structures of elements in the environment. They usually can be extracted from measurements and mathematically described. Good features are always perceivable and easily detectable from the environment. We distinguish between low-level features (geometric primitives) like lines, circles or polygons and high-level features (objects) such as edges, doors, tables or a trash can. At one extreme, raw sensor data provides a large volume of data, but with low distinctiveness of each individual quantum of data. Making use of raw data has the potential advantage that every bit of information is fully used, and thus there is a high conservation of information. Low level

features are abstractions of raw data, and as such provide a lower volume of data while increasing the distinctiveness of each feature. The hope, when one incorporates low level features, is that the features are filtering out poor or useless data, but of course it is also likely that some valid information will be lost as a result of the feature extraction process. High level features provide maximum abstraction from the raw data, thereby reducing the volume of data as much as possible while providing highly distinctive resulting features. Once again, the abstraction process has the risk of filtering away important information, potentially lowering data utilization.

Although features must have some spatial locality, their geometric extent can range widely. For example, a corner feature inhabits a specific coordinate location in the geometric world. In contrast, a visual "fingerprint" identifying a specific room in an office building applies to the entire room, but has a location that is spatially limited to the one, particular room. In mobile robotics, features play an especially important role in the creation of environmental models. They enable more compact and robust descriptions of the environment, helping a mobile robot during both map-building and localization. When designing a mobile robot, a critical decision revolves around choosing the appropriate features for the robot to use. A number of factors are essential to this decision:

Target Environment For geometric features to be useful, the target geometries must be readily detected in the actual environment. For example, line features are extremely useful in office building environments due to the abundance of straight walls segments while the same feature is virtually useless when navigating Mars.

Available Sensors Obviously the specific sensors and sensor uncertainty of the robot impacts the appropriateness of various features. Armed with a laser rangefinder, a robot is well qualified to use geometrically detailed features such as corner features due to the high quality angular and depth resolution of the laser scanner. In contrast, a sonar-equipped robot may not have the appropriate tools for corner feature extraction.

Computational Power Vision-based feature extraction can effect a significant computational cost, particularly in robots where the vision sensor processing is performed by one of the robot's main processors.

Environment representation Feature extraction is an important step toward scene interpretation, and by this token the features extracted must provide information that is consonant with the representation used for the environment model. For example, non-geometric vision-based features are of little value in purely geometric environment models but can be of great value in topological models of the environment. Figure 4.35 shows the application of two different representations to the task of modeling an office building hallway. Each approach has advantages and disadvantages, but extraction of line and corner features has much more relevance to the representation on the left. Refer to Chapter 5, Section 5.5 for a close look at map representations and their relative tradeoffs. In the following two sections, we present specific feature extraction techniques based on the two most popular sensing modalities of mobile robotics: range sensing and visual appearance-based sensing.

2.4.2 Using Range Data

Most of today's features extracted from ranging sensors are geometric primitives such as line segments or circles. The main reason for this is that for most other geometric primitives the parametric description of the features becomes too complex and no closed form solution exists. Here we will describe line extraction in detail, demonstrating how the uncertainty models presented above can be applied to the problem of combining multiple sensor measurements. Afterwards, we briefly present another very successful feature for indoor mobile robots, the corner feature, and demonstrate how these features can be combined in a single representation.

Line Extraction

Geometric feature extraction is usually the process of comparing and matching measured sensor data against a predefined description, or template, of the expected feature. Usually, the system is overdetermined in that the number of sensor measurements exceeds the number of feature parameters to be estimated. Since the sensor measurements all have some error, there is no perfectly consistent solution and, instead, the problem is one of optimization. One can, for example, extract the feature that minimizes the discrepancy with all sensor measurements used (e.g. least squares estimation). In this section we present an optimization-based solution to the problem of extracting a line feature from a set of uncertain sensor measurements. For greater detail than is presented below, refer to [19], pp. 15 and 221.

Probabilistic Line Extraction

4.36. There is uncertainty associated with each of the noisy range sensor measurements, and so there is no single line that passes through the set. Instead, we wish to select the best possible match, given some optimization criterion. More formally, suppose n ranging measurement points in polar coordinates $x = (\rho, \theta)$ are produced by the robot's sensors. We know that there is uncertainty associated with each measurement, and so we can model each measurement using two random variables $X = (P, Q)$. In this analysis we assume that uncertainty with respect to the actual value θ of P and Q are independent. Based on Equation (4.56) we can state this formally: Furthermore, we will assume that each random variable is subject to a Gaussian probability density curve, with a mean at the true value and with some specified variance: Given some measurement point (ρ, θ) , we can calculate the corresponding Euclidean coordinates $x = (\cos \theta, \sin \theta)$. If there were no error, we would want to find a line for which all measurements lie on that line: Of course there is measurement error, and so this quantity will not be zero. When it is non-zero, this is a measure of the error between the measurement point (ρ, θ) and the line, specifically in terms of the minimum orthogonal distance between the point and the line. It is always important to understand how the error that shall be minimized is being measured. For example a number of line extraction techniques do not minimize this orthogonal point-line distance, but instead the distance parallel to the y -axis between the point and the line. A good illustration of the variety of

optimization criteria is available in [18] where several algorithms for fitting circles and ellipses are presented which minimize algebraic and geo-metric distances. For each specific (x_i, y_i) , we can write the orthogonal distance d between (x_i, y_i) and ℓ the line as:

Chapter 3

Theory of Probability

Table of Contents

3.1	Introduction	65
3.2	Experiments & Outcomes	69
3.3	Probability	70
3.4	Permutations & Combinations	75
3.5	Random Variables and Probability Distributions	79
3.6	Mean and Variance of a Distribution	83
3.7	Binomial, Poisson, and Hyper-geometric Distributions	86
3.8	Normal Distribution	90
3.9	Distribution of Several Random Variables	93

3.1 Introduction

When the data we are working are influenced by “**chance**”, by factors whose effect we cannot predict exactly¹, we have to rely on **probability theory**. The application of this theory nowadays appears in numerous fields such as from studying a game of cards to the global financial market and allow us to model processes of chance called **random experiments**.

¹This could be weather data, stock prices, life spans or ties, etc.

In such an experiment we observe a **random variable** X , that is, a function whose values in a trial² occur “by chance” according to a **probability distribution** which gives the individual probabilities, which possible values of X may occur in the long run.

²a performance of an experiment.

i.e., each of the six faces of a die should occur with the same probability, $1/6$.

Or we may simultaneously observe more than one random variable, for instance, height and weight of persons or hardness and tensile strength of steel. But enough about spoiling all the fun and let's

begin with looking at data.

Representing Data

Data can be represented numerically or graphically in different ways

i.e., a news website may contain tables of stock prices and currency exchange rates, curves or bar charts illustrating economical or political developments, or pie charts showing how inflation is calculated.

And there are numerous other representations of data for special purposes. In this section, we will discuss the use of standard representations of data in statistics³.

³There are various software dedicated to analyse and visualise statistical data. Some of these include: R, a statistical programming language, Python, MATLAB, ...

Exercise 3.1: Recording Data

Sample values, such as observations and measurements, should be recorded in the order in which they occur. Sorting, that is, ordering the sample values by size, is done as a first step of investigating properties of the sample and graphing it.

As an example let's look at super alloys.

Super alloys is a collective name for alloys used in jet engines and rocket motors, requiring high temperature (typically 1000° C), high strength, and excellent resistance to oxidation.

Thirty (30) specimens of Hastelloy C (nickel-based steel, investment cast) had the tensile strength (in 1000 lb>sq in.), recorded in the order obtained and rounded to integer values.

$$\begin{array}{cccccccccccccccccccc} 89 & 77 & 88 & 91 & 88 & 93 & 99 & 79 & 87 & 84 & 86 & 82 & 88 & 89 & 78 \\ 90 & 91 & 81 & 90 & 83 & 83 & 92 & 87 & 89 & 86 & 89 & 81 & 87 & 84 & 89 \end{array} \quad (3.1)$$

Of course depending on the need the data needs to be sorted which is shown below:

$$\begin{array}{cccccccccccccccccccc} 77 & 78 & 79 & 81 & 81 & 82 & 83 & 83 & 84 & 84 & 86 & 86 & 87 & 87 & 87 \\ 88 & 88 & 88 & 89 & 89 & 89 & 89 & 89 & 90 & 90 & 91 & 91 & 92 & 93 & 99 \end{array}$$

Graphic Representation of Data

Let's now use the data we have seen in Example 1 and see the methods we can use for graphic representations.

Exercise 3.2: Leaf Plots

One of the simplest yet most useful representations of data [49]. For Eq. (3.1) it is shown in Table ??.

The numbers in Eq. (3.1) range from 78 to 99; which you can also see this in the sorted list. To visualise this data feature, we divide these numbers into five (5) groups:

75-79, 80-84, 85-89, 90-94, 95-99.

The integers in the tens position of the groups are 7, 8, 8,

9, 9. These form the stem which can be seen in Table ??.
The first leaf is 789, representing 77, 78, 79. The second leaf is 1123344, representing 81, 81, 82, 83, 83, 84, 84. And so on. The number of times a value occurs is called its **absolute frequency**.

Therefore in this example, 78 has absolute frequency 1, the value 89 has absolute frequency 5, etc. ■

Exercise 3.3: Histogram

For large sets of data, histograms are better in displaying the distribution of data than stem-and-leaf plots. The principle is explained in Fig. 3.1.

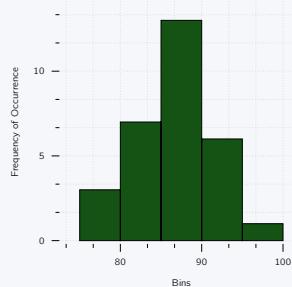


Figure 3.1: The histogram of the data given in Exercise 1.

The bases of the rectangles in seen in Fig. 3.1 are the x-intervals⁴ where there rage is:

$$74.5 - 79.5, \quad 79.5 - 84.5, \quad 84.5 - 89.5, \\ 89.5 - 94.5, \quad 94.5 - 99.5,$$

whose midpoints, known as **class marks**, are

$$x = 77, 82, 87, 92, 97,$$

respectively. The height of a rectangle with class mark x is the relative class frequency $f_{\text{rel}}(x)$, defined as the number of data values in that class interval, divided by n ($= 30$ in our case). Hence the areas of the rectangles are proportional to these relative frequencies,

$$0.10, 0.23, 0.43, 0.17, 0.07,$$

so that histograms give a good impression of the distribution of data.

⁴known as class intervals.

Mean, Standard Deviation, and Variance

Medians and quartiles are easily obtained by ordering and counting⁵.

⁵This can be done without the need of calculators.

However this method does not give full information on data as you can change data values to some extent without changing the median.

The average size of the data values can be measured in a more refined way by the mean:

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j = \frac{1}{n} (x_1 + x_2 + \dots + x_n). \quad (3.2)$$

This is the **arithmetic mean** of the data values, obtained by taking their sum and dividing by the data size (n). Therefore the arithmetic mean for Eq. (3.1) is:

$$\bar{x} = \frac{1}{30} (89 + 77 + \dots + 89) = \frac{260}{3} \approx 86.7 \quad \blacksquare$$

As we can see every data value contributes, and changing one of them will change the mean. Similarly, the spread⁶ of the data values can be measured in a more refined way by the **standard deviation** s or by its square, the **variance**⁷

$$s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2 = \frac{1}{n-1} [(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2] \quad (3.3)$$

⁶also known as variability.

⁷The symbol for variance is interesting as each domain have their own definition, as s^2 , σ^2 and $\text{Var}()$ are all acceptable symbols.

Therefore, to obtain the variance of the data, take the difference (i.e., $x_j - \bar{x}$) of each data value from the mean, square it, take the sum of these n squares, and divide it by $n - 1$.

To get the standard deviation s , take the square root of s^2 .

⁸which we calculated previously Returning back to our super alloy example, using $\bar{x} = 260/3^8$, we get for the data given in Eq. (3.1) the variance:

$$s^2 = \frac{1}{29} \left[\left(89 - \frac{260}{3} \right)^2 + \left(77 - \frac{260}{3} \right)^2 + \dots + \left(89 - \frac{260}{3} \right)^2 \right] = \frac{2006}{87} \approx 23.06 \blacksquare$$

Therefore, the standard deviation is calculated to be:

$$s = \sqrt{2006/87} \approx 4.802$$

The standard deviation has the same dimension as the data values, which is an advantage, whereas, the variance is preferable to the standard deviation in developing statistical methods.

Empirical Rule

For any round-shaped symmetric distribution of data the intervals:

$$\bar{x} \pm s, \quad \bar{x} \pm 2s, \quad \bar{x} \pm 3s, \quad \text{contain about } 68\%, \quad 95\%, \quad 99.7\%.$$

respectively, of the data points. This information is quite useful in doing quick calculation of statistical properties such as the quality of production which will be the focus in Chapter 4.

Exercise 3.4: Empirical Rule Outliers and z-Score

For the data set given in Example 1.1, with $\bar{x} = 86.7$ and $s = 4.8$, the three (3) intervals in the Rule are:

$$81.9 \leq x \leq 91.5, \quad 77.1 \leq x \leq 96.3, \quad 72.3 \leq x \leq 101.1$$

and contain 73% (22 values remain, 5 are too small, and 5 too large), 93% (28 values, 1 too small, and 1 too large), and 100%, respectively.

If we reduce the sample by omitting the outlier value of 99, mean and standard deviation reduce to $\bar{x}_{\text{red}} = 86.2$, and $s_{\text{red}} = 4.3$, approximately, and the percentage values become 67% (5 and 5 values outside), 93% (1 and 1 outside), and 100%.

Finally, the relative position of a value x in a set of mean \bar{x} and standard deviation s can be measured by the **z-score**:

$$z(s) = \frac{x - \bar{x}}{s}$$

This is the distance of x from the mean \bar{x} measured in multiples of s . For instance:

$$z(s) = \frac{(83 - 86.7)}{4.8} = -0.77$$

This is negative because 83 lies below the mean. By the empirical rule, the extreme z-values are about -3 and 3. \blacksquare

3.2 Experiments & Outcomes

Now we have the basis covered, it is time to look at probability theory⁹. This theory has the purpose of providing mathematical models of situations affected or even governed by **change effects**, for instance, in weather forecasting, life insurance, quality of technical products (computers, batteries, steel sheets, etc.), traffic problems, and, of course, games of chance with cards or dice, and the accuracy of these models can be tested by suitable observations or experiments.

⁹Sometimes known as probability calculus.

Let's start by defining some standard terms:

experiment A process of measurement or observation, in a laboratory, in a factory, ...

randomness Situation where absolute prediction is not possible.

trial A single performance of an experiment

outcome The result of a trial¹⁰

¹⁰also known as sample point.

sample space Defined as S , is the set of all possible outcomes of an experiment.

Exercise 3.5: Sample Spaces of Random Experiments & Events

- Inspecting a lightbulb | $S = \{\text{Defective, Non-defective}\}$.
- Rolling a die | $S = \{1, 2, 3, 4, 5, 6\}$
 - events are
 - $A = 1, 3, 5$ ("Odd number")
 - $B = 2, 4, 6$ ("Even number"), etc.
- Counting daily traffic accidents in Vienna | $S = \{\text{the integers in some interval}\}$.

3.3 Probability

The **probability** of an event A in an experiment is to measure **how frequently** A is roughly to occur if we make many trials. If we flip a coin, then heads H and tails T will appear **about** equally¹¹ often.

¹¹on the condition, the measurements are done for a long time.

we say that H and T are "**equally likely**."

¹²called a fair dice Similarly, for a regularly shaped die of homogeneous material¹² each of the six (6) outcomes $1, \dots, 6$ will be equally likely. These are examples of experiments in which the sample space S consists of finitely many outcomes (points) that for reasons of some symmetry can be regarded as equally likely.

Let's formulate this in a theory.

Theory 3.1: First Definition of Probability

If the sample space S of an experiment consists of **finitely** many outcomes (points) being equally likely, the probability $P(A)$ of an event A is defined to be:

$$P(A) = \frac{\text{Number of points in } A}{\text{Number of points in } S}.$$

From this definition it follows immediately, in particular, the probability of all events occurring in the sample space S is:

$$P(S) = 1.$$

Exercise 3.6: Fair Die

In rolling a fair die once:

1. What is the probability $P(A)$ of A of obtaining a 5 or a 6?
2. The probability of B : "Even number"?

Solution

The six outcomes are equally likely, so that each has probability $1/6$. Therefore:

$$P(A) = \frac{2}{6} = \frac{1}{3} \quad \text{and} \quad P(B) = \frac{3}{6} = \frac{1}{2} \blacksquare$$

The above theory takes care of many games as well as some practical applications, but not of all experiments, as in many problems we do not have finitely many equally likely outcomes. To arrive at a more general definition of probability, we regard probability as the counterpart of **relative frequency**:

$$f_{\text{rel}}(A) = \frac{f(A)}{n} = \frac{\text{Number of times } A \text{ occurs}}{\text{Number of trials}} \quad (3.4)$$

Now if A did not occur, then $f(A) = 0$. If A always occurred, then $f(A) = n$. These are of course extreme cases. Division by n gives:

$$0 \leq f_{\text{rel}}(A) \leq 1 \quad (3.5)$$

¹³meaning that some event always occurs

In particular, for $A = S$ we have $f(S) = n$ as S always occurs¹³. Division by n gives:

$$f_{\text{rel}}(S) = 1 \quad (3.6)$$

Finally, if A and B are **mutually exclusive**, they cannot occur together. Therefore the absolute frequency of their union $A = B$ must equal the sum of the absolute frequencies of A and B . Division

by n gives the same relation for the relative frequencies:

$$f_{\text{rel}}(A \cup B) = f_{\text{rel}}(A) + f_{\text{rel}}(B) \quad (3.7)$$

We can now extend the definition of probability to experiments in which equally likely outcomes are not available.

Theory 3.2: General Definition of Probability

Given a sample space S , with each event A of S (A being a subset of S) there is associated a number $P(A)$, called the **probability** of A , such the following **axioms of probability** are satisfied.

- For every A in S ,

$$0 \leq P(A) \leq 1. \quad (3.8)$$

- The entire sample space S has the probability

$$P(S) = 1. \quad (3.9)$$

- For **mutually exclusive** events A and B :

$$P(A \cup B) = P(A) + P(B) \quad (A \cap B = \emptyset). \quad (3.10)$$

- If S is **infinite**¹⁴, the previous statement has to be replaced by Eq. (3.4), where for mutually exclusive events A_1, A_2, \dots ,

$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots \quad (3.11)$$

¹⁴i.e., has infinitely many points.

In the infinite case the subsets of S on which $P(A)$ is defined are restricted to form a so-called σ -algebra.

Basic Theorems of Probability

We will see that the axioms of probability will enable us to build up probability theory and its application to statistics. We begin with three (3) basic theorems. The first one is useful if we can get the probability of the complement A^c more easily than $P(A)$ itself.

Theory 3.3: Complementation Rule

For an event A and its complement A^c in a sample space S ,

$$P(A^c) = 1 - P(A) \quad (3.12)$$

Exercise 3.7: Coin Tossing

Five (5) coins are tossed simultaneously.

Find the probability of the event A :

At least one head turns up. Assume that the coins are fair.

Solution

As each coin can turn up either heads or tails, the sample space consists of $2^5 = 32$ outcomes. Given the coins are fair, we may assign the same probability ($1/32$) to each outcome. Then the event A^c (No heads turn up) consists of only 1 outcome. Hence $P(A^c) = 1/32$, and the answer is:

$$P(A) = 1 - P(A^c) = \frac{31}{32} \quad \blacksquare$$

Theory 3.4: Addition Rule for Mutually Exclusive Events

For **mutually exclusive events** A_1, \dots, A_m in a sample space S ,

$$P(A_1 \cup A_2 \cup \dots \cup A_m) = P(A_1) + P(A_2) + \dots + P(A_m). \quad (3.13)$$

Exercise 3.8: Mutually Exclusive Events

If the probability that on any workday a garage will get 10-20, 21-30, 31-40, over 40 cars to service is 0.20, 0.35, 0.25, 0.12, respectively, what is the probability that on a given workday the garage gets at least 21 cars to service?

Solution

As these are mutually exclusive events, the answer is:

$$0.35 + 0.25 + 0.12 = 0.72 \blacksquare$$

However, most situations, events will **NOT** be mutually exclusive. Then we have the following theorem to formalise the previous statement.

Theory 3.5: Addition Rule for Arbitrary Events

For events A and B in a sample space, their union is defined as:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B). \quad (3.14)$$

For **mutually exclusive** events A and B we have $A \cap B = \emptyset$ by definition:

$$P(\emptyset) = 0 \quad (3.15)$$

Exercise 3.9: Union of Arbitrary Events

In tossing a fair die, what is the probability of getting an odd number or a number less than 4?

Solution

Let A be the event "Odd number" and B the event "Number less than 4." As these events are linked we can write:

$$P(A \cup B) = \frac{3}{6} + \frac{3}{6} - \frac{2}{6} = \frac{2}{3}$$

as $A \cup B = \text{Odd number less than 4} = \{1, 3\}$ ■

Conditional Probability and Independent Events

It is often required to find the probability of an event B given the condition of an event A occurs. This probability is called the **conditional probability** of B given A and is denoted by $P(B|A)$.

In this case A serves as a new, reduced, sample space, and that probability is the fraction of $P(A)$ which corresponds to $A \cap B$. Therefore,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{where} \quad P(B) \neq 0 \quad (3.16)$$

Similarly, the conditional probability of A given B is:

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad \text{where} \quad P(A) \neq 0 \quad (3.17)$$

Theory 3.6: Multiplication Rule

Given A and B are events defined in a sample space S and $P(A) \neq 0, P(B) \neq 0$, then

$$P(A \cap B) = P(A) P(B|A) = P(B) P(A|B). \quad (3.18)$$

Exercise 3.10: Multiplication Rule

In producing screws, let:

- A mean "screw too slim",
- B mean "screw too short."

Let $P(A) = 0.1$ and let the conditional probability that a slim screw is also too short be $P(B|A) = 0.2$. What is the probability that a screw that we pick randomly from the lot produced will be both too slim and too short?

Solution

$$P(A \cap B) = P(A) P(B|A) = 0.1 \times 0.2 = 0.02 = 2\% \quad \blacksquare$$

Independent Events

If events A and B are such that

$$P(A \cap B) = P(A) P(B), \quad (3.19)$$

they are called **independent events**. Assuming $P(A) \neq 0, P(B) \neq 0$, we see from Eq. (3.16) - Eq. (3.18):

$$P(A|B) = P(A), \quad P(B|A) = P(B).$$

This means that the probability of A does not depend on the occurrence or nonoccurrence of B, and conversely. This justifies the term **independent**.

Independence of m Events

Similarly, m events A_1, \dots, A_m are called **independent** if:

$$P(A_1 \cap \dots \cap A_m) = P(A_1) \dots P(A_m) \quad (3.20)$$

as well as for every k different events $A_{j_1}, A_{j_2}, \dots, A_{j_k}$.

$$P(A_{j_1} \cap A_{j_2} \cap \dots \cap A_{j_k}) = P(A_{j_1}) P(A_{j_2}) \dots P(A_{j_k}) \quad (3.21)$$

where $k = 2, 3, \dots, m - 1$. Accordingly, three events A, B, C are independent if and only if

$$P(A \cap B) = P(A) P(B), \quad (3.22)$$

$$P(B \cap C) = P(B) P(C), \quad (3.23)$$

$$P(C \cap A) = P(C) P(A), \quad (3.24)$$

$$P(A \cap B \cap C) = P(A) P(B) P(C). \quad (3.25)$$

Sampling

Our next example has to do with randomly drawing objects, *one at a time*, from a given set of objects. This is called **sampling from a population**, and there are two ways of sampling, as follows.

■ **In sampling with replacement**, the object that was drawn at random is placed back to the given set and the set is mixed thoroughly. Then we draw the next object at random.

■ **In sampling without replacement** the object that was drawn is put aside.

Exercise 3.11: Sampling w/o Replacement

A box contains 10 screws, three (3) of which are defective. Two screws are drawn at random. Find the probability that neither of the two screws is defective.

Solution

We consider the events

- A First drawn screw non-defective,
- B Second drawn screw non-defective.

We can see:

$$P(A) = \frac{1}{10}$$

as 7 of the 10 screws are non-defective and we sample at random, so that each screw has the same probability ($\frac{1}{10}$) of being picked.

If we sample with replacement, the situation before the second drawing is the same as at the beginning, and $P(B) = \frac{7}{10}$. The events are independent, and the answer is

$$P(A \cap B) = P(A) P(B) = 0 \cdot 7 \cdot 0.7 = 0.49\%.$$

If we sample without replacement, then $P(A) = \frac{7}{10}$, as before. If A has occurred, then there are 9 screws left in the box, 3 of which are defective.

Thus $P(B|A) = \frac{6}{9} = \frac{2}{3}$, therefore:

$$P(A \cap B) = \frac{7}{10} \cdot \frac{2}{3} = 47\% \blacksquare$$

3.4 Permutations & Combinations

Permutations and combinations help in finding probabilities $P(A) = a/k$ by systematically counting the number a of points of which an event A consists.

where, k is the number of points of the sample space S .

The practical difficulty is that a may often be surprisingly large, so that actual counting becomes hopeless. For example, if in assembling some instrument you need 10 different screws in a certain order and you want to draw them randomly from a box¹⁵ the probability of obtaining them in the required order is only 1/3,628,800 because there are exactly:

$$10! = 1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7 \cdot 8 \cdot 9 \cdot 10 = 3,628,800$$

¹⁵Of course, this goes without saying, there is nothing but screws in this imaginary box.

orders in which they can be drawn. Similarly, in many other situations the numbers of orders, arrangements, etc. are often incredibly large.

3.4.1 Permutations

A **permutation** of given things¹⁶ is an arrangement of these things in a row in some order.

¹⁶such as *elements* or *objects*.

i.e., for three (3) letters a, b, c there are $3! = 1 \cdot 2 \cdot 3 = 6$ permutations: abc, acb, bca, cab, cba

Let's write this behaviour down as a theory:

Theory 3.7: Permutations

Different things

The number of permutations of n different things taken all at a time is

$$n! = 1 \cdot 2 \cdot 3, \dots, n. \quad (3.26)$$

Classes of Equal Things

If n given things can be divided into c classes of alike things differing from class to class, then the number of permutations of these things taken all at a time is

$$\frac{n!}{n_1!n_2!\cdots n_c!} \quad \text{where} \quad n_1 + n_2 + \cdots + n_c = n, \quad (3.27)$$

where n_j is the number of things in the j^{th} class.

Permutation of n things taken k at a time

A permutation containing only k of the n given things. Two such permutations consisting of the same k elements, in a different order, are different, by definition.

i.e., there are 6 different permutations of the three letters a, b, c , taken two letters at a time, ab, ac, bc, ba, ca, cb .

Permutation of n things taken k at a time with repetitions

An arrangement obtained by putting any given thing in the first position, any given thing, including a repetition of the one just used, in the second, and continuing until k positions are filled.

i.e., there are $3^2 = 9$ different such permutations of a, b, c taken 2 letters at a time, namely, the preceding 6 permutations and aa, bb, cc .

Theory 3.8: Permutations

The number of different permutations of n different things taken k at a time **without repetitions** is

$$n(n-1)(n-2)\cdots(n-k+1) = \frac{n!}{(n-k)!}, \quad (3.28)$$

and **with repetitions** is,

$$n^k. \quad (3.29)$$

Exercise 3.12: An Encrypted Message

In an encrypted message the letters are arranged in groups of five (5) letters, called words. Knowing the letter can be repeated, we see that the number of different such words is

$$26^5 = 11,881,376 \blacksquare$$

For the case of different such words containing each letter no more than once is

$$\frac{26!}{(26-5)!} = 26 \cdot 25 \cdot 24 \cdot 23 \cdot 22 = 7,893,600 \blacksquare$$

3.4.2 Combinations

In a permutation, the **order of the selected things is essential**. In contrast, a **combination** of a given things means any selection of one or more things **without regard to order**. There are two (2) kinds of combinations, as follows:

1. The number of **combinations of n different things, taken k at a time, without repetitions** is the number of sets that can be made up from the n given things, each set containing k different things and no two (2) sets containing exactly the same k things.
2. The number of **combinations of n different things, taken k at a time, with repetitions** is the number of sets that can be made up of k things chosen from the given n things, each being used as often as desired.

i.e, there are three (3) combinations of the three (3) letters a, b, c , taken two (2) letters at a time, without repetitions, namely, ab, ac, bc , and six such combinations with repetitions, namely, ab, ac, bc, ca, bb, cc .

Theory 3.9: Combinations

The number of different combinations of n different things taken, k at a time, **without repetitions**, is:

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{n(n-1)\cdots(n-k+1)}{1\cdot2\cdots k}, \quad (3.30)$$

and the number of those combinations **with repetitions** is:

$$\binom{n+k-1}{k}. \quad (3.31)$$

Exercise 3.13: Sampling Light-bulbs

The number of samples of five (5) light-bulbs that can be selected from a lot of 500 bulbs is

$$\binom{500}{5} = \frac{500!}{5!495!} = \frac{500 \cdot 499 \cdot 498 \cdot 497 \cdot 476}{1 \cdot 2 \cdot 3 \cdot 4 \cdot 5} = 255,244,687,600 \blacksquare$$

3.4.3 Factorial Function

In Eq. (3.26)-Eq. (3.31) the **factorial function** is relatively straightforward. By definition¹⁷,

$$0! = 1.$$

Values may be computed recursively from given values by

$$(n+1)! = (n+1)n!.$$

For large n the function is very large and hard to keep track of. A convenient approximation for large n is the **Stirling formula**, defined as:

$$n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \quad \text{where} \quad e = 2.718\cdots \quad (3.32)$$

where \sim is read **asymptotically equal**¹⁸ and means that the ratio of the two sides of Eq. (3.32) approaches 1 as n approaches infinity.

¹⁷This is done by convention. An intuitive way to look at it is $n!$ counts the number of ways to arrange distinct objects in a line, and there is only one way to arrange nothing.

¹⁸it means the percentage difference between the vertical distances between points on the two graphs approaches 0.

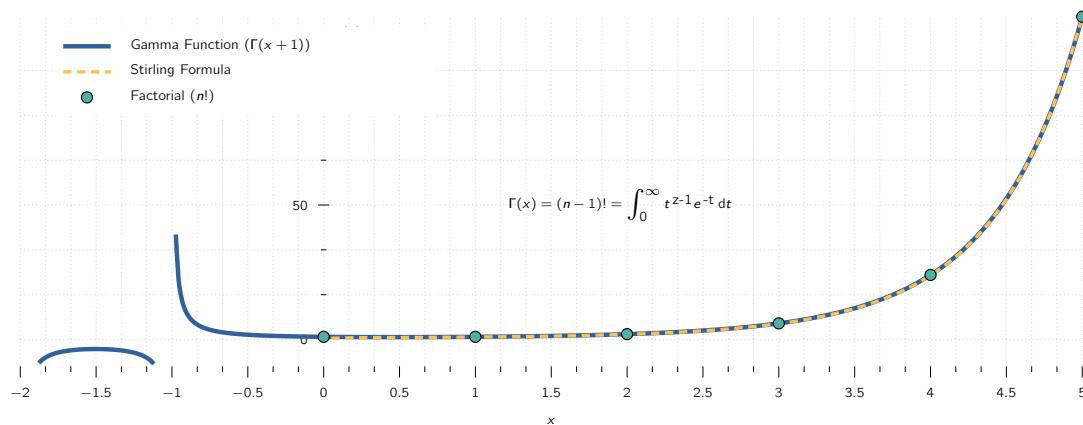


Figure 3.2: A visual comparison of the Stirling formula and the actual values of the factorial function.

3.4.4 Binomial Coefficients

The **binomial coefficients** are defined by the following formula:

$$\binom{a}{k} = \frac{(a)(a-1)(a-2)\cdots(a-k+1)}{k!} \quad \text{where } (k \geq 0, \text{ integer}) \quad (3.33)$$

The numerator has k factors. Furthermore, we define

$$\binom{a}{0} = 1, \quad \text{in particular,} \quad \binom{0}{0} = 1.$$

For integer $a = n$ we obtain from Eq. (3.33):

$$\binom{n}{k} = \binom{n}{n-k} \quad (n \geq 0 \quad \text{and} \quad 0 \leq k \leq n).$$

Binomial coefficients may be computed recursively, because

$$\binom{a}{k} + \binom{a}{k+1} = \binom{a+1}{k+1} \quad (k \geq 0, \text{ integer}).$$

Formula Eq. (3.33) also gives:

$$\binom{-m}{k} = (-1)^k \binom{m+k-1}{k} \quad \text{where} \quad k \geq 0, \text{ integer} \quad \text{and} \quad m > 0.$$

There are two (2) important relations worth mentioning:

$$\sum_{s=0}^{n-1} \binom{k+s}{k} = \binom{n+k}{k+1} \quad (k \geq 0 \quad \text{and} \quad n \geq 1)$$

and

$$\sum_{k=0}^r \binom{p}{k} \binom{q}{r-k} = \binom{p+q}{r} \quad (r \geq 0, \text{ integer}).$$

3.5 Random Variables and Probability Distributions

In the beginning of this chapter we considered frequency distributions of data¹⁹. These distributions show the **absolute** or **relative** frequency of the data values.

¹⁹Remember we did a histogram and a stem-and-leaf plot.

Similarly, a **probability distribution** or, a **distribution**, shows the probabilities of events in an experiment. The quantity we observe in an experiment will be denoted by X and called a **random variable**²⁰ as the value it will assume in the next trial depends on the **stochastic process**

²⁰or **stochastic variable** if you want to be pedantic.

i.e., if you roll a die, you get one of the numbers from 1 to 6, but you don't know which one will show up next. An example would be, $X = \text{Number a die turns up}$, which is a random variable.

If we count²¹, we have a **discrete random variable and distribution**. If we measure (electric voltage, rainfall, hardness of steel), we have a **continuous random variable and distribution**. For both cases (discrete, discontinuous), the distribution of X is determined by the **distribution function**:

$$F(x) = P(X \leq x) \quad (3.34)$$

This is the probability that in a trial, X will assume any value not exceeding x .

The terminology is unfortunately **NOT** uniform across the field as $F(x)$ is sometimes also called the **cumulative distribution function**.

For Eq. (3.34) to make sense in both the discrete and the continuous case we formulate conditions as follows.

Theory 3.10: Random Variable

A **random variable** X is a function defined on the sample space S of an experiment. Its values are real numbers. For every number a the probability:

$$P(X = a),$$

with which X assumes a is defined. Similarly, for any interval I , the probability

$$P(X \in I),$$

with which X assumes any value in I is defined²².

²²Although this definition is very general, in practice only a very small number of distributions will occur over and over again in applications.

From Eq. (3.34) we can define the fundamental formula for the probability corresponding to an interval $a < x \leq b$:

$$P(a < X \leq b) = F(b) - F(a). \quad (3.35)$$

This follows because $X \leq a$ (X assumes any value **NOT** exceeding a) and $a < X \leq b$ (X assumes any value in the interval $a < x \leq b$) are **mutually exclusive** events, so based on Eq. (3.34):

$$\begin{aligned} F(b) &= P(X \leq b) = P(X \leq a) + P(a < X \leq b) \\ &= F(a) + P(a < X \leq b) \end{aligned}$$

and subtraction of $F(a)$ on both sides gives Eq. (3.35).

3.5.1 Discrete Random Variables and Distributions

By definition, a random variable X and its distribution are **discrete** if X assumes only **finitely** many or at most countably many values x_1, x_2, x_3, \dots , called the **possible values** of X , with positive probabilities,

$$p_1 = P(X = x_1), p_2 = P(X = x_2), p_3 = P(X = x_3), \dots$$

whereas the probability $P(X \in I)$ is zero for any interval I containing no possible value. Clearly, the discrete distribution of X is also determined by the **probability function** $f(x)$ of X , defined by

$$f(x) = \begin{cases} p_j & \text{if } x = x_j \\ 0 & \text{otherwise} \end{cases} \quad \text{where } j = 1, 2, \dots, \quad (3.36)$$

From this we get the values of the **distribution function** $F(x)$ by taking sums,

$$F(x) = \sum_{x_j \leq x} f(x_j) = \sum_{x_j \leq x} p_j \quad (3.37)$$

where for any given x we sum all the probabilities p_j for which x_j is smaller than or equal to that of x . This is a **step function** with upward jumps of size p_j at the possible values x_j of X and constant in between. The two (2) useful formulas for discrete distributions are readily obtained as follows. For the probability corresponding to intervals we have from Eq. (3.35) and Eq. (3.37):

$$P(a < X \leq b) = F(b) - F(a) = \sum_{a < x_j \leq b} p_j \quad (3.38)$$

²³Be careful about $<$ and \leq as the former means it is NOT included and the latter means it is.

This is the sum of all probabilities p_j for which x_j satisfies $a < x_j \leq b$ ²³. From this and $P(S) = 1$ we obtain the following formula.

$$\sum_j p_j = 1 \quad (\text{sum of all probabilities}). \quad (3.39)$$

Exercise 3.14: Waiting Time Problem

In tossing a fair coin, let X be the Number of trials until the first head appears. Then, by independence of events we get (where H is heads, and T is tails):

$$\begin{aligned} P(X = 1) &= P(H) = \frac{1}{2} \\ P(X = 2) &= P(TH) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \\ P(X = 3) &= P(TTH) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8} \end{aligned}$$

and in general, $P(X = n) = \left(\frac{1}{2}\right)^n$, $n = 1, 2, 3, \dots$ which when all possible event are summed up will always give 1.

3.5.2 Continuous Random Variables and Distributions

Discrete random variables appear in experiments in which we count²⁴. Continuous random variables appear in experiments in which we measure (lengths of screws, voltage in a power line, etc.). By definition, a random variable X and its distribution are of *continuous type* or, briefly, **continuous**, if its distribution function $F(x)$, defined in Eq. (3.34), can be given by an integral²⁵:

$$F(x) = \int_{-\infty}^x f(v) dv \quad (3.40)$$

²⁴defectives in a production, days of sunshine in Kufstein, customers in a line, etc.

²⁵we write v as a toss-away variable because x is needed as the upper limit of the integral.

whose integrand $f(x)$, called the **density** of the distribution, is **non-negative**, and is continuous, perhaps except for finitely many x -values. Differentiation gives the relation of f to F as

$$f(x) = F'(x) \quad (3.41)$$

for every x at which $f(x)$ is continuous.

From Eq. (3.35) and Eq. (3.40) we obtain the very important formula for the probability corresponding to an interval²⁶:

²⁶This is an analog of Eq. (3.38)

$$P(a < X \leq b) = F(b) - F(a) = \int_a^b f(v) dv \quad (3.42)$$

Which can be seen visually in **Fig. 3.3**. From Eq. (3.40) and $P(S) = 1$ we also have the analogue of Eq. (3.39):

$$\int_{-\infty}^{\infty} f(v) dv = 1. \quad (3.43)$$

Continuous random variables are **simpler than discrete ones** with respect to intervals as, in the continuous case the four probabilities corresponding to $a < X \leq b$, $a < X < b$, $a \leq X \leq b$,

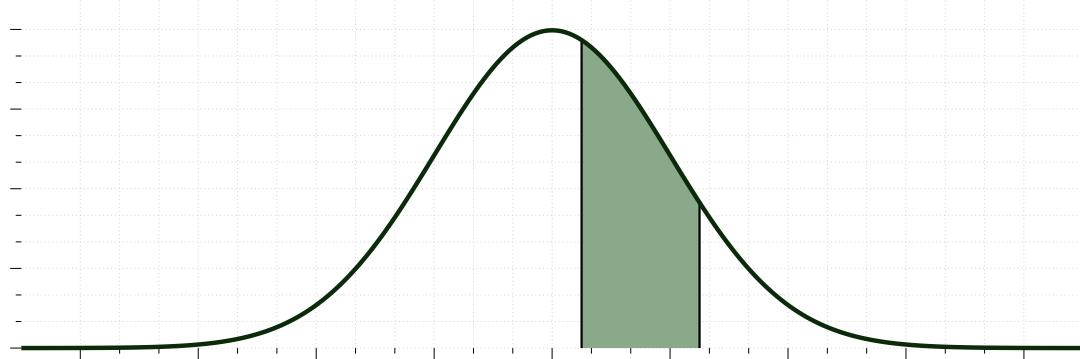


Figure 3.3: A visual representation of the Eq. (3.42).

and $a \leq X \leq b$ with any fixed a and b ($> a$) are all the same.

The next example illustrates notations and typical applications of our present formulas.

Exercise 3.15: Continuous Distribution

Let X have the density function:

$$f(x) = 0.75(1 - x^2) \quad \text{if} \quad -1 \leq x \leq 1,$$

and zero otherwise. Find:

1. The distribution function.
2. Find the probabilities $P\left(-\frac{1}{2} \leq X \leq \frac{1}{2}\right)$ and $P\left(\frac{1}{2} \leq X \leq 2\right)$
3. Find x such that $P(X \leq x) = 0.95$.

Solution

From Eq. (3.40), we obtain $F(x) = 0$ if $x \leq -1$,

$$F(x) = 0.75 \int_{-1}^x (1 - v^2) dv = 0.5 + 0.75x - 0.25x^3 \quad \text{if} \quad -1 < x \leq 1,$$

and $F(x) = 1$ if $x > 1$. From this and Eq. (3.42) we get:

$$P\left(-\frac{1}{2} \leq X \leq \frac{1}{2}\right) = F\left(\frac{1}{2}\right) - F\left(-\frac{1}{2}\right) = 0.75 \int_{-1/2}^{1/2} (1 - v^2) dv = 68.75\%$$

because $P\left(-\frac{1}{2} \leq X \leq \frac{1}{2}\right) = P\left(-\frac{1}{2} < X \leq \frac{1}{2}\right)$ for a continuous distribution we can write:

$$P\left(\frac{1}{4} \leq X \leq 2\right) = F(2) - F\left(\frac{1}{4}\right) = 0.75 \int_{1/4}^1 (1 - v^2) dv = 31.64\%.$$

Note that the upper limit of integration is 1, not 2. Finally,

$$P(X \leq x) = F(x) = 0.5 + 0.75x - 0.25x^3 = 0.95.$$

Algebraic simplification gives $3x - x^3 = 1.8$. A solution is $x = 0.73$, approximately ■

3.6 Mean and Variance of a Distribution

The mean μ and variance σ^2 of a random variable X and of its distribution are the theoretical counterparts of the mean \bar{x} and variance s^2 of a frequency distribution and serve a similar purpose.

The mean characterises the central location and the variance the spread (the variability) of the distribution. The **mean** μ is defined by:

$$(a) \quad \mu = \sum_j x_j f(x_j) \quad (\text{Discrete distribution}) \quad (3.44a)$$

$$(b) \quad \mu = \int_{-\infty}^{\infty} x f(x) dx \quad (\text{Continuous distribution}) \quad (3.44b)$$

and the **variance** σ^2 by:

$$(a) \quad \sigma^2 = \sum_j (x_j - \mu)^2 f(x_j) \quad (\text{Discrete distribution}) \quad (3.45a)$$

$$(b) \quad \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx \quad (\text{Continuous distribution}) \quad (3.45b)$$

σ (the positive square root of σ^2) is called the standard deviation²⁷ of X and its distribution. f is the probability function or the density, respectively, in (a) and (b).

²⁷Sometimes it is known as $\text{Var}(x)$

The mean μ is also denoted by $E(X)$ and is called the **expectation of X** because it gives the average value of X to be expected in many trials.

Quantities such as μ and σ^2 that measure certain properties of a distribution are called **parameters**. μ and σ^2 are the two (2) most important ones.

From Eq. (3.45a) and Eq. (3.45b), we see that²⁸:

$$\sigma^2 > 0$$

²⁸except for a discrete distribution with only one possible value.

We assume that μ and σ^2 exist²⁹, as is the case for practically all distributions that are useful in applications.

²⁹and finite.

Exercise 3.16: Mean and Variance

The random variable X , *Number of heads in a single toss of a fair coin*, has the possible values $X = 0$ and $X = 1$ with probabilities $P(X = 0) = \frac{1}{2}$ and $P(X = 1) = \frac{1}{2}$. From Eq. (3.44a) we thus obtain the mean:

$$\mu = 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{1}{2},$$

and Eq. (3.45a) gives the variance:

$$\sigma^2 = (0 - \frac{1}{2})^2 \cdot \frac{1}{2} + (1 - \frac{1}{2})^2 \cdot \frac{1}{2} = \frac{1}{4} \blacksquare$$

Symmetry

We can obtain the mean μ without calculation if a distribution is symmetric. Indeed, we can write:

Theory 3.11: Mean of a Symmetric Distribution

If a distribution is **symmetric** with respect to $x = c$, that is,

$$f(c - x) = f(c + x)$$

then $\mu = c$.

Transformation of Mean and Variance

Given a random variable X with mean μ and variance σ^2 , we want to calculate the mean and variance of $X^* = a_1 + a_2X$, where a_1 and a_2 are given constants.

This problem is important in statistics, where it often appears.

Theory 3.12: Transformation of Mean and Variance

If a random variable X has mean μ and variance σ^2 , then the random variable:

$$X^* = a_1 + a_2X \quad \text{where} \quad a_2 > 0$$

has the mean μ^* and variance σ^{*2} , where

$$\mu^* = a_1 + a_2\mu \quad \text{and} \quad \sigma^{*2} = a_2^2\sigma^2.$$

In particular, the **standardised random variable** Z corresponding to X , given by:

$$Z = \frac{X - \mu}{\sigma}$$

has the mean 0 and the variance 1.

Expectation & Moments

³⁰the value of X to be expected on the average If we recall, Eq. (3.44a) and Eq. (3.44b) define the mean of X ³⁰, written $\mu = E(X)$. More generally, if $g(x)$ is **non-constant** and continuous for all x , then $g(X)$ is a random variable. Therefore its **mathematical expectation** or, briefly, its expectation $E(g(X))$ is the value of $g(X)$ to be expected on the average, defined by:

$$E(g(X)) = \sum_j g(x_j) f(x_j) \quad \text{or} \quad E(g(X)) = \int_{-\infty}^{\infty} g(x) f(x) dx$$

In the formula on the Left Hand Side (LHS), f is the probability function of the discrete random variable X . In the formula on the Right Hand Side (RHS), f is the density of the continuous random variable X . Important special cases are the k^{th} of X (where $k = 1, 2, \dots$)

$$E(X^k) = \sum_j x_j^k f(x_j) \quad \text{or} \quad \int_{-\infty}^{\infty} x^k f(x) dx$$

and the k^{th} of X ($k = 1, 2, \dots$)

$$E([X - \mu]^k) = \sum_j (x_j - \mu)^k f(x_j) \quad \text{or} \quad \int_{-\infty}^{\infty} (x - \mu)^k f(x) dx.$$

This includes the first moment, the **mean** of X

$$\mu = E(X) \quad \text{where} \quad k = 1 \quad (3.46)$$

It also includes the second central moment, the **variance** of X

$$\sigma^2 = E((X - \mu)^2) \quad \text{where} \quad k = 2. \quad (3.47)$$

3.7 Binomial, Poisson, and Hyper-geometric Distributions

These are the three (3) most important **discrete** distributions, with numerous applications therefore are worth of a bit of a detailed look.

Of course these are not the only distributions present. There are as many distributions as there are problems with some distributions used in wide variety of fields (Gaussian) whereas some are used only in a very narrow field (Nakagami).

Binomial Distribution

The **binomial distribution** occurs in problems involving chance³¹.

What we are interested is in the number of times an event A occurs in n **independent** trials. In each trial, the event A has the same probability $P(A) = p$. Then in a trial, A will **NOT** occur with probability $q = 1 - p$. In n trials the random variable that interests us is:

$$X = \text{Number of times the event } A \text{ occurs in } n \text{ trials.} \quad (3.48)$$

X can assume the values $0, 1, \dots, n$, and we want to determine the corresponding probabilities. Now $X = x$ means that A occurs in x trials and in $n - x$ trials it does not occur. We can write this down as follows:

$$\underbrace{A \ A \ \dots A}_{x \text{ times}} \quad \text{and} \quad \underbrace{B \ B \ \dots B}_{n - x \text{ times}} \quad (3.49)$$

Here $B = A^c$ is the complement of A , meaning that A does not occur. We now use the assumption that the trials are independent³². Hence Eq. (3.49) has the probability:

$$\underbrace{p \ p \ \dots p}_{x \text{ times}} \cdot \underbrace{q \ q \ \dots q}_{n - x \text{ times}} = p^x q^{n-x} \quad (3.50)$$

Now Eq. (3.49) is just one order of arranging xA 's and $n - xB$'s. We will now calculate the number of permutations of n things³³ consisting of two (2) classes;

³³the n outcomes of the n trials

1. class 1 containing the $n_1 = x$ A 's
2. class 2 containing the $n - n_1 = n - x$ B 's

This number is:

$$\frac{n!}{x!(n - x)!} = \binom{n}{x}. \quad (3.51)$$

Accordingly, Eq. (3.50), multiplied by this binomial coefficient, gives the probability $P(X = x)$ of $X = x$, that is, of obtaining A precisely x times in n trials. Hence X has the probability function:

$$f(x) = \binom{n}{x} p^x q^{n-x} \quad (x = 0, 1, \dots, n) \quad (3.52)$$

and $f(x) = 0$ otherwise. The distribution of X with probability function (2) is called the **binomial distribution** or *Bernoulli distribution*. The occurrence of A is called *success*³⁴ and the non-occurrence of A is called *failure*.

The mean of the binomial distribution is:

$$\mu = np$$

and the variance is:

$$\sigma^2 = npq.$$

For the *symmetric case* of equal chance of success and failure ($p = q = \frac{1}{2}$) this gives the mean $n/2$, the variance $n/4$, and the probability function

$$f(x) = \binom{n}{x} \left(\frac{1}{2}\right)^x \quad (x = 0, 1, \dots, n).$$

³⁴regardless of what it actually is; it may mean that you miss your plane or lose your watch

Exercise 3.17: Binomial Distribution

Calculate the probability of obtaining at least two (2) "six" in rolling a fair die 4 times.

Solution

$p = P(A) = P(\text{six}) = \frac{1}{6}$, $q = \frac{5}{6}$, $n = 4$. The event "At least two (2) "six" occurs if we obtain 2 or 3 or 4 "six". Hence the answer is:

$$\begin{aligned} P &= f(2) + f(3) + f(4) = \binom{4}{2} \left(\frac{1}{6}\right)^2 \left(\frac{5}{6}\right)^2 + \binom{4}{3} \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right) + \binom{4}{4} \left(\frac{1}{6}\right)^4 \\ &= \frac{1}{6^4} (6 \cdot 25 + 4 \cdot 5 + 1) = \frac{171}{1296} = 13.2\%. \end{aligned}$$

Poisson Distribution

The discrete distribution with infinitely many possible values and probability function:

$$f(x) = \frac{\mu^x}{x!} e^{-\mu} \quad \text{where} \quad x = 0, 1, \dots \quad (3.53)$$

is called the **Poisson distribution**, named after *S. D. Poisson*. **Fig. 3.4** shows Eq. (3.53) for some values of μ ³⁵.

³⁵While μ is used here, some textbook use λ

It can be proved that this distribution is obtained as a limiting case of the binomial distribution, if we let $p \rightarrow 0$ and $n \rightarrow \infty$ so that the mean $\mu = np$ approaches a finite value. The Poisson distribution has the mean μ and the variance:

$$\sigma^2 = \mu. \quad (3.54)$$

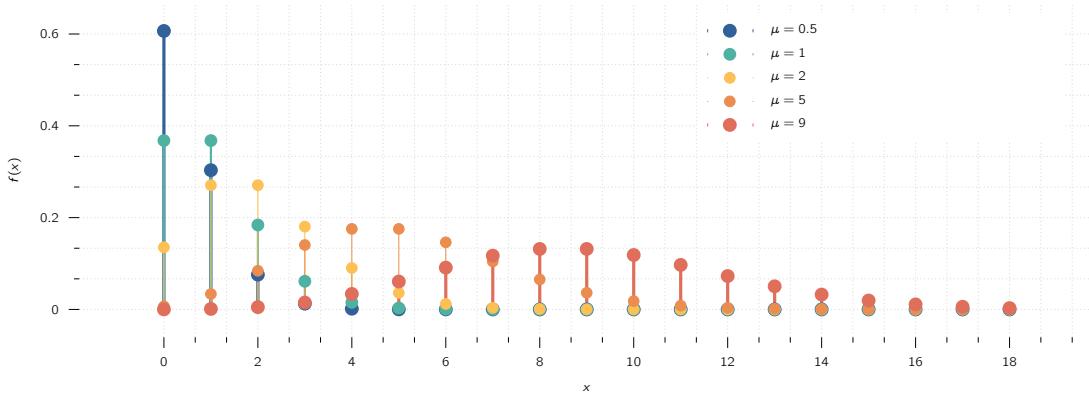
Figure 3.4: The Poisson distribution with different mean (μ) values.

Fig. 3.4 gives the impression that, with increasing mean, the spread of the distribution increases, thereby illustrating formula Eq. (3.54), and that the distribution becomes more and more symmetric³⁶

³⁶approximately

Exercise 3.18: Poisson Distribution

If the probability of producing a defective screw is $p = 0.01$, what is the probability that a lot of 100 screws will contain more than 2 defectives?

Solution

The complementary event is A^c . No more than 2 defectives. For its probability we get, from the binomial distribution with mean $\mu = np = 1$, the value.

$$P(A^c) = \binom{100}{0} 0.99^{100} + \binom{100}{1} 0.01 \cdot 0.99^{99} + \binom{100}{2} 0.01^2 \cdot 0.99^{98}.$$

Since p is very small, we can approximate this by the much more convenient Poisson distribution with mean $\mu = np = 100 \cdot 0.01 = 1$, obtaining.

$$P(A^c) = e^{-1} \left(1 + 1 + \frac{1}{2} \right) = 91.97\%.$$

Thus $P(A) = 8.03\%$. Show that the binomial distribution gives $P(A) = 7.94\%$, so that the Poisson approximation is quite good ■

Exercise 3.19: The Parking Problem

If on the average, 2 cars enter a certain parking lot per minute, what is the probability that during any given minute four (4) or more cars will enter the lot?

Solution

To understand that the Poisson distribution is a model of the situation, we imagine the minute to be divided into very many short time intervals. Let p be the (constant) probability that a car will enter the lot during any such short interval, and assume independence of the events that happen during those intervals. Then, we are dealing with a binomial distribution with very large n and very small p , which we can approximate by the Poisson distribution with

$$\mu = np = 2$$

because 2 cars enter on the average, the complementary event of the event "4 cars or more during a given minute" is "3 cars or fewer enter the lot" and has the probability

$$f(0) + f(1) + f(2) + f(3) = e^{-2} \left(\frac{2^0}{0!} + \frac{2^1}{1!} + \frac{2^2}{2!} + \frac{2^3}{3!} \right) = 0.857.$$

Which means the result is 14.3% ■

3.7.1 Sampling with Replacement

This means that we draw things from a given set one by one, and after each trial we replace the thing drawn³⁷ before we draw the next thing. This guarantees **independence of trials** and leads to the **binomial distribution**. Indeed, if a box contains N things, for example, screws, M of which are defective, the probability of drawing a defective screw in a trial is $p = M/N$. Hence the probability of drawing a nondefective screw is $q = 1 - p = 1 - M/N$, and Eq. (3.52) gives the probability of drawing x defectives in n trials in the form:

$$f(x) = \binom{n}{x} \left(\frac{M}{N}\right)^x \left(1 - \frac{M}{N}\right)^{n-x} \quad (x = 0, 1, \dots, n). \quad (3.55)$$

³⁷put it back to the given set and mix.

3.7.2 Sampling without Replacement: Hyper-geometric Distribution

Sampling without replacement means that we return no screw to the box. Then we no longer have independence of trials, and instead of Eq. (3.55) the probability of drawing x defectives in n trials is:

$$f(x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}} \quad \text{where} \quad x = 1, 2, \dots, n. \quad (3.56)$$

The distribution with this probability function is called the **hyper-geometric distribution**³⁸.

The hypergeometric distribution has the mean:

$$\mu = n \frac{M}{N},$$

and the variance

$$\sigma^2 = \frac{nM(N-M)(N-n)}{N^2(N-1)}.$$

³⁸because its moment generating function can be expressed by the hypergeometric function, which is a fact only useful to write it in a margin.

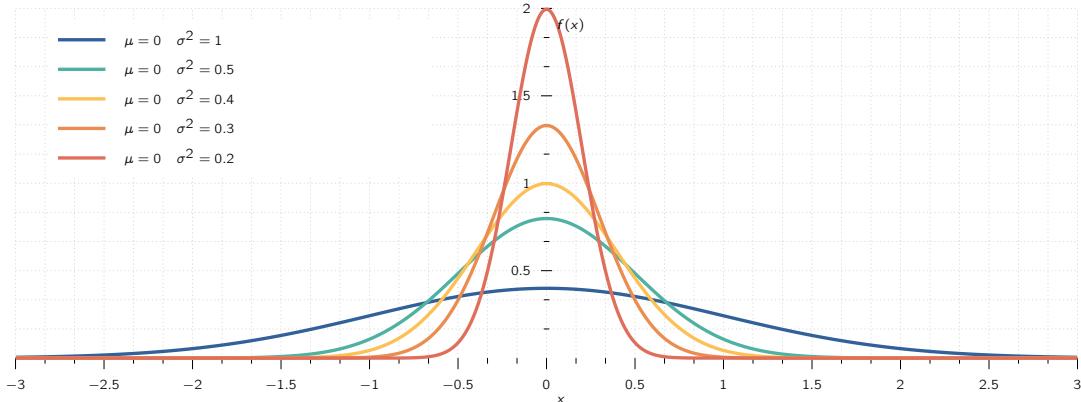


Figure 3.5: The poster child of probability and statistics, the normal distribution.

3.8 Normal Distribution

Turning from discrete to continuous distributions, in this section we discuss the normal distribution. This is the most important continuous distribution because in applications many random variables are normal random variables³⁹ or they are approximately normal or can be transformed into normal random variables in a relatively simple fashion. Furthermore, the normal distribution is a useful approximation of more complicated distributions, and it also occurs in the proofs of various statistical tests.

³⁹that is, they have a normal distribution.

The **normal distribution** or *Gauss distribution* is defined as the distribution with the density:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) \quad (3.57)$$

where \exp is the exponential function with base $e = 2.718\cdots$. This is simpler than it may at first look. $f(x)$ has these features (see also Fig. 3.5).

1. μ is the mean, and σ the standard deviation.
2. $1/(\sigma\sqrt{2\pi})$ is a constant factor that makes the area under the curve of $f(x)$ from $-\infty$ to ∞ equal to 1, as it must be⁴⁰.
3. The curve of $f(x)$ is symmetric with respect to $x = \mu$ because the exponent is quadratic. Hence for $\mu = 0$ it is symmetric with respect to the y -axis $x = 0$ ⁴¹.
4. The exponential function in Eq. (3.57) goes to zero very fast—the faster the smaller the standard deviation σ is, as it should be, as seen in Fig. 3.5.

⁴⁰Having a probability higher than 1 does NOT make sense

⁴¹This distribution is also known as bell-shaped curves.

3.8.1 Distribution Function

From Eq. (3.55) and Eq. (3.57) we see that the normal distribution has the **distribution function** of the following form:

$$F(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left[-\frac{1}{2}\left(\frac{v-\mu}{\sigma}\right)^2\right] dv. \quad (3.58)$$

Here we needed x as the upper limit of integration and wrote v (instead of x) in the integrand.

For the corresponding **standardised normal distribution** with mean 0 and standard deviation 1 we denote $F(x)$ by $\Phi(z)$. Then we simply have from Eq. (3.58).

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du. \quad (3.59)$$

This integral cannot be integrated by one of the methods of calculus.

But this is no serious handicap because its values can be obtained from standardised tables. These values are needed in working with the normal distribution. The curve of $\Phi(z)$ is *S*-shaped. It increases monotone from 0 to 1 and intersects the vertical axis at $\frac{1}{2}$, as shown in **Fig. 3.6**.

Theory 3.13: Relationship between PDF and CDF

The distribution function $F(x)$ of the normal distribution with any μ and σ is related to the standardised distribution function $\Phi(z)$ in Eq. (3.59) by the formula

$$F(x) = \Phi\left(\frac{x-\mu}{\sigma}\right).$$

Theory 3.14: Normal Probabilities for Intervals

The probability a normal random variable X with mean μ and standard deviation σ assume any value in an interval $a < x \equiv b$ is:

$$P(a < X \leq b) = F(b) - F(a) = \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right).$$

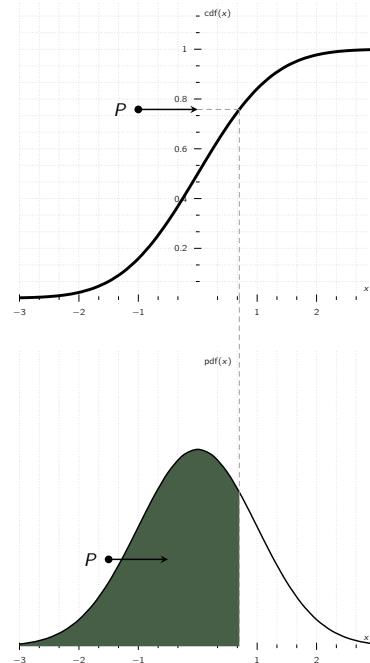


Figure 3.6: A visual representation between the relationship of PDF and CDF.

3.8.2 Numeric Values

In practical work with the normal distribution it is good to remember that about 67% of all values of X to be observed will be between $\mu \pm \sigma$, about 95% between $\mu \pm 2\sigma$, and practically all between

the **three-sigma limits** $\mu \pm 3\sigma$:

$$P(\mu - \sigma < X \leq \mu + \sigma) \approx 68\% \quad (3.60a)$$

$$P(\mu - 2\sigma < X \leq \mu + 2\sigma) \approx 95.5\% \quad (3.60b)$$

$$P(\mu - 3\sigma < X \leq \mu + 3\sigma) \approx 99.7\%. \quad (3.60c)$$

The aforementioned formulas show that a value deviating from μ by more than σ , 2σ , or 3σ will occur in one of about 3, 20, and 300 trials, respectively.

⁴²Which we shall cover in Chapter 4. In tests⁴², we shall ask, conversely, for the intervals that correspond to certain given probabilities; practically most important use the probabilities of 95%, 99%, and 99.9%. For these, the answers are $\mu \pm 2\sigma$, $\mu \pm 2.6\sigma$, and $\mu \pm 3.3\sigma$, respectively.

More precisely,

$$P(\mu - 1.96\sigma < X \leq \mu + 1.96\sigma) \approx 95\% \quad (3.61a)$$

$$P(\mu - 2.58\sigma < X \leq \mu + 2.58\sigma) \approx 99\% \quad (3.61b)$$

$$P(\mu - 3.29\sigma < X \leq \mu + 3.29\sigma) \approx 99.9\%. \quad (3.61c)$$

3.8.3 Normal Approximation of the Binomial Distribution

The probability function of the binomial distribution, as a reminder, is:

$$f(x) = \binom{n}{x} p^x q^{n-x} \quad (x = 0, 1, \dots, n). \quad (3.62)$$

If n is large, the binomial coefficients and powers become very inconvenient. It is of great practical⁴³ importance that, in this case, the normal distribution provides a good approximation of the binomial distribution, according to the following theorem, one of the most important theorems in all probability theory.

Theory 3.15: Limit Theorem of De Moivre and Laplace

For large n ,

$$f(x) \sim f^*(x) \quad \text{where} \quad x = 0, 1, \dots, n$$

Here f is given by Eq. (3.62). The function

$$f^*(z) = \frac{1}{\sqrt{2\pi}\sqrt{npq}} \exp\left(-\frac{z^2}{2}\right), \quad \text{and} \quad z = \frac{x - np}{\sqrt{npq}}$$

is the density of the normal distribution with mean $\mu = np$ and variance $\sigma^2 = npq$ (the mean and variance of the binomial distribution). Furthermore, for any nonnegative integers a and b ($> a$):

$$P(a \leq X \leq b) = \sum_{x=a}^b \binom{n}{x} p^x q^{n-x} \sim \Phi(\beta) - \Phi(\alpha)$$

where,

$$\alpha = \frac{a - np - 0.5}{\sqrt{npq}} \quad \text{and} \quad \beta = \frac{b - np + 0.5}{\sqrt{npq}}$$

3.9 Distribution of Several Random Variables

Distributions of two (2) or more random variables are of interest for two (2) reasons:

1. They occur in experiments in which we observe several random variables, for example, carbon content X and hardness Y of steel, amount of fertiliser X and yield of corn Y , height X_1 , weight X_2 , and blood pressure X_3 of persons, and so on.
2. They will be needed in the mathematical justification of the methods of statistics in Chapter 4.

In this section we consider two (2) random variables X and Y or, as we also say, a **two-dimensional random variable** (X, Y) . For (X, Y) the outcome of a trial is a pair of numbers $X = x, Y = y$, briefly $(X, Y) = (x, y)$, which we can plot as a point in the XY -plane.

The **two-dimensional probability distribution** of the random variable (X, Y) is given by the **distribution function**

$$F(x, y) = P(X \leq x, Y \leq y). \quad (3.63)$$

This is the probability that in a trial, X will assume any value not greater than x and in the same trial, Y will assume any value not greater than y . $F(x, y)$ determines the probability distribution **uniquely**, because extending the analogy we developed previously, $P(a < X \leq b) = F(b) - F(a)$, we now have for a rectangle defined using the following equation:

$$P(a_1 < X \leq b_1, a_2 < Y \leq b_2) = F(b_1, b_2) - F(a_1, b_2) - F(b_1, a_2) + F(a_1, a_2). \quad (3.64)$$

As before, in the two-dimensional case we shall also have **discrete** and **continuous** random variables and distributions.

3.9.1 Discrete Two-Dimensional Distribution

In analogy to the case of a single random variable, we call (X, Y) and its distribution **discrete** if (X, Y) can assume only finitely many or at most countably infinitely many pairs of values $(x_1, y_1), (x_2, y_2), \dots$ with positive probabilities, whereas the probability for any domain containing none of those values of (X, Y) is zero.

Let (x_i, y_i) be any of those values and let $P(X = x_i, Y = y_j) = p_{ij}$ (where we admit that p_{ij} may be 0 for certain pairs of subscripts i). Then we define the **probability function** $f(x, y)$ of (X, Y) by:

$$f(x, y) = p_{ij} \quad \text{if} \quad x = x_i, y = y_j \quad \text{and} \quad f(x, y) = 0 \quad \text{otherwise};$$

where, $i = 1, 2, \dots$ and $j = 1, 2, \dots$ independently. In analogy to Eq. (3.37), we now have for the distribution function the formula:

$$F(x, y) = \sum_{x_i \leq x} \sum_{y_j \leq y} f(x_i, y_j).$$

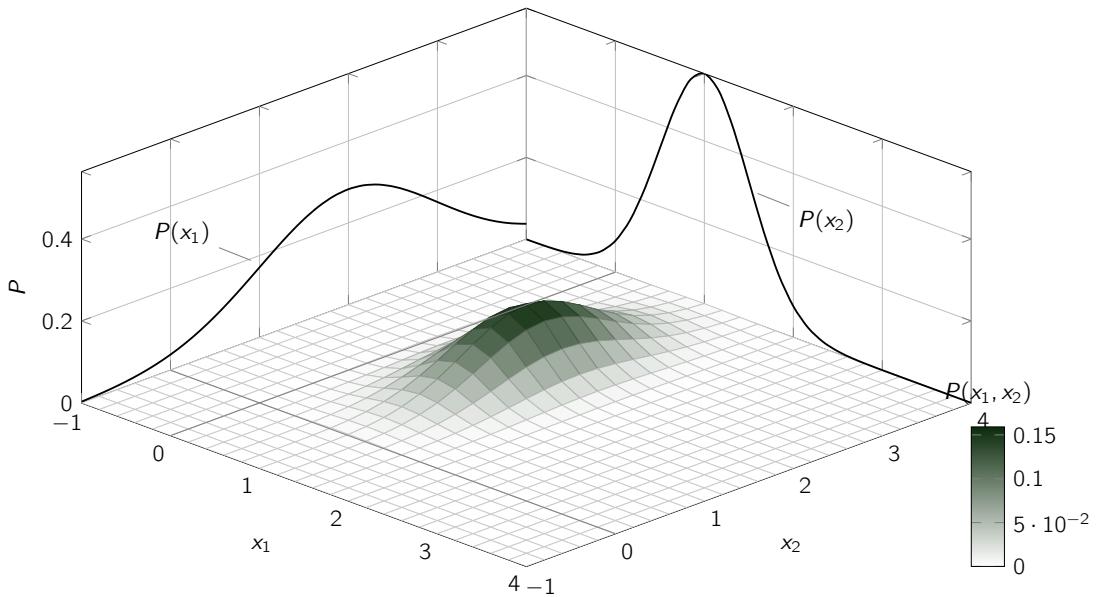


Figure 3.7: Many samples from a bivariate normal distribution. The marginal distributions are shown on the z-axis. The marginal distribution of X is also approximated by creating a histogram of the X coordinates without consideration of the Y coordinates.

Instead of Eq. (3.39), we now have the condition:

$$\sum_i \sum_j f(x_i, y_j) = 1.$$

3.9.2 Continuous Two-Dimensional Distribution

In analogy to the case of a single random variable, we call (X, Y) and its distribution **continuous** if the corresponding distribution function $F(x, y)$ can be given by a double integral:

$$F(x, y) = \int_{-\infty}^y \int_{-\infty}^x f(x^*, y^*) dx^* dy^* \quad (3.65)$$

whose integrand f , called the **density** of (X, Y) , is non-negative everywhere, and is continuous, possibly except on finitely many curves.

From Eq. (3.65) we obtain the probability that (X, Y) assume any value in a rectangle (Fig. 523) given by the formula:

$$P(a_1 < X \leq b_1, a_2 < Y \leq b_2) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} f(x, y) dx dy$$

3.9.3 Marginal Distributions of a Discrete Distribution

This is a rather natural idea, without counterpart for a single random variable.

It amounts to being interested only in one of the two variables in (X, Y) , say, X , and asking for its distribution, called the **marginal distribution** of X in (X, Y) . So we ask for the probability $P(X = x, Y \text{ arbitrary})$.

Since (X, Y) is discrete, so is X . We get its probability function, call it $f_1(x)$, from the probability function $f(x, y)$ of (X, Y) by summing over y :

$$f_1(x) = P(X = x, Y, \text{arbitrary}) = \sum_y f(x, y) \quad (3.66)$$

where we sum all the values of $f(x, y)$ that are not 0 for that x .

From Eq. (3.66) we see that the distribution function of the marginal distribution of X is

$$F_1(x) = P(X \leq x, Y, \text{arbitrary}) = \sum_{x^* \leq x} f_1(x^*).$$

Similarly, the probability function

$$f_2(y) = P(X, \text{arbitrary}, Y \equiv y) = \sum_x f(x, y)$$

determines the **marginal distribution** of Y in (X, Y) . Here we sum all the values of $f(x, y)$ that are not zero for the corresponding y . The distribution function of this marginal distribution is

$$F_2(y) = P(X, \text{arbitrary}, Y \equiv y) = \sum_{y^* \equiv y} f_2(y^*).$$

Exercise 3.20: Marginal Distributions of a Discrete Two-Dimensional Random Variable

In drawing 3 cards with replacement from a bridge deck let us consider

$$(X, Y) \quad \text{where } X = \text{Number of queens} \quad \text{and} \quad Y = \text{Number of kings or aces.}$$

The deck has 52 cards. These include 4 queens, 4 kings, and 4 aces. Therefore, in a single trial a queen has probability:

$$\frac{4}{52} = \frac{1}{13}$$

and a king or ace:

$$\frac{8}{52} = \frac{2}{13}$$

This gives the probability function of (X, Y) as:

$$f(x, y) = \frac{3!}{x!y!(3-x-y)!} \left(\frac{1}{13}\right)^x \left(\frac{2}{13}\right)^y \left(\frac{10}{13}\right)^{3-x-y} \quad \text{where } (x + y \leq 3)$$

and $f(x, y) = 0$ otherwise.

3.9.4 Marginal Distributions of a Continuous Distribution

This is conceptually the same as for discrete distributions, with probability functions and sums replaced by densities and integrals. For a continuous random variable (X, Y) with density $f(x, y)$ we now have the **marginal distribution** of X in (X, Y) , defined by the distribution function

$$F_1(x) = P(X \leq x, -\infty < Y < \infty) = \int_{-\infty}^x f_1(x^*) dx^*$$

with the density f_1 of X obtained from $f(x, y)$ by integration over y ,

$$f_1(x) = \int_{-\infty}^{\infty} f(x, y) dy.$$

Interchanging the roles of X and Y , we obtain the **marginal distribution** of Y in (X, Y) with the distribution function

$$F_2(y) = P(-\infty < X < \infty, Y \leq y) = \int_{-\infty}^y f_2(y^*) dy^*$$

and density

$$f_2(y) = \int_{-\infty}^{\infty} f(x, y) dx.$$

3.9.5 Independence of Random Variables

X and Y in a, discrete or continuous, random variable (X, Y) are said to be **independent** if

$$F(x, y) = F_1(x)F_2(y)$$

holds for all (x, y) . Otherwise these random variables are said to be **dependent**. Necessary and sufficient for independence is

$$f(x, y) = f_1(x)f_2(y)$$

for all x and y . Here the f 's are the above probability functions if (X, Y) is discrete or those densities if (X, Y) is continuous.

Exercise 3.21: Independence and Dependence

In tossing a 50 cent and a 20 cent coin, with X being the number of heads on the 50 cent, and Y number of heads on the 20 cent, we may assume the values 0 or 1 and are independent.

Extension of Independence to n -Dimensional Random Variables. This will be needed throughout Chapter 4. The distribution of such a random variable $\mathbf{X} = (X_1, \dots, X_n)$ is determined by a **distribution function** of the form

$$F(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n).$$

The random variables X_1, \dots, X_n are said to be **independent** if

$$F(x_1, \dots, x_n) = F_1(x_1)F_2(x_2) \cdots F_n(x_n)$$

for all (x_1, \dots, x_n) . Here $F_j(x_j)$ is the distribution function of the marginal distribution of X_j in \mathbf{X} , that is,

$$F_j(x_j) = P(X_j \leq x_j, X_k \text{ arbitrary}, k = 1, \dots, n, k \neq j).$$

Otherwise these random variables are said to be **dependent**.

3.9.6 Functions of Random Variables

When $n = 2$, we write $X_1 = X$, $X_2 = Y$, $x_1 = x$, $x_2 = y$. Taking a non-constant continuous function $g(x, y)$ defined for all x, y , we obtain a random variable $Z = g(X, Y)$.

For example, if we roll two (2) dice and X and Y are the numbers the dice turn up in a trial, then $Z = X + Y$ is the sum of those two (2) numbers.

In the case of a **discrete** random variable (X, Y) we may obtain the probability function $f(z)$ of $Z = g(X, Y)$ by summing all $f(x, y)$ for which $g(x, y)$ equals the value of z considered; thus

$$f(z) = P(Z = z) = \sum_{g(x,y)=z} \sum f(x, y).$$

Hence the distribution function of Z is

$$F(z) = P(Z \leq z) = \sum_{g(x,y) \leq z} \sum f(x, y),$$

where we sum all values of $f(x, y)$ for which $g(x, y) \leq z$.

In the case of a **continuous** random variable (X, Y) we similarly have

$$F(z) = P(Z \leq z) = \iint_{g(x,y) \leq z} f(x, y) dx dy$$

where for each z we integrate the density $f(x, y)$ of (X, Y) over the region $g(x, y) \leq z$ in the xy -plane, the boundary curve of this region being $g(x, y) = z$.

3.9.7 Addition of Means

The number

$$E(g(X, Y)) = \begin{cases} \sum_x \sum_y g(x, y) f(x, y) & \text{where } X, Y \text{ are discrete} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dx dy & \text{where } X, Y \text{ are continuous} \end{cases} \quad (3.67)$$

is called the **mathematical expectation** or, briefly, the **expectation of $g(X, Y)$** . Here it is assumed that the double series converges absolutely and the integral of $|g(x, y)|/(x, y)$ over the y -plane exists⁴⁴. Since summation and integration are linear processes, we have from Eq. (3.67):

⁴⁴meaning it is finite.

$$E(ag(X, Y) + bh(X, Y)) = aE(g(X, Y)) + bE(h(X, Y))$$

An important special case is

$$E(X + Y) = E(X) + E(Y),$$

and by induction we have the following result.

Theory 3.16: Addition of Means

The mean (expectation) of a sum of random variables equals the sum of the means (expectations), that is,

$$E(X_1 + X_2 + \dots + X_n) = E(X_1) + E(X_2) + \dots + E(X_n).$$

We can also deduce the following statement:

Theory 3.17: Multiplication of Means

The mean (expectation) of the product of independent random variables equals the product of the means (expectations), that is,

$$E(X_1 X_2 \dots X_n) = E(X_1) E(X_2) \dots E(X_n).$$

⁴⁵This is left as an exercise to the reader.

3.9.8 Addition of Variances

A final matter to cover is how we can sum up variances. Similar to before, let $Z = X + Y$ and denote the mean and variance of Z by μ and σ^2 .

Then we first have:

$$\sigma^2 = E([Z - \mu]^2) = E(Z^2) - [E(Z)]^2$$

From (24) we see that the first term on the right equals

$$E(Z^2) = E(X^2 + 2XY + Y^2) = E(X^2) + 2E(XY) + E(Y^2).$$

For the second term on the right we obtain from Theorem 1

$$[E(Z)]^2 = [E(X) + E(Y)]^2 = [E(X)]^2 + 2E(X)E(Y) + [E(Y)]^2$$

By substituting these expressions into the formula for σ^2 we have

$$\begin{aligned} \sigma^2 &= E(X^2) - [E(X)]^2 + E(Y^2) - [E(Y)]^2 \\ &\quad + 2[E(XY) - E(X)E(Y)]. \end{aligned}$$

the expression in the first line on the right is the sum of the variances of X and Y , which we denote by σ_1^2 and σ_2^2 , respectively.

The quantity in the second line (except for the factor 2) is:

$$\sigma_{XY} = E(XY) - E(X)E(Y), \tag{3.68}$$

and is called the **covariance** of X and Y . Consequently, our result is

$$\sigma^2 = \sigma_1^2 + \sigma_2^2 + 2\sigma_{XY}.$$

If X and Y are **independent**, then

$$E(XY) = E(X)E(Y);$$

hence $\sigma_{XY} = 0$, and

$$\sigma^2 = \sigma_1^2 + \sigma_2^2$$

Extension to more than two variables gives the basic

Theory 3.18: Addition of Variances

The variance of the sum of independent random variables equals the sum of the variances of these variables.

Chapter 4

Statistical Methods

Table of Contents

4.1	Introduction	101
4.2	Point Estimation of Parameters	104
4.3	Confidence Intervals	108
4.4	Testing of Hypotheses and Making Decisions	115
4.5	Goodness of Fit	121

4.1 Introduction

Statistical¹ methods consists of a wide range of tools for designing and evaluating random experiments to obtain information about practical problems:

¹The word is derived from New Latin *statistica* or *statisticus* ("of the state")

such as exploring the relation between iron content and density of iron ore, the quality of raw material or manufactured products, the efficiency of air-conditioning systems, the performance of certain cars, the effect of advertising, the reactions of consumers to a new product, etc.

Therefore, the diameter of screws is a random variable X and we have non-defective screws, with diameter between given tolerance limits, and defective screws, with diameter outside those limits. We can ask for the distribution of X , for the percentage of defective screws to be expected, and for necessary improvements of the production process.

Samples are selected from populations:

20 screws from 1000 screws, 100 of 5000 voters, 8 behaviours in a wildlife observation.

²It would be inconceivable for a company who produces over a billion light bulbs to test all their products. That is why we have return policies.

³of being drawn when we sample.

as inspecting the entire sample, would be expensive, time-consuming, impossible or even senseless.²

To obtain a meaningful sense of information, samples must be **random selections**. Each of the 1000 screws must have the same chance of being sampled;³ at least approximately. Only then will the sample mean:

$$\bar{x} = \frac{1}{20} (x_1 + \cdots + x_{20}) \quad \text{where} \quad n = 20,$$

will be a **good approximation** of the population mean μ , and the accuracy of the approximation will generally improve with increasing n , as we shall see.

This is also applicable to other statistical quantities such as standard deviation, variance, etc.

Independent sample values will be obtained in experiments with an infinite sample space S certainly for the **normal distribution**. This is also true in sampling with replacement. It is approximately true in drawing **small samples** from a large finite population.⁴ However, if we sample without replacement from a small population, the effect of dependence of sample values may be considerable.

⁴for instance, 5 or 10 of 1000 items.

Random numbers help in obtaining samples that are in fact random selections. This is sometimes not easy to accomplish as there are numerous subtle factors which can bias sampling.⁵ Random numbers can be obtained from a **random number generator**

It is important to state that the numbers generated by a computer are **NOT** truly random, as are calculated by a tricky formula that produces numbers that do have practically all the essential features of true randomness. Because these numbers eventually repeat, they must not be used in cryptography, for example, where true randomness is required.

Exercise 4.1: Generating Random Numbers

To select a sample of size $n = 10$ from 80 given ball bearings, we number the bearings from 1 to 80. We then let the generator randomly produce 10 of the integers from 1 to 80 and include the bearings with the numbers obtained in our sample, for example,

44 55 57 03 61 51 68 22 34 77

or whichever number pops up in your head.⁶

⁵Such as by personal interviews, by poorly working machines, by the choice of non-typical observation conditions, etc.

⁶Of course in a professional setting you can't just write numbers like that as there is also a pattern when we make successive random numbers. Before the prevalence of computers there used to be books containing random numbers which people consulted.

Representing and processing data were considered in the previous chapter in connection with **frequency distributions**. These are the **empirical counterparts** of probability distributions and helped motivating axioms and properties in probability theory. The new aspect in this chapter is **randomness**:

i.e., the data are samples selected **randomly** from a population.

Accordingly, we can already use the plots we have used in probability, such as stem-and-leaf plots, box plots, and histograms.

In this chapter, the mean \bar{x} we defined previously, will now be referred as **sample mean**.

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j = \frac{1}{n} (x_1 + x_2 + \dots + x_n) . \quad (4.1)$$

We call n the **sample size**, and similar to mean, the variance s^2 is called the **sample variance**:

$$s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2 = \frac{1}{n-1} [(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2] , \quad (4.2)$$

and its positive square root, s is the **sample standard deviation**.

\bar{x}, s^2, s are called **sample parameters** of a dataset.

4.2 Point Estimation of Parameters

Before we dive deep into statistics, let's spend some time to learn the most basic practical tasks in statistics and corresponding statistical methods to accomplish them. The first is point **estimation of parameters**, that is, of **quantities** appearing in distributions:

such as p in the binomial distribution and μ and σ in the normal distribution.

⁷which is a point on the real line.

A **point estimate** of a parameter is a number,⁷ which is computed from a given sample and serves as an **approximation of the unknown exact value** of the parameter of the population. An interval estimate is an interval⁸ obtained from a sample. Think of it as a value which is a sensible guess for that parameter.

Estimation of parameters is of great practical importance in many applications.

⁹To describe something which is an approximation or an educated guess, we use hat (i.e., \hat{x}) notation. This is applicable for fields in statistics, machine learning or data science.

As an approximation⁹ of the mean of a population we may take the mean \bar{x} of a corresponding sample. This gives the estimate $\hat{\mu} = \bar{x}$ for μ , that is,

$$\hat{\mu} = \bar{x} = \frac{1}{n} (x_1 + \dots + x_n), \quad (4.3)$$

where n is the sample size. Similarly, an estimate $\hat{\sigma}^2$ for the variance of a population is the variance s^2 of a corresponding sample, that is:

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2. \quad (4.4)$$

As can be seen, both Eq. (4.3) and Eq. (4.4) are **estimates** of parameters for distributions in which μ or σ^2 appear explicitly as parameters, such as the normal and Poisson distributions.

An estimator is not expected to estimate the population parameter without error. We do not expect \bar{x} to estimate μ exactly, but we certainly hope that it is not far off.

For the binomial distribution, $p = \mu/n$. From Eq. (4.3) we obtain for p the estimate:

$$\hat{p} = \frac{\bar{x}}{n}. \quad (4.5)$$

It is important to mention Eq. (4.3) is a special case of the so-called **method of moments**. Here, the parameters to be estimated are expressed in terms of the moments of the distribution. In the resulting formulas, those moments of the distribution are replaced by the corresponding moments of the sample, which gives the estimates. Here the k^{th} moment of a sample x_1, \dots, x_n is:

$$m_k = \frac{1}{n} \sum_{j=1}^n x_j^k. \quad (4.6)$$

4.2.1 Maximum Likelihood Method

Another method for obtaining estimates is

the so-called **maximum likelihood method** conceived by R. A. Fisher.¹⁰ To explain it, we consider a discrete (or continuous) random variable X whose probability function (or density) $f(x)$ depends on a single parameter θ . We take a corresponding sample of n independent values x_1, \dots, x_n . Then in the discrete case the probability given a sample of size n consists precisely of those n values is

$$I = f(x_1) f(x_2) \cdots f(x_n). \quad (4.7)$$



¹⁰Considered the father of modern statistics. For his work in statistics, he has been described as "a genius who almost single-handedly created the foundations for modern statistical science" and "the single most important figure in 20th century statistics". Fisher has also been praised as a pioneer of the Information Age. His work on a mathematical theory of information ran parallel to the work of Claude Shannon and Norbert Wiener, though based on statistical theory.

¹¹not at the boundary.

In the continuous case the probability that the sample consists of values in the small intervals $x_j \leq x \leq x_j + \Delta x (j = 1, 2, \dots, n)$ is

$$f(x_1) \Delta x f(x_2) \Delta x \cdots f(x_n) \Delta x = I(\Delta x)^n \quad (4.8)$$

As $f(x_j)$ depends on θ , the function I in Eq. (4.8) given by Eq. (4.7) depends on x_1, \dots, x_n and θ .

We imagine x_1, \dots, x_n to be given and fixed.

Then I is a function of θ , which is called the **likelihood function**. The basic idea of the maximum likelihood method is quite simple, as follows.

We choose an approximation for the unknown value of θ for which I is as large as possible.

If I is a differentiable function of θ , a necessary condition for I to have a maximum in an interval¹¹ is

$$\frac{\partial I}{\partial \theta} = 0 \quad (4.9)$$

A solution of Eq. (4.9) depending on x_1, \dots, x_n is called a **maximum likelihood estimate** for θ .

We may replace Eq. (4.9) by:

$$\frac{\partial \ln I}{\partial \theta} = 0 \quad (4.10)$$

as $f(x_j) > 0$, a maximum of I is in general positive, and $\ln I$ is a monotone increasing function of I . This often simplifies calculations.

Several Parameters

If the distribution of X involves r parameters $\theta_1, \dots, \theta_r$, then instead of Eq. (4.9) we have the r conditions $\partial \ln I / \partial \theta_1, \dots, \partial \ln I / \partial \theta_r = 0$, and instead of Eq. (4.10) we have:

$$\frac{\partial \ln I}{\partial \theta_1} = 0, \dots, \frac{\partial \ln I}{\partial \theta_r} = 0. \quad (4.11)$$

Exercise 4.2: Maximum Likelihood of Gaussian Distribution

Find maximum likelihood estimates for $\theta_1 = \mu$ and $\theta_2 = \sigma$ in the case of the normal distribution.

Solution

We obtain the likelihood function:

$$L = \left(\frac{1}{\sqrt{2\pi}} \right)^n \left(\frac{1}{\sigma} \right)^n e^{-h}$$

where $h = \frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2$.

Taking logarithms, we have

$$\ln L = -n \ln \sqrt{2\pi} - n \ln \sigma - h.$$

The first equation in Eq. (4.11) is $\frac{\partial \ln L}{\partial \mu} = 0$, written out:

$$\frac{\partial \ln L}{\partial \mu} = -\frac{\partial h}{\partial \mu} = \frac{1}{\sigma^2} \sum_{j=1}^n (x_j - \mu) = 0,$$

therefore $\sum_{j=1}^n x_j - n\mu = 0$.

The solution is the desired estimate $\hat{\mu}$ for μ : we find

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n x_j = \bar{x}.$$

The second equation in Eq. (4.11) is $\frac{\partial \ln L}{\partial \sigma} = 0$, written out

$$\frac{\partial \ln L}{\partial \sigma} = -\frac{n}{\sigma} - \frac{\partial h}{\partial \sigma} = -\frac{1}{\sigma} + \frac{1}{\sigma^3} \sum_{j=1}^n (x_j - \mu)^2 = 0.$$

Replacing μ by $\hat{\mu}$ and solving for σ^2 , we obtain the estimate:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 \quad \blacksquare$$

Exercise 4.3: Maximum Likelihood of Poisson Distribution

Consider a Poisson distribution:

$$f(x|\mu) = \frac{e^{-\mu} \mu^x}{x!}, \quad x = 0, 1, 2, \dots$$

Suppose that a random sample x_1, x_2, x_n is taken from the distribution. What is the maximum likelihood estimate of μ ?

Solution

The likelihood function is

$$L(x_1, x_2, \dots, x_n; \mu) = \prod_{i=1}^n f(x_i | \mu) = \frac{e^{-n\mu} \sum_{i=1}^n x_i}{\prod_{i=1}^n x_i!}.$$

Now consider

$$\ln L(x_1, x_2, \dots, x_n; \mu) = -n\mu + \sum_{i=1}^n x_i \ln \mu - \ln \prod_{i=1}^n x_i!$$

$$\frac{\partial \ln L(x_1, x_2, \dots, x_n; \mu)}{\partial \mu} = -n + \sum_{i=1}^n \frac{x_i}{\mu}.$$

Solving for $\hat{\mu}$, the maximum likelihood estimator, involves setting the derivative to zero and solving for the parameter. Thus,

$$\hat{\mu} = \sum_{i=1}^n \frac{x_i}{n} = \bar{x}.$$

The second derivative of the log-likelihood function is negative, which implies that the solution above indeed is a maximum. Since μ is the mean of the Poisson distribution, the sample average would certainly seem like a reasonable estimator.

Exercise 4.4: For Science

Suppose ten (10) rats are used in a biomedical study where they are injected with cancer cells and then given a cancer drug that is designed to increase their survival rate. The survival times, in months, are:

14 17 27 18 12 8 22 13 19 12

Assume exponential distribution applies which is given as:

$$f(x, \beta) = \begin{cases} 1/\beta \exp^{-x/\beta}, & x > 0, \\ 0, & \text{elsewhere.} \end{cases}$$

Give a maximum likelihood estimate of the mean survival time.

Solution

We know that the probability density function for the exponential random variable X . Therefore, the log-likelihood function for the data, given $n = 10$, is:

$$\ln L(x_1, x_2, \dots, x_{10}; \beta) = -10 \ln \beta - \frac{1}{\beta} \sum_{i=1}^{10} x_i.$$

Setting

$$\frac{\partial \ln L}{\partial \beta} = -\frac{10}{\beta} + \frac{1}{\beta^2} \sum_{i=1}^{10} x_i = 0$$

implies that

$$\hat{\beta} = \frac{1}{10} \sum_{i=1}^{10} x_i = \bar{x} = 16.2.$$

Evaluating the second derivative of the log-likelihood function at the value $\hat{\beta}$ above yields a negative value. As a result, the estimator of the parameter β , the population mean, is the sample average \bar{x} .

Exercise 4.5: Sampling the Population

It is known that a sample consisting of the values:

12 11.2 13.5 12.3 13.8 11.9

comes from a population with the density function:

$$f(x; \theta) = \begin{cases} \frac{\theta}{\theta x + 1}, & x > 1, \\ 0, & \text{elsewhere,} \end{cases} \quad (4.12)$$

where $\theta > 0$. Find the maximum likelihood estimate of θ .

Solution

The likelihood function of n observations from this population can be written as:

$$L(x_1, x_2, \dots, x_{10}; \theta) = \prod_{i=1}^n \frac{\theta}{x_i^{\theta+1}} = \frac{\theta^n}{(\prod_{i=1}^n x_i)^{\theta+1}},$$

which implies that

$$\ln L(\cdot) | x_1, x_2, \dots, x_{10}; \theta = n \ln \theta - (\theta + 1) \sum_{i=1}^n \ln x_i.$$

Setting $0 = \frac{\partial \ln L}{\partial \theta} = \frac{n}{\theta} - \sum_{i=1}^n \ln(x_i)$ results in

$$\hat{\theta} = \frac{n}{\sum_{i=1}^n \ln x_i} = 0.3970 \quad \blacksquare.$$

Since the second derivative of L is $-n/\theta^2$, which is always negative, the likelihood function does achieve its maximum value at $\hat{\theta}$.

4.3 Confidence Intervals



Confidence intervals¹² for an unknown parameter θ of some distribution (e.g., $\theta = \mu$) are intervals $\theta_1 \leq \theta \leq \theta_2$ which contain θ , not with certainty but with a **high probability** γ , which we can choose.¹³ Such an interval is calculated from a sample. $\gamma = 95\%$ means probability $1 - \gamma = 5\% = 1/20$ of being wrong.¹⁴ Instead of writing $\theta_1 \leq \theta \leq \theta_2$, we denote this more **distinctly** by writing:

$$\text{CONF}_{\gamma} \{ \theta_1 \leq \theta \leq \theta_2 \} \quad (4.13)$$

¹²Established by Jerzy Neyman. He proposed and studied randomised experiments in 1923.

Furthermore, his paper *On the Two Different Aspects of the Representative Method: The Method of Stratified Sampling and the Method of Purposive Selection*, given at the Royal Statistical Society on 19 June 1934, was the groundbreaking event leading to modern scientific sampling. He introduced the confidence interval in his paper in 1937. Another noted contribution is the Neyman-Pearson lemma, the basis of hypothesis testing.

¹³95% and 99% are popular

¹⁴one of about 20 such intervals will **NOT** contain θ .

¹⁵not in the strict sense of numerical means, but except for an error whose probability we know.

¹⁶with the same distribution, namely, that of X

Such a special symbol, CONF, seems worthwhile to avoid the misunderstanding that θ **must** lie between θ_1 and θ_2 .

γ is called the **confidence level**, and θ_1 and θ_2 are called the **lower** and **upper confidence limits**, respectively and **depend** on the γ value. The larger we **choose** γ , the smaller is the error probability $1 - \gamma$, but the longer is the confidence interval.

If $\gamma \rightarrow 1$, then its length goes to infinity.

The choice of γ depends on the kind of application.

In taking no umbrella, a 5% chance of getting wet is **NOT** a problem. In a medical decision of life or death, a 5% chance of being wrong may be too large and a 1% chance of being wrong ($\gamma = 99\%$) may be more desirable.

Confidence intervals are more valuable than point estimates. We can take the midpoint of Eq. (4.13) as an approximation of θ and half the length of Eq. (4.13) as an error bound.¹⁵

θ_1 and θ_2 in Eq. (4.13) are calculated from a sample x_1, \dots, x_n . These are n observations of a random variable X . Now comes a **standard trick**.

We regard x_1, \dots, x_n as single observations of n random variables X_1, \dots, X_n ¹⁶. Then $\theta_1 = \theta_1(x_1, \dots, x_n)$ and $\theta_2 = \theta_2(x_1, \dots, x_n)$ in Eq. (4.13) are observed values of two random variables $\Theta_1 = \Theta_1(X_1, \dots, X_n)$ and $\Theta_2 = \Theta_2(X_1, \dots, X_n)$. The condition Eq. (4.13) involving γ can now be written

$$P(\Theta_1 \leq \theta \leq \Theta_2) = \gamma. \quad (4.14)$$

Let us see what all this means in concrete practical cases.

In each case in this section we shall first state the steps of obtaining a confidence interval in the form of a table, then consider a typical example, and finally justify those steps theoretically.

For Mean with known Variance in Normal Distribution

The method of tackling is this problem is as follows:

1. Choose a confidence level for γ ¹⁷.

¹⁷ 95%, 99%, depending on the application.

2. Determine the corresponding c :

γ	0.90	0.95	0.99	0.999
c	1.645	1.960	2.576	3.291

Table 4.1: Useful c values based on a given confidence (γ) value.

3. Compute the mean \bar{x} of the sample x_1, \dots, x_n .
4. Compute $k = c\sigma/\sqrt{n}$. The confidence interval for μ is

$$\text{CONF}_\gamma \{\bar{x} - k \leq \mu \leq \bar{x} + k\}. \quad (4.15)$$

Exercise 4.6: Confidence Interval for mean with known variance in Normal Distribution

Determine 95% confidence interval for the mean of a normal distribution with variance $\sigma^2 = 9$, using a sample of $n = 100$ values with mean $\bar{x} = 5$.

Solution

1. First we define γ as 0.95.
2. Then looking at the table find the corresponding c which equals 1.960.
3. $\hat{x} = 5$ is given.

4. We need:

$$k = c \frac{\sigma}{\sqrt{n}} = 1.960 \frac{3}{\sqrt{100}} = 0.588$$

Therefore

$$\hat{x} - k = 4.412 \quad \text{and} \quad \hat{x} + k = 5.588$$

and the confidence interval is:

$$\text{CONF}_{0.95} \{4.412 \leq \mu \leq 5.588\} \blacksquare$$

Theory 4.19: Sum of Independent Normal Random Variables

Let X_1, \dots, X_n be independent normal random variables each of which has mean μ and variance σ^2 .

Then the following holds:

- a. The sum $X_1 + \dots + X_n$ is normal with mean $n\mu$ and variance $n\sigma^2$.
- b. The following random variable \bar{X} is normal with mean μ and variance σ^2/n .

$$\bar{X} = \frac{1}{n} (X_1 + \dots + X_n) \quad (4.16)$$

- c. The following random variable Z is normal with mean 0 and variance 1.

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

Exercise 4.7: Sample Size Needed for a Confidence Interval of Prescribed Length

How large must n be in Example 4.3 if we want to obtain a 95% confidence interval of length $L = 0.4$?

Solution

The interval in Eq. (4.15) has the length:

$$L = 2k = 2c\sigma/\sqrt{n}$$

Solving for n , we obtain

$$n = \left(\frac{2c\sigma}{L} \right)^2$$

In the present case the answer is:

$$n = \left(\frac{2 \times 1.96 \times 3}{0.4} \right)^2 \approx 870 \blacksquare$$

For Mean of the Normal Distribution with Unknown Variance

For practical applications, σ^2 is frequently **unknown**. Then the method described previously does **NOT** help and the whole theory changes, although the steps of determining a confidence interval for μ remain quite similar.

We see that k differs from previous method, namely, the sample standard deviation s has taken the place of the unknown standard deviation σ of the population as it is now the variance we are trying to estimate, and c now depends on the sample size n and must be determined from Table XX given in the Appendix. That table lists values z for given values of the distribution function.

$$F(z) = K_m \int_{-\infty}^z \left(1 + \frac{u^2}{m} \right)^{-(m+1)/2} du \quad (4.17)$$

of the t -distribution. Here, $m = 1, 2, \dots$ is a parameter, called the **number of degrees of freedom**

¹⁸abbreviated d.f. of the distribution.¹⁸ In the present case, $m = n - 1$ where n is the number of sample we have to determine variance. The constant K_m is such that $F(\infty) = 1$. By integration it turns out that

$$K_m = \frac{\Gamma\left(\frac{1}{2}m + \frac{1}{2}\right)}{\sqrt{m\pi}\Gamma\left(\frac{1}{2}m\right)},$$

¹⁹Do not worry if these equations do not make sense as it is here for literature purposes.

where Γ is the gamma function.¹⁹

The method of tackling is this problem is as follows:

²⁰95%, 99%, or the like.

1. Choose a confidence level γ .²⁰

2. Determine the solution c of the equation,

$$F(c) = \frac{1}{2}(1 + \gamma)$$

from the table of the t -distribution with $m = n - 1$ degrees of freedom

3. Compute the mean \bar{x} and the variance s^2 of the sample x_1, \dots, x_n .
4. Compute $k = cs/\sqrt{n}$. The confidence interval is:

$$\text{CONF}_\gamma\{\bar{x} - k \leq \mu \leq \bar{x} + k\}.$$

This illustrates that Table XX²¹ provides shorter confidence intervals than Table XX. This is confirmed in, which also gives an idea of the gain by increasing the sample size.

²¹which uses more information, namely, the known value of σ^2

Exercise 4.8: Confidence Interval for Mean of Normal Distribution with Unknown Variance

The five (5) independent measurements of flash point of Diesel oil (D-2) gave the values (in °F):

144 147 146 142 144

If we assume normality, determine a 99% confidence interval for the mean.

Solution

1. $\gamma = 0.99$ is required.
2. $F(c) = \frac{1}{2}(1 + \gamma) = 0.99$ and looking at the reference table with $n - 1 = 4$ d.f., giving $c = 4.60$.
3. $\bar{x} = 144.6$ and $s = 3.8$,

4. $k = \sqrt{3.8} \times 4.60/\sqrt{5} = 4.01$. Therefore the confidence interval is:

$$\text{CONF}_{0.99} \{140.5 \leq \mu \leq 148.7\} \blacksquare$$

If the variance σ^2 were known and equal to the sample variance s^2 , thus $\sigma^2 = 3.8$, then the Reference Table would give:

$$k = \frac{c\sigma}{\sqrt{n}} = 2.576 \frac{\sqrt{3.8}}{\sqrt{3}} = 2.25$$

and

$$\text{CONF}_{0.99} \{140.5 \leq \mu \leq 148.7\}$$

We see that the present interval is almost twice as long as that with a known variance $\sigma^2 = 3.8$.

Theory 4.20: Student's t-Distribution

Let X_1, \dots, X_n be independent normal random variables with the same mean μ and the same variance σ^2 . Then the random variable:

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \quad (4.18)$$

has a t-distribution²² with $n - 1$ degrees of freedom (d.f.); here \bar{X} is given by Eq. (4.16) and

$$S^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2.$$

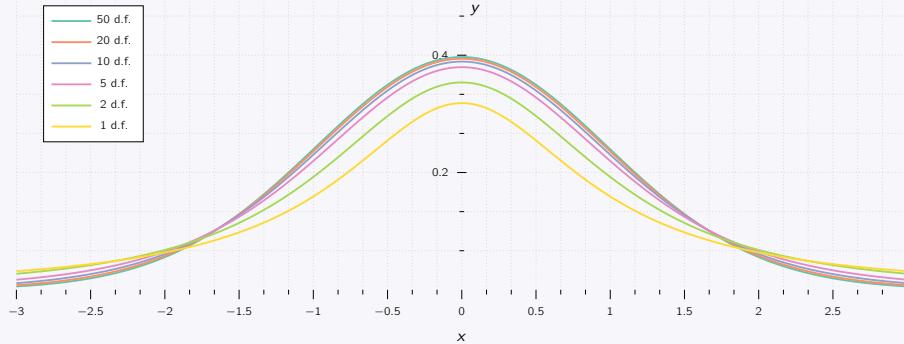


Figure 4.1: The student-t distribution with different degrees of freedom m .



²²William Gosset (13 June 1876 – 16 October 1937) was an English statistician, chemist and brewer who served as Head Brewer of Guinness and Head Experimental Brewer of Guinness and was a pioneer of modern statistics. He published his results under the pen name student.

For the Variance of the Normal Distribution

The method for calculating the confidence interval is similar to the previous methods, with slight change in some steps which are as follows:

²³as usual this can be 95%, 99%, or the like.

1. Choose a confidence level γ .²³
2. Determine solutions c_1 and c_2 of the equations:

$$F(c_1) = \frac{1}{2}(1 - \gamma) \quad \text{and} \quad F(c_2) = \frac{1}{2}(1 + \gamma).$$

where the necessary values are calculated from the table of the chi-square distribution with $n - 1$ degrees of freedom.

3. Compute $(n - 1)s^2$, where s^2 is the variance of the sample x_1, \dots, x_n .
4. Compute $k_1 = (n - 1)s^2/c_1$ and $k_2 = (n - 1)s^2/c_2$. The confidence interval is

$$\text{CONF}_\gamma\{k_2 \cong \sigma^2 \cong k_1\}. \quad (4.19)$$

Exercise 4.9: Confidence Interval for the Variance of the Normal Distribution

Determine a 95% confidence interval Eq. (4.19) for the variance, using Table 25.3 and a sample (tensile strength of sheet steel in kg mm^{-2} , rounded to integer values)

89 84 87 81 89 86 91 90 78 89 87 99 83 89

Solution

1. $\gamma = 0.95$ is required.
 2. For $n - 1 = 13$ we find
- $$c_1 = 5.01 \quad \text{and} \quad c_2 = 24.74.$$
3. $13s^2 = 326.9$
 4. $13s^2/c_1 = 65.25$ and $13s^2/c_2 = 13.21$
 5. This makes the confidence interval as:

$$\text{CONF}_{0.05}\left\{13.21 \leq \sigma^2 \leq 65.25\right\}.$$

This is rather large, and for obtaining a more precise result, one would need a much larger sample ■.

Theory 4.21: Chi-Square Distribution

Under the assumptions in Theorem 2 the random variable

$$Y = (n - 1) \frac{S^2}{\sigma^2}$$

with S^2 given by (12) has a chi-square distribution with $n - 1$ degrees of freedom.

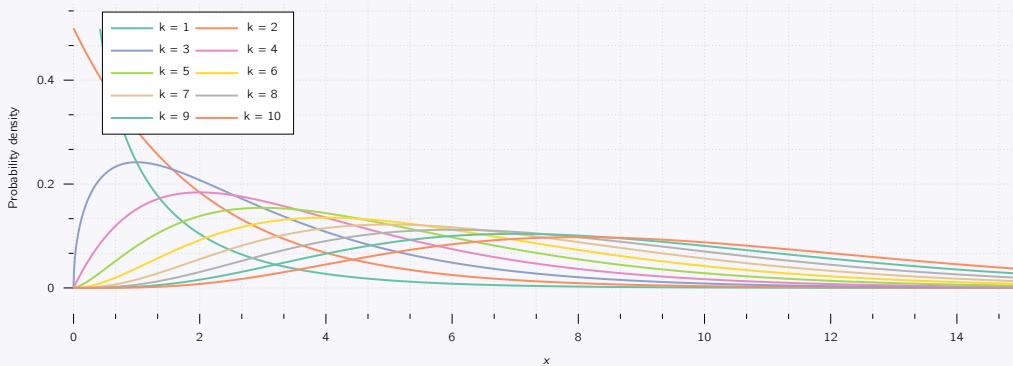


Figure 4.2: Chi-square distribution with different degrees of freedom.

The chi-squared distribution, which can be seen in Fig. 4.2 is used primarily in **hypothesis testing**, and to a lesser extent for confidence intervals for population variance when the underlying distribution is normal. Unlike more widely known distributions such as the normal distribution and the exponential distribution, the chi-squared distribution is not as often applied in the direct modelling of natural phenomena.

The primary reason for which the chi-squared distribution is extensively used in hypothesis testing is its relationship to the normal distribution. Many hypothesis tests use a test statistic, such as the t-statistic in a t-test. For these hypothesis tests, as the sample size n increases, the sampling distribution of the test statistic approaches the normal distribution.²⁴ Because the test statistic (t) is asymptotically normally distributed, provided the sample size is sufficiently large, the distribution used for hypothesis testing may be approximated by a normal distribution.

²⁴This is the result of the central limit theorem.

So wherever a normal distribution could be used for a hypothesis test, a chi-squared distribution could be used.

Confidence Internals for Parameters of Other Distributions

The methods mentioned previously for confidence intervals for μ and σ^2 are designed for the **normal distribution**. We will see it here that they can also be applied to other distributions if we **use large samples**.

We know that if X_1, \dots, X_n are independent random variables with the same mean μ and the same variance σ^2 , then their sum $Y_n = X_1 + \dots + X_n$ has the following properties:

- Y_n has the mean $n\mu$ and the variance $n\sigma^2$,
- If those variables are normal, then Y_n is normal.

If those random variables are **not normal**, then second property is **NOT** applicable. However, for large n the random variable Y_n is still **approximately** normal.

This follows from the **central limit theorem**, which is one of the most fundamental results in probability theory.

Theory 4.22: Central Limit Theorem

Let X_1, \dots, X_n be independent random variables having the same distribution function and therefore the same mean μ and

variance σ^2 . Let $Y_n = X_1 + \dots + X_n$, then the random variable

$$Z_n = \frac{Y_n - n\mu}{\sigma\sqrt{n}}$$

is **asymptotically normal** with mean 0 and variance 1. That is, the distribution function $F(x)$ of Z_n satisfies:

$$\lim_{n \rightarrow \infty} F(x) = \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du.$$

This theorem basically boils down to the following statement:

Under appropriate conditions, the distribution of a normalised version of the sample mean converges to a standard normal distribution. This holds even if the original variables themselves are not normally distributed.

Therefore, when applying the previous confidence interval methods to a **non-normal distribution**, we must use sufficiently large samples.

As a rule of thumb, if the sample indicates that the skewness of the distribution is small, use at least $n = 20$ for the mean and at least $n = 50$ for the variance.

4.4 Testing of Hypotheses and Making Decisions

The ideas of confidence intervals and of tests²⁵ are the two (2) most important ideas in modern statistics. In a statistical test we make inference from sample to population through testing a **hypothesis**, resulting from experience or observations, from a theory or a quality requirement, and so on.

In many cases the result of a test is used as a basis for a **decision**:

to buy, or not to buy a certain model of car, depending on a test of the fuel efficiency (km L^{-1}), or, to apply some medication, depending on a test of its effect; to proceed with a marketing strategy, depending on a test of consumer reactions, etc.

As with most abstract mathematical concepts, it is better to explain such a test in terms of a typical example and then introduce the corresponding standard notions of statistical testing.

Exercise 4.10: Test of a Hypothesis

Let's say we want to buy 100 coils of a certain kind of wire, provided we can verify the manufacturer's claim that the wire has a specific strength of $\mu = \mu_0 = 200 \text{ kN m kg}^{-1}$, or more.

This is a test of the hypothesis:²⁶ $\mu = \mu_0 = 200$. We shall **NOT** buy the wire if the statistical tests shows that actually $\mu = \mu_1 < \mu_0$, the wire is weaker, the claim does **NOT** hold. μ_1 is called the **alternative** of the test.²⁷ We shall **accept** the hypothesis if the test suggests that it is true, except for a small error probability α , called the **significance level** of the test.

Otherwise we reject the hypothesis.

Hence α is the probability of rejecting a hypothesis although it is true. The choice of α is up to us, 5% and 1% are popular values.

For the test we need a sample. We randomly select 25 coils of the wire, cut a piece from each coil, and determine the breaking limit experimentally. Suppose that this sample of $n = 25$ values of the breaking limit has the mean $\bar{x} = 197 \text{ kN m kg}^{-1}$, which is somewhat less than the claim, and the standard deviation $s = 6 \text{ kN m kg}^{-1}$.

At this point we could only speculate when this difference $197 - 200 = -3$ is due to randomness, is a chance effect, or whether it is **significant**, due to the actual inferior quality of the wire. To continue beyond speculation requires probability theory, as follows.

We assume that the blocking limit is normally distributed. Then

$$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

with $\mu = \mu_0$ has a **t-distribution** with $n - 1$ degrees of freedom ($n - 1 = 24$ for our sample). Also $\bar{x} = 197$ and $s = 6$ are observed values of \bar{X} and S to be used later. We can now choose a significance level, say, $\alpha = 95\%$. From the Reference Table, we then obtain a critical value c such that $P(T \leq c) = \alpha = 5\%$. For $P(T \leq \bar{c}) = 1 - \alpha = 95\%$ the table gives $\bar{c} = 1.71$, so that $c = -\bar{c} = -1.71$ because of the symmetry of the distribution shown in Fig. 4.3.

We now reason as follows—this is the crucial idea of the test. If the hypothesis is true, we have a chance of only $\alpha (= 5\%)$ that we observe a value t of T (calculated from a sample) that will fall between $-\infty$ and -1.71 . Hence, if we nevertheless do observe such a t , we start that the hypothesis cannot be true and we reject it.

A simple calculation gives:

$$T = \frac{(107 - 200)}{6/\sqrt{25}} = -2.5,$$

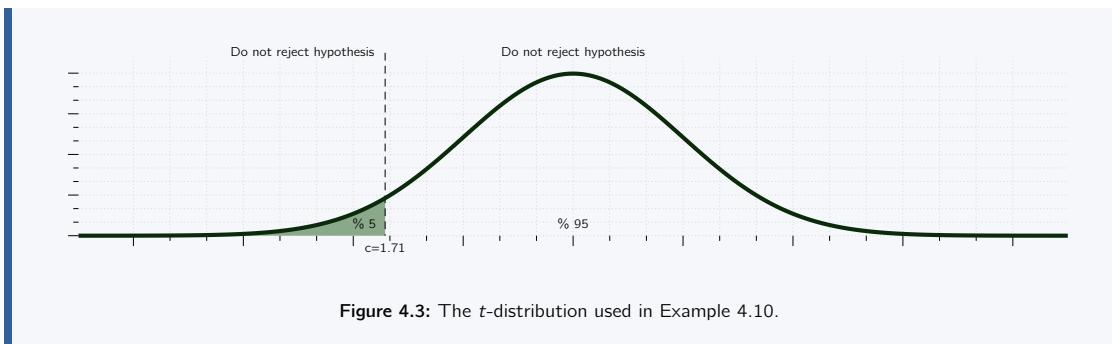
as an observed value of T . Since $-2.5 < -1.71$, we reject the hypothesis, the manufacturer's claim, and accept the alternative result of $\mu = \mu_1 < 200$, which means the wire seems to be weaker than claimed ■



²⁵The modern development of tests are generally attributed to Egon Sharpe Pearson and Neymar whom was mentioned previously. Egon Sharpe was one of three children of Karl Pearson and Maria, and, like his father, a British statistician. He is known throughout the world as co-author of the Neyman-Pearson theory of testing statistical hypotheses, and responsible for many important contributions to problems of statistical inference and methodology, especially in the development and use of the likelihood ratio criterion.

²⁶also called **null hypothesis**.

²⁷or **alternative hypothesis**



This aforementioned example perfectly captures the **steps of a test**:

1. Formulate the **hypothesis** $\theta = \theta_0$ to be tested. In our previous example it is $\theta_0 = \mu_0$.
2. Formulate an **alternative** $\theta = \theta_1$, which in our example is $\theta_1 = \mu_1$.
3. Choose a **significance level** α with values such as 5%, 1%, or, 0.1%.
4. Use a random variable $\hat{\Theta} = g(X_1, \dots, X_n)$ whose distribution depends on the hypothesis and on the alternative, and this distribution is known in both cases.

Determine a critical value c from the distribution of $\hat{\Theta}$, assuming the hypothesis to be true. In the example, $\hat{\Theta} = T$, and c is, obtained from $P(T \leq c) = \alpha$.

5. Use a sample x_1, \dots, x_n to determine an observed value $\hat{\theta} = g(x_1, \dots, x_n)$ of $\hat{\Theta}$, where in our example it is t .
6. Accept or reject the hypothesis, depending on the size of $\hat{\theta}$ relative to c .

There are two (2) important facts require further discussion and careful attention.

1. The choice of an alternative. In the example, $\mu_1 < \mu_0$, but other applications may require $\mu_1 > \mu_0$ or $\mu_1 \neq \mu_0$.
2. Addressing errors. We know that α , the significance level of the test, is the probability of reflecting a **true** hypothesis. And we shall discuss the probability β of accepting a false hypothesis.

One-Sided and Two-Sided Alternatives

Let θ be an **unknown parameter** in a distribution, and suppose we want to test the hypothesis $\theta = \theta_0$.

Then there are three (3) main kinds of alternatives, namely,

$$\theta > \theta_0 \quad (4.20)$$

$$\theta < \theta_0 \quad (4.21)$$

$$\theta \neq \theta_0 \quad (4.22)$$

Here Eq. (4.20), and Eq. (4.21) are **one-sided alternatives**, and Eq. (4.22) is a **two-sided alternative**.

We call rejection region²⁸ the region such that we reject the hypothesis if the observed value in the test falls in this region. In [1] the critical c lies to the right of θ_0 because so does the alternative. Hence the rejection region extends to the right. This is called a **right-sided test**. In [2] the critical c lies to the left of θ_0 (as in Example 1), the rejection region extends to the left, and we have a **left-sided test**. These are one-sided tests. In [3]

²⁸or called the critical region

All three kinds of alternatives occur in practical problems. For example, Eq. (4.20) may arise if θ_0 is the maximum tolerable inaccuracy of a voltmeter or some other instrument. Alternative Eq. (4.21) may occur in testing strength of material, as in Example A. Finally, θ_0 in Eq. (4.22) may be the diameter of axle-shafts, and shafts that are too thin or too thick are equally undesirable, so that we have to watch for deviations in both directions.

4.4.1 Errors in Tests

Tests always involve **risks of making false decisions**:

I Rejecting a true hypothesis (Type I error)

■ α = Probability of making a Type I error.

II Accepting a false hypothesis (Type II error).

■ β = Probability of making a Type II error.

Clearly, we cannot avoid these errors.

No absolutely certain conclusions about populations can be drawn from samples.

But we show there are ways and means of choosing suitable levels of risks, that is, of values α and β . The choice of α depends on the nature of the problem.²⁹

²⁹e.g., a small risk
 $\alpha = 1\%$ is used if it is a matter of life or death.

Let us discuss this systematically for a test of a hypothesis $\theta = \theta_0$ against an alternative that is a single number θ_1 , for simplicity. We let $\theta_1 > \theta_0$, so that we have a **right-sided test**. For a left-sided or a two-sided test the discussion is quite similar.

We choose a critical $c > \theta_0$ ³⁰. From a given sample x_1, \dots, x_n we then compute a value:

$$\hat{\theta} = g(x_1, \dots, x_n)$$

³⁰as in the upper part of Fig. 533, by methods discussed below

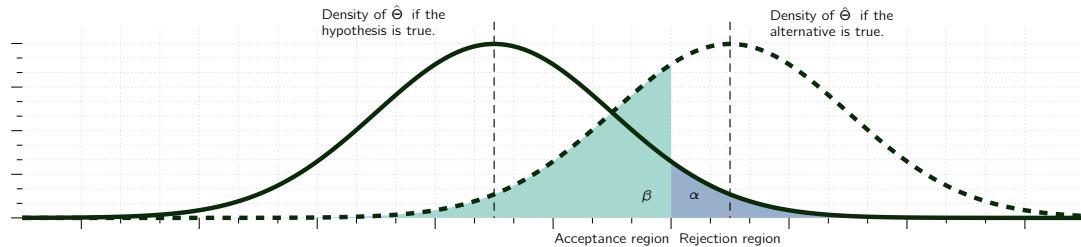


Figure 4.4: Illustration of Type I and II errors in testing a hypothesis $\theta = \theta_0$ against an alternative $\theta = \theta_0$.

with a suitable g .

whose choice will be a main point of our further discussion; for instance, take $g = (x_1 + \dots + x_n)/n$ in the case in which θ is the mean.

If $\hat{\theta} > c$, we **reject the hypothesis**. If $\hat{\theta} \leq c$, we accept it. Here, the value $\hat{\theta}$ can be regarded as an observed value of the random variable

$$\hat{\theta} = g(X_1, \dots, X_n)$$

because x_j may be regarded as an observed value of X_j where $j = 1, \dots, n$. In this test there are two (2) possibilities of making an error, as follows.

³¹hence the alternative is accepted.

Type I Error The hypothesis is true but is rejected³¹ because $\hat{\theta}$ assumes a value $\hat{\theta} > c$. Obviously, the probability of making such an error equals

$$P(\hat{\theta} > c)_{\theta} = \theta_0 = \alpha. \quad (4.23)$$

α is called the **significance level** of the test, as mentioned before.

Type II Error The hypothesis is false but is accepted because $\hat{\theta}$ assumes a value $\hat{\theta} \leq c$. The probability of making such an error is denoted by β ; thus

$$P(\hat{\theta} \leq c)_{\theta=\theta_0} = \beta. \quad (4.24)$$

$\eta = 1 - \beta$ is called the **power** of the test. Obviously, the power η is the probability of avoiding a Type II error.

Formulas Eq. (4.23) and Eq. (4.24) show that both α and β depend on c , and we would like to choose c so that these probabilities of making errors are as small as possible. But the important Fig. 4.4 shows that these are conflicting requirements because to let α decrease we must shift c to the right, but then β increases. In practice we first choose α (5%, sometimes 1%), then determine c , and finally compute β . If β is large so that the power $\eta = 1 - \beta$ is small, we should repeat the test, choosing a larger sample, for reasons that will appear shortly. If the alternative is **NOT** a single

number but is of the form Eq. (4.20)-Eq. (4.22), then β becomes a function of θ . This function $\beta(\theta)$ is called the operating characteristic (OC) of the test and its curve the OC curve. Clearly, in this case $\eta = 1 - \beta$ also depends on θ . This function $\eta(\theta)$ is called the **power function** of the test.

Of course, from a test that leads to the acceptance of a certain hypothesis θ_0 , it does **NOT** follow that this is the only possible hypothesis or the best possible hypothesis. Hence the terms “not reject” or “fail to reject” are perhaps better than the term “accept”.

The following example explains the three (3) kinds of hypotheses.

Exercise 4.11: Test for the Mean of the Normal Distribution with Known Variance

Let X be a normal random variable with variance $\sigma^2 = 9$. Using a sample of size $n = 10$ with mean \bar{x} , test the hypothesis $\mu = \mu_0 = 24$ against the three (3) kinds of alternatives, namely,

$$(a) \mu > \mu_0 \quad (b) \mu < \mu_0 \quad (c) \mu \neq \mu_0$$

Solution

We choose the significance level $\alpha = 0.05$. An estimate of the mean will be obtained from:

$$\bar{X} = \frac{1}{n} (X_1 + \dots + X_n).$$

If the hypothesis is true, \bar{X} is normal with mean $\mu = 24$ and variance $\sigma^2/n = 0.9$. Hence we may obtain the critical value c from Table 48 in App. 5.

a. **Right-Sided Test** We determine c from

$$P(\bar{X} > c)_{\mu=24} = \alpha = 0.05$$

that is,

$$P(\bar{X} \leq c)_{\mu=24} = \Phi\left(\frac{c - 24}{\sqrt{0.9}}\right) = 1 - \alpha = 0.95.$$

Using Table A8 in App. 5 gives $(c - 24)/\sqrt{0.9} = 1.645$, and $c = 25.56$, which is greater than μ_0 . If $\bar{x} \leq 25.56$, the hypothesis is **accepted**. If $\bar{x} > 25.56$, it is rejected.

The power function of the test is:

$$\begin{aligned} \eta(\mu) &= P(\bar{X} > 25.56)_{\mu} = 1 - P(\bar{X} \leq 25.56)_{\mu} \\ &= 1 - \Phi\left(\frac{25.56 - \mu}{\sqrt{0.9}}\right) = 1 - \Phi(26.94 - 1.05\mu) \end{aligned}$$

b. **Left-Sided Test** The critical value c is obtained from the equation

$$P(\bar{X} \leq c)_{\mu=24} = \Phi\left(\frac{c - 24}{\sqrt{0.9}}\right) = \alpha = 0.05.$$

Table A8 in App. 5 yields $c = 24 - 1.56 = 22.44$. If $\bar{x} \geq 22.44$, we accept the hypothesis. If $\bar{x} < 22.44$, we reject it. The power function of the test is

$$\eta(\mu) = P(\bar{x} \geq 22.44)_{\mu} = \Phi\left(\frac{22.44 - \mu}{\sqrt{0.9}}\right) = \Phi(23.65 - 1.05\mu).$$

c. **Two-Sided Test** Since the normal distribution is symmetric, we choose c_1 and c_2 equidistant from $\mu = 24$, say, $c_1 = 24 - k$ and $c_2 = 24 + k$, and determine k from

$$P(24 - k \leq \bar{X} \leq 24 + k)_{\mu=24} = \Phi\left(\frac{k}{\sqrt{0.9}}\right) - \Phi\left(-\frac{k}{\sqrt{0.9}}\right) = 1 - \alpha = 0.95.$$

Table A8 in App. 5 gives $k/\sqrt{0.9} = 1.960$, hence $k = 1.86$. This gives the values $c_1 = 24 - 1.86 = 22.14$ and $c_2 = 24 + 1.86 = 25.86$. If \bar{x} is not smaller than c_1 and not greater than c_2 , we accept the hypothesis. Otherwise, we reject it. The power function of the test is (Fig. 535)

$$\eta(\mu) = P(\bar{x} < 22.14)_{\mu} + P(\bar{x} > 25.86)_{\mu} = P(\bar{x} < 22.14)_{\mu} + 1 - P(\bar{x} \leq 25.86)_{\mu}$$

$$\begin{aligned}
 &= 1 + \Phi\left(\frac{22.14 - \mu}{\sqrt{0.5}}\right) - \Phi\left(\frac{25.86 - \mu}{\sqrt{0.5}}\right) \\
 &= 1 + \Phi(23.34 - 1.05\mu) - \Phi(27.26 - 1.05\mu).
 \end{aligned}$$

Consequently, the operating characteristic $\beta(\mu) = 1 - \eta(\mu)$ (see before) is (Fig. 536).

Exercise 4.12: Comparison of the Means of Two Normal Distributions

Using a sample x_1, \dots, x_m from a normal distribution with unknown mean μ_x and a sample y_1, \dots, y_m from another normal distribution with unknown mean μ_y , we want to test the hypothesis that the means are equal, $\mu_x = \mu_y$, again an alternative, say, $\mu_x > \mu_y$. The variances need not be known but are assumed to be equal.

Solution

Two cases of comparing means are of practical importance:

Case A

The samples have the same size. Furthermore, each value of the first sample corresponds to precisely one value of the other, because corresponding values result from the same person or thing (**paired comparison**)

for example, two measurements of the same thing by two different methods or two measurements from the two cycles of the same person.

More generally, they may result from a circular individual or things, for example, the lower reason is that the lower reason is that the lower value of the second sum of the differences of corresponding values used as the previous that the population corresponds to the differences has mean 0. using the method in Example 3. If we have a choice, this method is better than the following.

Case B

The two samples are independent and not necessarily of the same size. Then we may proceed as follows. Suppose that the alternative is $\frac{1}{2} u^n > u_{PV}$ we choose a significance level α . Then we compute the sample means \bar{x} and \bar{y} as well as $(n_1 - 1)u_{PV}^2$ and $(n_2 - 1)u_{PV}^2$ where \bar{x}^2 and \bar{y}^2 are the sample variances. Using Table 9 in App.5 with $n_1 + n_2 - 2$ degrees of freedom, we now determine c from

4.5 Goodness of Fit

To test for goodness of fit³² means that we wish to test that a certain function $F(x)$ is the distribution function of a distribution from which we have a sample x_1, \dots, x_n . Then we test whether the **sample distribution function** $\tilde{F}(x)$ defined as:

$$\tilde{F}(x) = \text{Sum of the relative frequencies of all sample values } x_j \text{ not exceeding } x, \quad (4.25)$$

fits $\tilde{F}(x)$ **sufficiently well**. If this is so, we shall **accept** the hypothesis that $\tilde{F}(x)$ is the distribution function of the population; else, we shall **reject the hypothesis**.

This test is of considerable practical importance, and it differs in character from the tests for parameters (μ, σ^2 , etc.) considered thus far.

To test in that fashion, we have to know how much $\tilde{F}(x)$ can differ from $F(x)$ if the hypothesis is **true**. Hence we must first introduce a quantity which measures the deviation of $\tilde{F}(x)$ from $F(x)$, and we must know the probability distribution of this quantity under the assumption that the hypothesis is true.

Then we proceed as follows.

We determine a number, lets use c , such that, if the hypothesis is **true**, a deviation greater than c has a small preassigned probability. If, nevertheless, a deviation greater than c occurs, we have reason to doubt that the hypothesis is true and we reject it. On the other hand, if the deviation does not exceed c , so that $\tilde{F}(x)$ approximates $F(x)$ sufficiently well, we accept the hypothesis. Of course, if we accept the hypothesis, this means that we have insufficient evidence to reject it, and this does not exclude the possibility that there are other functions that would not be rejected in the test.

In this respect the situation is quite similar to hypothesis testing we talked previously.

The following text-block shows a test of that type, which was introduced by *R. A. Fisher*. This test is justified by the fact that if the hypothesis is true, then χ^2_0 is an observed value of a random variable whose distribution function approaches that of the chi-square distribution with $K - 1$ degrees of freedom³³ as n approaches infinity. The requirement that at least five (5) sample values lie in each interval results from the fact that for finite n that random variable has only approximately a chi-square distribution.

If the sample is so small that the requirement cannot be satisfied, one may continue with the test, but then use the result **with caution**.

³²In literature, this method also means χ^2 -Test.

Historical Anecdote

During the 19th century, statistical analytical methods were mainly applied to biological data and it was customary for researchers to assume observations followed a **normal distribution**, such as Sir George Airy and Mansfield Merriman, whose works were criticized by Karl Pearson in his 1900 paper.

At the end of the 19th century, Pearson noticed the existence of significant skewness within some biological observations. To model the observations regardless of being normal or skewed, Pearson, in a series of articles published from 1893 to 1916,[3][4][5][6] devised the Pearson distribution, a family of continuous probability distributions, which includes the normal distribution and many skewed distributions, and proposed a method of statistical analysis consisting of using the Pearson distribution to model the observation and performing a test of goodness of fit to determine how well the model really fits to the observations.

³³or $K - r - 1$ degrees of freedom if r parameters are estimated.

Chi-square Test for $F(x)$ being the Distribution Function of a Population

- Subdivide the x -axis into n intervals I_1, \dots, I_n such that each interval contains at least five (5) values of the given sample x_1, \dots, x_n .

Determine the number b_j of sample values in the interval I_j , where $j = 1, \dots, K$. If a sample value lies at a common boundary point of two (2) intervals, add 0.5 to each of the two (2) corresponding b_j .

- Using $F(x)$, calculate the probability p_j that the random variable X under consideration assumes any value in the interval I_j , where $j = 1, \dots, K$. Then, calculate

$$e_j = np_j.$$

This is the number of sample values **theoretically expected** in I_j if the hypothesis is true.

- Compute the deviation:

$$\chi^2_0 = \sum_{j=1}^K \frac{(b_j - e_j)^2}{e_j}.$$

- Choose a significance level such as 5%, 1%, or the like.

- Determine the solution c of the equation

$$P(\chi^2 \leq c) = 1 - \alpha.$$

from the table of the chi-square distribution with $K - 1$ degrees of freedom (Table A10 in App. 5).

If r parameters of $F(x)$ are unknown and their maximum likelihood estimates are used, then use $K - r - 1$ degrees of freedom, instead of $K - 1$.

If $\chi^2_0 \leq c$, accept the hypothesis. If $\chi^2_0 > c$, reject the hypothesis.

Exercise 4.13: Test of Normality

Test whether the population from which the sample in table given below was taken is normal.

320	380	340	410	380	340	360	350	320	370
350	340	350	360	370	350	380	370	300	420
370	390	390	440	330	390	330	360	400	370
320	350	360	340	340	350	350	390	380	340
400	360	350	390	400	350	360	340	370	420
420	400	350	370	330	320	390	380	400	370
390	330	360	380	350	330	360	300	360	360
360	390	350	370	370	350	390	370	370	340
370	400	360	350	380	380	360	340	330	370
340	360	390	400	370	410	360	400	340	360

Solution

The table given in the question shows the values, column by column, in the order obtained in the experiment. The next table gives the frequency distribution and Fig. 542 the histogram.

The maximum likelihood estimates for μ and σ^2 are $\hat{\mu} = \bar{x} = 364.7$ and $\hat{\sigma}^2 = 712.9$. The computation in Table 25.10 yields $\bar{x}_0^2 = 2.688$. It is very interesting that the interval $375 \dots 385$ contributes over 50% of \bar{x}_0^2 . From the histogram we see that the corresponding frequency looks much too small. The second largest contribution comes from $395 \dots 405$, and the histogram shows that the frequency seems somewhat too large, which is perhaps not obvious from inspection.

Tensile Strength	Absolute Freq.	Relative Freq.	Cumulative Absolute Freq.	Cumulative Relative Freq.
300	2	0.02	2	0.02
310	0	0.00	2	0.02
320	4	0.04	6	0.06
330	6	0.06	12	0.12
340	11	0.11	23	0.23
350	14	0.14	37	0.37
360	16	0.16	53	0.53
370	15	0.15	68	0.68
380	8	0.08	76	0.76
390	10	0.10	86	0.86
400	8	0.08	94	0.94
410	2	0.02	96	0.96
420	3	0.03	99	0.99
430	0	0.00	99	0.99
440	1	0.01	100	1.00

Table 4.2: Frequency table of the sample given in the question.

We choose $\alpha = 5\%$. Since $K = 10$ and we estimated $r = 2$ parameters we have to use Table A10 in App. 5 with $K - r - 1 = 7$ degrees of freedom. We find $c = 14.07$ as the solution of $P(\chi^2 \geq c) = 95\%$. Since $\bar{x}_0^2 < c$, we accept the hypothesis that the population is normal.

Part II

Localisation and Mapping

Part III

GNU/Linux Operating System

Chapter 5

Welcome to Linux

Table of Contents

5.1 Learning the Linux Command Line	129
-----------------------------------------------	-----

5.1 Learning the Linux Command Line

Working with a text-based Command Line Environment (CLI), without the graphical user interface can be intimidating at first glance, as most are accustomed to using a graphical user interface (GUI). But understanding the command line environment will show how powerful and efficient it is.

Most senior programmers in the industry and veteran Linux system administrators will exclusively use CLI as their day to day interaction with the computer. The reason is, the GUI being designed for simplifying human interaction with computers rather than improving the computer's efficiency.

This document aims to introduce the fundamentals of working with the **Linux command line** using a very common shell called Bash as it will be important in the future when working with ROS (Robot Operating System) or in any future endeavour the reader may pursue in the fields related to computer science.

- Work on what the command line is and how it works,
- Look at working with files and folders,
- How Linux protects files from unauthorised access with permissions,
- Common commands to be familiar with and how to connect commands together with pipes,
- Introduction to some complex command line tasks.

This document aims to give practical knowledge on working with the widely used Bash shell, in case you choose to extend your learning into user management, network configuration, programming and development, or system administration.

5.1.1 A History of Command Line Interface

The command-line interface came from a form of dialogue once conducted by humans over teleprinter (TTY) machines, in which human operators remotely exchanged information instead of a human communicating with another human over a teleprinter. Early computer systems often used teleprinter machines as the means of interaction with a human operator. The computer became one end of the human-to-human teleprinter model.

The mechanical teleprinter was replaced by a terminal, a keyboard and screen emulating the teleprinter. "Smart" terminals permitted additional functions, such as cursor movement over the entire screen, or local editing of data on the terminal for transmission to the computer.

As the microcomputer revolution replaced the traditional systems, hardware terminals were replaced by terminal emulators - PC software that interpreted terminal signals sent through the PC's serial ports. These were typically used to interface an organization's new PC's with their existing mini- or mainframe computers, or to connect PC to PC. Some of these PCs were running Bulletin Board System software.

Early operating system CLIs were implemented as part of resident monitor programs, and could not easily be replaced. The first implementation of the shell as a replaceable component was part of the Multics time-sharing operating system. In 1964, MIT Computation Center staff member Louis Pouzin developed the RUNCOM tool for executing command scripts while allowing argument substitution.

Pouzin coined the term "shell" to describe the technique of using commands like a programming language, and wrote a paper about how to implement the idea in the Multics operating system. Pouzin returned to his native France in 1965, and the first Multics shell was developed by Glenda Schroeder. At Nokia Bell Labs headquarters the first Unix shell, the V6 shell, was developed by Ken Thompson in 1971 and was modelled after Schroeder's Multics shell. The Bourne shell was introduced in 1977 as a replacement for the V6 shell. Although it is used as an interactive command interpreter, it was also intended as a scripting language and contains most of the features that are commonly considered to produce structured programs.

The Bourne shell led to the development of the KornShell (ksh), Almquist shell (ash), and the popular Bourne-again shell (or Bash). Early microcomputers themselves were based on a command-line interface such as CP/M , DOS or AppleSoft BASIC. During the 1980s and 1990s, the introduction of the Apple Macintosh and of Microsoft Windows on PCs saw the command line interface as the primary user interface replaced by the Graphical User Interface. The command line remained available as an alternative user interface, often used by system administrators and other advanced

users for system administration, computer programming and batch processing.

```

-rwxr-xr-x 1 bin      18296 Jun  8 1979 fsck
-rwxr-xr-x 1 bin      1458 Jun  8 1979 getty
-rw-r--r-- 1 root     49 Jun  8 1979 group
-rwxr-xr-x 1 bin      2482 Jun  8 1979 init
-rwxr-xr-x 1 bin      8484 Jun  8 1979 mkfs
-rwxr-xr-x 1 bin      3642 Jun  8 1979 mknod
-rwxr-xr-x 1 bin      3976 Jun  8 1979 mount
-rw-r--r-- 1 root     141 Jun  8 1979 passwd
-rw-r--r-- 1 bin      366 Jun  8 1979 rc
-rw-r--r-- 1 bin      266 Jun  8 1979 ttys
-rwxr-xr-x 1 bin      3794 Jun  8 1979 umount
-rwxr-xr-x 1 bin      634 Jun  8 1979 update
-rw-r--r-- 1 bin      40 Sep 22 05:49 utmp
-rwxr-xr-x 1 root     4520 Jun  8 1979 wall
# ls -l /*unix*
-rwxr-xr-x 1 sys      53302 Jun  8 1979 /hptunix
-rwxr-xr-x 1 sys      52850 Jun  8 1979 /hptmunix
-rwxr-xr-x 1 root     50990 Jun  8 1979 /rkunix
-rwxr-xr-x 1 root     51982 Jun  8 1979 /rl2unix
-rwxr-xr-x 1 sys      51790 Jun  8 1979 /rphptunix
-rwxr-xr-x 1 sys      51274 Jun  8 1979 /rptmunix
# ls -l /bin/sh
-rwxr-xr-x 1 bin      17310 Jun  8 1979 /bin/sh
#

```

Figure 5.1: Bourne shell interaction on Version 7 Unix (Original).

Shells in other Operating Systems

Windows

In November 2006, Microsoft released version 1.0 of Windows PowerShell, which combined features of traditional Unix shells with their proprietary object-oriented .NET Framework. MinGW and Cygwin are open-source packages for Windows that offer a Unix-like CLI. Microsoft provides MKS Inc.'s ksh implementation MKS Korn shell for Windows through their Services for UNIX add-on.

Macintosh

Since 2001, the Macintosh operating system macOS has been based on a Unix-like operating system called Darwin. On these computers, users can access a Unix-like command-line interface by running the terminal emulator program called Terminal, which is found in the Utilities sub-folder of the Applications folder, or by remotely logging into the machine using ssh. Z shell is the default shell for macOS (as of macOS Catalina) with bash, tcsh, and the KornShell also provided. Before macOS Catalina, bash was the default.

5.1.2 Linux is a Nutshell

A Brief Description of What Linux Does

Linux is a general purpose computer operating system, originally released in 1991 by Linus Torvalds and began as a personal project of him [50]. It was to create a new **free** operating system kernel which the resulting kernel has been marked by constant growth throughout its history¹.

¹ There were alternative OSs on the market such as MINIX but it was under a proprietary license which was later became open-source in 2000

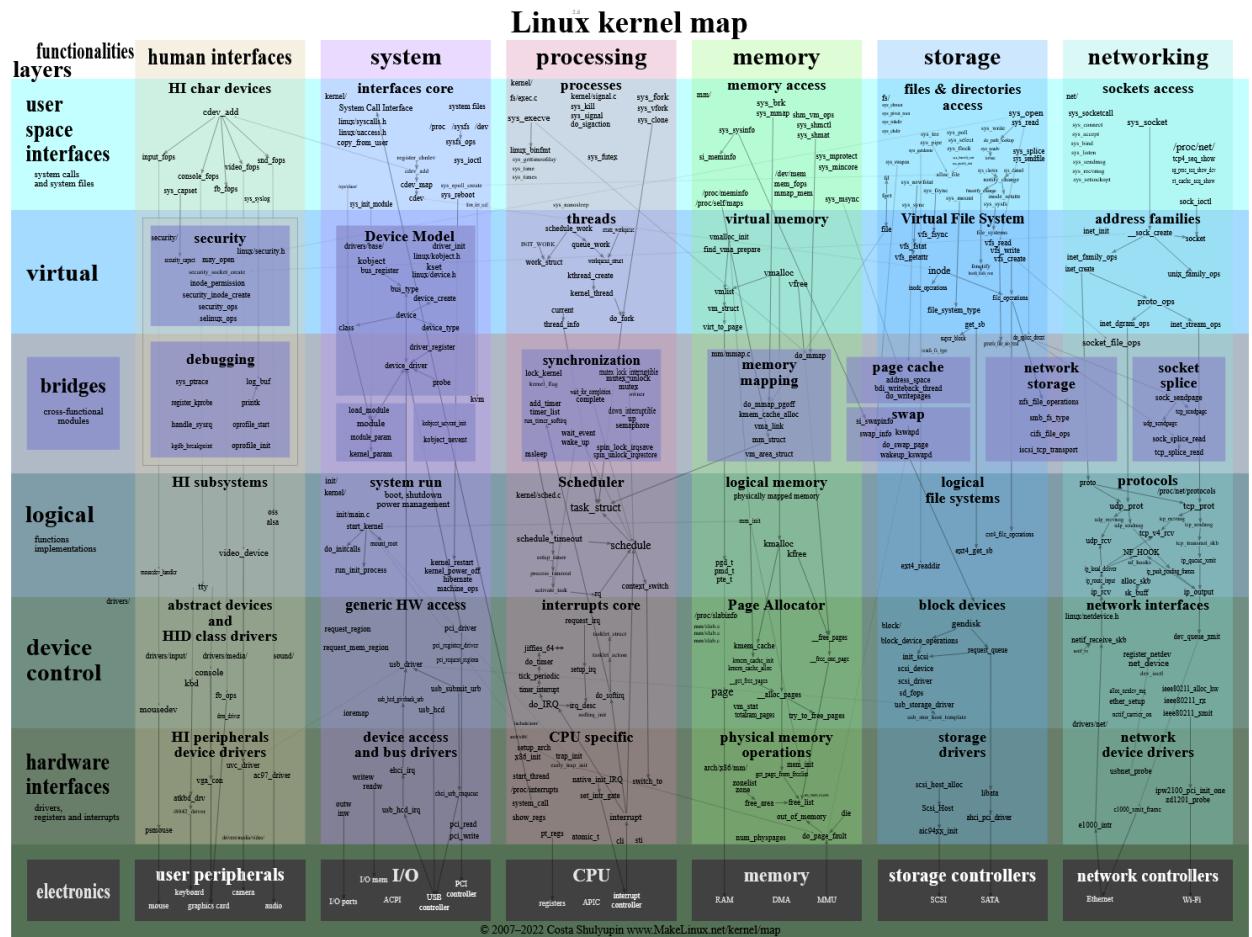


Figure 5.2: The kernel mapping of the Linux operating system.

Linux is defined by its kernel, called the **Linux kernel**, which is the core component of the system. This kernel interacts with the computer hardware to allow software and other hardware to exchange information, which you can see in **Fig. 5.2**.

Imagine the kernel as the middle-man between your software and the hardware

As linux is an open-source project and is probably one of the greatest collaborative software work in history, it has a rich history. It was inspired by MINIX which, in turn, was inspired by UNIX with UNIX being the first portable operating system ever designed [51] as it was mostly written with the C programming language [52].

Open Source v. Closed Source

In programming there are two (2) main approaches when it comes to sharing code:

- It can be closed source, which means, you are not allowed to edit the code the program is

running on,

- Open source which you are free to edit and share the code as you see fit.

Linux is based on a philosophy of software and operating systems being **free**; free of cost and freely modifiable. The software license which allows this, in the case of the Linux kernel, is called the GNU General Public License². This emphasis of freedom, both, of cost and modification has helped Linux to become popular for many different applications and purposes from tinkering programming to being used in massive databases of major companies.

²are a series of widely used free software licenses, or copyleft licenses, that guarantee end users the freedoms to run, study, share, or modify the software.

Linux has popped up everywhere from the majority of the servers that run web services we all use, to super computers, to Wi-Fi routers, in cars, mobile phones, and everywhere in between. Odds are that you are closed to a device that uses some part of the Linux kernel. In the midst of all these different kinds of Linux installations, the most important distinction you'll need to be aware of is one of the genealogy of Linux.

5.1.3 Linux Distributions

While the Linux kernel is more or less the same across nearly all installations of Linux, the software that surrounds the kernel that provides capabilities like *software package management*, *control of services*, and the *location of configuration files* differs between them. Many of the tools that come packaged with Linux come from the GNU Project and aren't actually a part of Linux and, taken together, the combination of the kernel and these common tools is often referred to as **GNU Linux**. Different groups of software and configuration choices that are maintained by individuals or groups of people are called distributions, or distro's. Most major distributions of Linux fall into categories based on the original distribution from which they were derived. These are:

Distribution	Advantages
Linux Mint	Superb collection of "minty" tools developed in-house, hundreds of user-friendly enhancements, inclusion of multimedia codecs, and a wide range of desktop environments.
Ubuntu	Fixed release cycle and support period; long-term support (LTS) variants with five years of security updates; novice-friendly; and a large community of users.
Arch Linux	Excellent software management infrastructure; unparalleled customisation and tweaking options; superb on-line documentation and a strong focus on system security.
Gentoo	Highly flexible, endlessly customizable, able to use a range of compile-time configurations, init systems and run on many architectures.
Slackware Linux	Considered highly stable, clean and largely bug-free, strong adherence to UNIX principles.
Debian	Very stable; remarkable quality control; includes over 30,000 software packages; supports more processor architectures than most distributions.
Fedora	Highly innovative; outstanding security features; large number of supported packages; strict adherence to the free software philosophy.
openSUSE	Comprehensive and intuitive configuration tool; large repository of software packages, excellent web site infrastructure and package management system.
Red Hat	Long-term, commercial support of ten years or more. Stability.
FreeBSD	Fast and stable; availability of over 24,000 software applications (or "ports") for installation; very good documentation; native support for many hardware platforms.

Table 5.1: Most popular distributions used according to distrowatch.com

Depending the readers future work or study area, it is likely to end up learning to use the command line on a system that inherits from one of these distributions. Most likely, it will be a distribution derived from Debian or Red Hat. Linux Mint, Ubuntu, Elementary OS, and Kali Linux are all derived from Debian. CentOS, Fedora, and Red Hat Enterprise Linux are derived from Red Hat.

The history of all of these different distributions of Linux is beyond the scope of this document. But,

what this means at its core is the need to be aware of what system is in use and the need to adapt to account for differences in distributions. As we begin working with Linux, through the command line, it will be apparent, most of what can be done is the same across the major distributions.

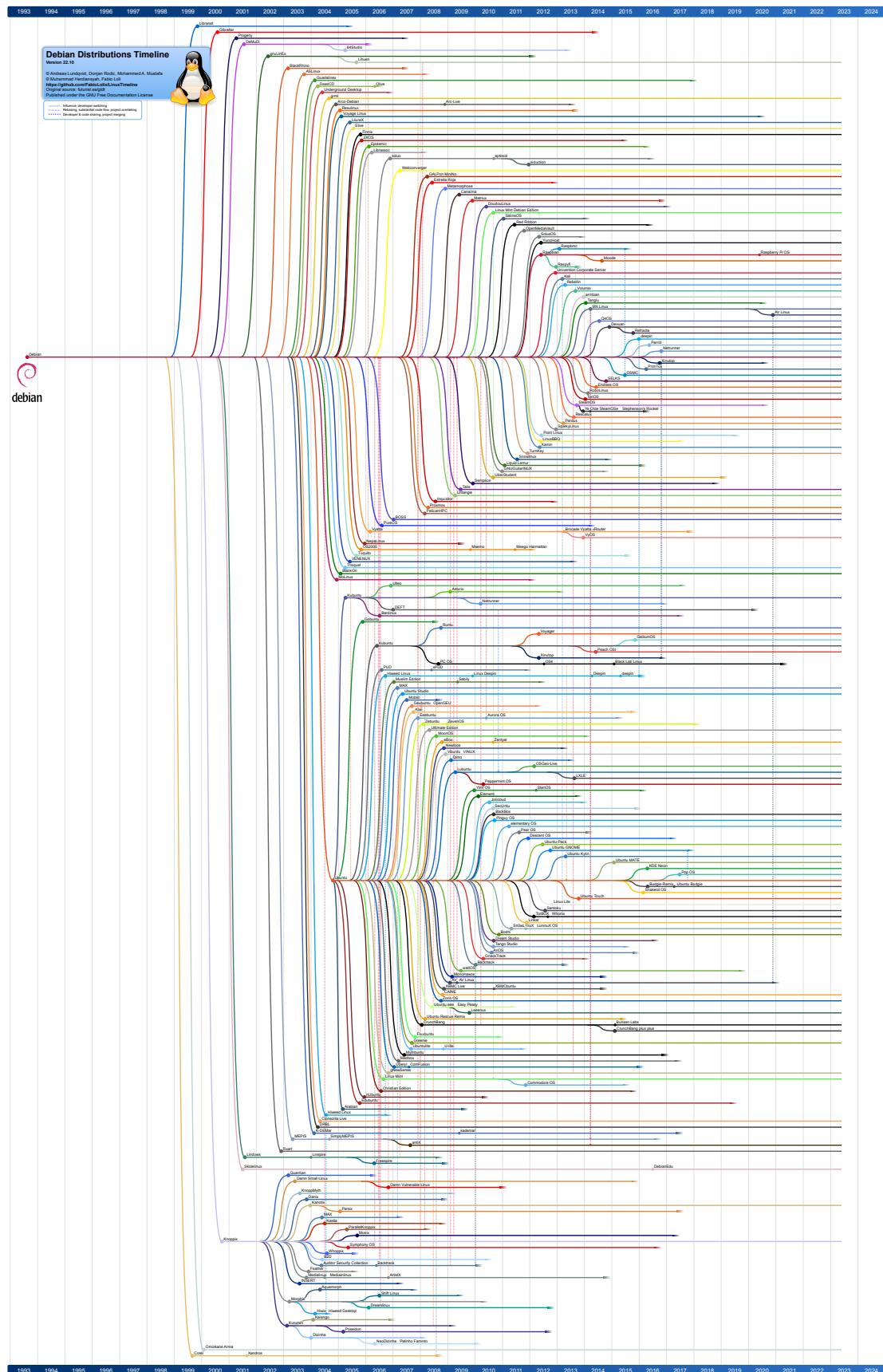


Figure 5.3: A Family tree of the debian branch of linux [53].

Chapter 6

Command Line Fundamentals

Table of Contents

6.1	Introduction	137
6.2	The Structure of Commands	140
6.3	Helpful Keyboard Shortcuts for the Terminal	143
6.4	When you need help with Commands	145
6.5	Additional Information	149

6.1 Introduction

In this day and age, it takes a certain level of skill to be alien to technology and as everyone has computers in their pockets, the interactions with them is almost uncountable. However, these interactions are done through what is called a Graphical User Interface (GUI). Devices running on Windows, MacOS, iOS, and Android all use this interface to interact with the user.

i.e., when clicked on an icon, the close, minimize and maximize buttons on the windows etc..

It must be stressed as these visual components are all for the benefit of the user. While these greatly simplify tasks like photo editing or video creation, some applications just completely omit the use of GUI and instead use a simpler version of it called Command-Line Interface (CLI). This is especially true for servers¹, embedded applications and in many other areas where either **memory is limited** or efficiency of the computer (e.g., such as limiting the CPU load etc.) is highly desired. Server software, utilities, and other programs usually only need some text-based information to do what they do. Many of these programs run on a server in a data centre somewhere without a monitor so the overhead of a GUI is completely **unnecessary**.

¹A server is a software or hardware offering a service to a user, usually referred to as client. As an example, a hardware server is a shared computer on a network, usually powerful and housed in a data centre.

UNIX	Windows
Bourne shell (sh)	COMMAND.COM, default in Windows 9x and provided for DOS compatibility in 32-bit versions of NT-based Windows via NTVDM.
Almquist shell (ash)	
Debian Almquist shell (dash)	
Bash (Unix shell) (bash)	
Korn shell (ksh)	<code>cmd.exe</code> , the default command-line interpreter of the Windows NT-family
Z Shell (zsh)	
C shell (csh)	Recovery Console
TENEX C shell (tcsh)	Windows PowerShell, based on .NET Framework
Ch shell (ch)	PowerShell, based on .NET Core
Emacs shell (eshell)	Hamilton C shell, a clone of the Unix C shell
Friendly interactive shell (fish)	
Powershell (pwsh)	4NT, a clone of CMD.EXE.
rc shell (rc)	Take Command, a newer incarnation of 4NT
Stand-alone shell (sash)	
Scheme Shell (scsh)	

Table 6.1: Types of shells used in industry and academia. For reference, the authors computer uses zsh.

One way we interact with these programs that don't have a GUI is through the CLI. This is a text-based interface where the commands to execute are typed and all actions are shown as text on a terminal screen, whether it is updating a software or moving files around. The environment we use is called a shell, or command-line interpreter, and there are many shells out there.

A list of Shells that can be encountered in industry and academia can be seen in **Table 6.1**.

The command-line interpreter was one of the earliest ways of interacting with the general-purpose computer, starting in 1971 with the Thompson shell for UNIX². As UNIX evolved and came to be replaced in many capacities by Linux, the shell environments evolved and improved as well.

Bash, or the Bourne-again shell is one of the most widely-used shells and odds are, it's the one to be encountered in industry or in academic work. Bash is the shell that comes enabled by default with most of the popular Linux distributions. It's also available on macOS³ and in Windows with the Windows subsystem for Linux.

²A family of multitasking, multi-user computer operating systems that derive from the original AT&T Unix, whose development started in 1969. It is considered one of the most groundbreaking software ever designed.

³newer versions have `zsh` shell instead but they are designed to be compatible.

The author of this work also uses `zsh` as his main driver.

In this document, Bash will be used. However, the reader is encouraged to explore some of the other shells out there once a working foundation in Bash is achieved.

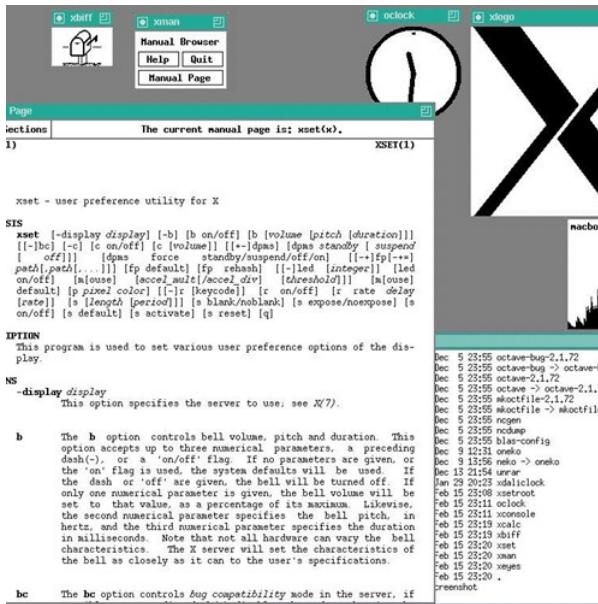


Figure 6.1: A graphical interface from the late 1980s, which features a TUI window for a man page, a shaped window (oclock) as well as several iconified windows. In the lower right we can see a terminal emulator running a Unix shell, in which the user can type commands as if they were sitting at a terminal. - *From Wikipedia*

6.2 The Structure of Commands

There are a few concepts and principles which needs to be understood to be a productive member of the CLI family. Before jumping into using commands though, have a look at how command line statements are structured with the following:

```
1  command [-flag(s)] [-option(s) [value]] [argument(s)]
```

C.R. 1

bash

This is the general form. The pattern is **command**, **options**, and then **arguments**. Here's a couple of common commands you'll see with options and arguments that are used with them.

```
1  ls -l /tmp
2  cd /usr/local
3  cat /etc/passwd
```

C.R. 2

bash

The details of what the aforementioned commands do will be the focus in the future. I just want to show you the structure of what we'll be working with before we get into what these actually do. Depending on the current action, you might just have a command or a command and one or more options or just a command with one or more arguments.

But there will always be a command.

Command is the **minimum** thing which can be done with a CLI. Think of it as the atom of any action you can take. The command is the program you're running or the action you're taking. To give

command to a UNIX system, type the name of the command, along with any associated information, such as a filename, and press the `\return` key.

The typed line is called the command line and UNIX uses a special program, called the shell or the command line interpreter, mentioned in the previous section, to interpret what you have typed into what you want to do.

The components of the command line are:

1. the command,
2. any options required by the command,
3. the command's arguments.⁴

⁴This is optional as some commands just don't have any options.

6.2.1 Some Rules Regarding the Syntax

Since the introduction of UNIX System V, Release 3 (released 1983), any new commands must obey a particular syntax governed by the following rules:

- Command names must be between 2 and 9 characters in length,
- Command names must be comprised of lowercase characters and digits,
- Option names must be one character in length,
- All options are preceded by a hyphen (-),
- Options without arguments may be grouped after the hyphen,
- The first option argument, following an option, must be preceded by white space,
i.e., `-o sfile` is valid but `-osfile` is **illegal**.
- Option arguments are **not optional**,
- If an option takes more than one argument then they must be separated by commas with no spaces, or if spaces are used the string must be included in double quotes (").
i.e., both are acceptable: `-f past,now,next` and `-f "past now next"`.
- All options must precede other arguments on the command line,
- A double hyphen -- may be used to indicate the end of the option list,
- The order of the options are order independent,

- The order of arguments may be important,
- A single hyphen – is used to mean standard input.

Options **must** come after the command and before arguments. Options **should not** appear after the main argument(s). However, some options can have their own arguments! Historically, UNIX commands have been fairly standard in the way that they use options but there are variations.

Bear in mind that commands established before System V, Release 3, do not conform to all of the above rules.

6.3 Helpful Keyboard Shortcuts for the Terminal

Before we moving on to more specific commands and get into CLI programming, there's a few other helpful things to know about working at the Command Line. The first is **Tab completion**, a wonderful feature of the Bash shell, and is also included in many others. This feature let's you skip typing out a whole file name or folder name when you're working at the Command Line.

When you're working in the command line it looks at all the information it has so far and makes a guess about what you mean. For example, I can type `ls -l De` and press `\tab`, and it completes the line with Desktop. Now type `ls -l Do` and nothing would happen when I press `\tab`. That's because `\tab` doesn't have one clear suggestion to return. As it can be either Documents or Downloads. However, pressing `\tab` again should give you a suggestion of which items can be completed to.

For reference, in the following page, there is a table for most useful keyboard shortcut for Linux Bash.

	Shortcut	Action
Navigation	<code>Ctrl + A</code>	Go to the beginning of the line.
	<code>Ctrl + E</code>	Go to the end of the line.
	<code>Alt + F</code>	Move the cursor forward one word.
	<code>Alt + B</code>	Move the cursor back one word.
	<code>Ctrl + F</code>	Move the cursor forward one character.
	<code>Ctrl + B</code>	Move the cursor back one character.
	<code>Ctrl + X</code>	Toggle between the current cursor position and the beginning of the line.
Editing	<code>Ctrl + A</code>	Undo! (That's an underscore, so you'll need to use <code>Shift</code> as well.).
	<code>Ctrl + X</code>	Edit the current command in your \$EDITOR.
	<code>Alt + D</code>	Delete the word after the cursor.
	<code>Alt</code>	Delete the word before the cursor.
	<code>Ctrl + D</code>	Delete the character beneath the cursor.
	<code>Ctrl + H</code>	Delete the character before the cursor (like backspace).
	<code>Ctrl + K</code>	Cut the line after the cursor to the clipboard.
	<code>Ctrl + U</code>	Cut the line before the cursor to the clipboard.
	<code>Ctrl + D</code>	Cut the word after the cursor to the clipboard.
	<code>Ctrl + W</code>	Cut the word before the cursor to the clipboard.
	<code>Ctrl + Y</code>	Paste the last item to be cut.
Processes	<code>Ctrl + L</code>	Clear the entire screen (like the clear command).
	<code>Ctrl + Z</code>	Place the currently running process into a suspended background process.
	<code>Ctrl + C</code>	Kill the currently running process by sending the SIGINT signal.
	<code>Ctrl + D</code>	Exit the current shell.
	<code>Return</code>	Exit a stalled SSH session.
History	<code>Ctrl + R</code>	Bring up the history search..
	<code>Ctrl + G</code>	Exit the history search.
	<code>Ctrl + P</code>	See the previous command in the history.
	<code>Ctrl + N</code>	See the next command in the history.

6.4 When you need help with Commands

If you ever see an experienced Linux user typing away at the command line in blazing speeds it can seem like memorising the ins and outs of commands and options is the only way to be productive and understand what's going on. But everybody starts somewhere, and even experienced command-line users don't memorize everything.

In the world of programming, it's not practical to try to memorise all of the syntax and options of command-line tools. Of course, it's important to remember the basics, but while you're getting started, you only need to remember a few commands. The first one is `man`, which stands for the **manual pages**.

```

1  MAN(1)                                Manual pager utils          MAN(1)      text
2
3  NAME
4      man - an interface to the system reference manuals
5
6  SYNOPSIS
7      man [man options] [[section] page ...] ...
8      man -k [apropos options] regexp ...
9      man -K [man options] [section] term ...
10     man -f [whatis options] page ...

```

A `man` page⁵ is a form of software documentation found on UNIX and Unix-like operating systems. Topics covered include programs, system libraries, system calls, and sometimes local system details.

⁵Stands for short for manual page.

Think of the man pages as a technical reference book for your Linux distribution⁶. If you know the name of a command, you can find out a wealth of information about what it does, what options it provides or what arguments it takes. To look up something in the manual pages, type `man`, followed by a command you want to learn. Open up the Terminal by `Ctrl + Alt + T`. Earlier, you saw the command `ls`, so let's look that up. Type `man ls` and press `\return`.

⁶i.e., Ubuntu, Mint

Some distributions or application specific installation of Linux remove the `man` pages to save up on space. In these system one must first do `unminimize` to install `man` pages.

```

1  man ls | head -10                                     C.R. 3
2
3  LS(1)                                User Commands          LS(1)      text
4
5  NAME
6      ls - list directory contents
7
8  SYNOPSIS
9      ls [OPTION]... [FILE]...

```

```

9  DESCRIPTION
10     List information about the FILEs (the current directory by default).

```

⁷Please ignore the `head - 10`, we will have a look at it later.

Here, you can see some information about the `ls` command⁷. You can see that it's for listing directory contents and in the synopsis section you get a quick overview of how to use the command. In this case it is `ls [OPTION]... [FILE]...`. We write `ls` followed by any of the options we need, and the file or folder path we want to use.

⁸for example `[OPTION]` and `[FILE]`

The terms in square brackets⁸ are optional. This basically means **you don't have to use these** for the command to work. You can just use the `ls` command by itself to see the default output of listing the directory. Here, below the description header, there is a bit more detailed information about the command, including its default behaviour and usage notes, and below, is a listing of the options that the command takes.

```

1  Usage: ls [OPTION]... [FILE]...
2  List information about the FILEs (the current directory by default).
3  Sort entries alphabetically if none of -cftuvSUX nor --sort is specified.
4
5  Mandatory arguments to long options are mandatory for short options too.
6      -a, --all            do not ignore entries starting with .
7      -A, --almost-all      do not list implied . and ..
8      --author             with -l, print the author of each file
9      -b, --escape          print C-style escapes for nongraphic characters
10     --block-size=SIZE     with -l, scale sizes by SIZE when printing them;

```

There are a lot of ways to use the `man` pages efficiently and is a powerful tool when you need to find what can a command do. There are other ways to learn about a command. Most of commands also have an option called `help`, which provides a brief amount of information about them. However, they usually refer you to the manual pages for more detailed documentation. Therefore, `help` will give you a brief information compared to the `man` command.

You can see if a command you're using has this feature available by typing `--help` after the command.

```

1  ls --help | head -10

```

C.R. 4
bash

```

1  Usage: ls [OPTION]... [FILE]...
2  List information about the FILEs (the current directory by default).
3  Sort entries alphabetically if none of -cftuvSUX nor --sort is specified.
4
5  Mandatory arguments to long options are mandatory for short options too.
6      -a, --all            do not ignore entries starting with .
7      -A, --almost-all      do not list implied . and ..
8      --author             with -l, print the author of each file
9      -b, --escape          print C-style escapes for nongraphic characters
10     --block-size=SIZE     with -l, scale sizes by SIZE when printing them;

```

Here you can scroll up and down to have a look at some of the information. There is another command that's useful when you're working in Bash, and that's just `help` by itself.

`help` Displays information about shell built-in commands.

```
1  help | head -10                                C.R. 5
                                         bash
```

```
1  GNU bash, version 5.2.21(1)-release (aarch64-unknown-linux-gnu)          text
2  These shell commands are defined internally. Type `help' to see this list.
3  Type `help name' to find out more about the function `name'.
4  Use `info bash' to find out more about the shell in general.
5  Use `man -k' or `info' to find out more about commands not in this list.
6
7  A star (*) next to a name means that the command is disabled.
8
9  job_spec [&]                      history [-c] [-d offset] [n] or hist>
10 (( expression ))                  if COMMANDS; then COMMANDS; [ elif C>
```

As we get into working with the Bash shell, the `help` tool can act as a handy reminder for the syntax of some Bash specific commands.

But what if you don't know the name of a command you are looking for?

In that case, you can use another program called `apropos` which searches a list of commands and their descriptions for text you provide as an argument.

`apropos` helps users find any command using its `man` pages.

So if you wanted to find out what can list things, I could type `apropos list` and see a number of results that match that word.

```
1  apropos list | head -10                                C.R. 6
                                         bash
```

```
1  port-contents(1)      - List the files installed by a given port          text
2  port-dependents(1), port-rdependents(1) - List ports that depend on a given (installed) port
3  port-deps(1), port-rdeps(1) - Display a dependency listing for the given port(s)
4  port-distfiles(1)      - Print a list of distribution files for a port
5  port-echo(1)           - Print the list of ports the argument expands to
6  port-installed(1)     - List installed versions of a given port, or all installed ports
7  port-list(1)           - List available ports
8  port-outdated(1)      - List outdated ports
9  port-variants(1)       - Print a list of variants with descriptions provided by a port
```

```
10 AllPlanes(3), BlackPixel(3), WhitePixel(3), ConnectionNumber(3), DefaultColormap(3),
    ↳ DefaultDepth(3), XListDepths(3), DefaultGC(3), DefaultRootWindow(3),
    ↳ DefaultScreenOfDisplay(3), DefaultScreen(3), DefaultVisual(3), DisplayCells(3),
    ↳ DisplayPlanes(3), DisplayString(3), XMaxRequestSize(3), XExtendedMaxRequestSize(3),
    ↳ LastKnownRequestProcessed(3), NextRequest(3), ProtocolVersion(3), ProtocolRevision(3),
    ↳ QLength(3), RootWindow(3), ScreenCount(3), ScreenOfDisplay(3), ServerVendor(3),
    ↳ VendorRelease(3) - Display macros and functions
```

Here's the command that can list directory contents we were looking for.

```
1 ls(1)                                - list directory contents          text
```

Searching for commands this way can be time-consuming, but if you know what you need to do but not the command to do it, `apropos` is very helpful and powerful.

6.5 Additional Information

6.5.1 Use Tab completion on the Shell

If you do not know the exact name of a command, then you can make use of tab completion. To use this action, launch the terminal by pressing **Ctrl + Alt + T** or just click on the terminal icon in the task bar. Just type the command name that you know in the terminal and then press **\tab** twice. For example, if we can't remember **man**, we can write **ma** and can choose one of the option the Bash shell presents us.

```
1 ~$ ls                                C.R. 7
                                         bash

1 macptopbm      make-ssl-cert          text
2 mag            mako-render
3 mailmail3       man
4 make           mandb
5 make4ht         manpath
6 makeconv        man-recode
7 makedtx        mapfile
8 make-first-existing-target  mapscrn
9 makeglossaries   match_parens
10 makeglossaries-lite  mathspic
11 makeindex      mattrib
12 makejvf        mawk
```

6.5.2 The info command

Some commands do not have their manuals written or they are either **incomplete**. To get help with those commands, we use **info**. To use this command, launch the terminal by pressing **Ctrl + Alt + T** or just click on the terminal icon in the task bar. Just type **info** in the terminal and with a space, type the name of the command whose manual does not exist and press **\return**.

```
1 info ls | head -10                         C.R. 8
                                         bash

1 File: coreutils.info,  Node: ls invocation,  Next: dir invocation,  Up: Directory listing text
2
3 10.1 ls: List directory contents
4 =====
5
6 The ls program lists information about files (of any type, including
7 directories). Options and file arguments can be intermixed arbitrarily,
8 as usual. Later options override earlier options that are incompatible.
9
```

¹⁰

For non-option command-line arguments that are directories, by

⁹A mostly a plain text transliteration of the Texinfo source, with the addition of a few control characters to separate nodes and provide navigational information, designed by the NU project.

¹⁰1 for commands, 2 for system calls, etc...

¹¹A major component of a document processing system developed by Bell Labs for the Unix operating system. It is mostly outdated.

The `info` command reads documentation in the `info` format⁹. It will give detailed information for a command when compared with the man page. The pages are made using the Texinfo tools which can link with other pages, create menus, and easy navigation.

Man v. Info

Man pages are the UNIX traditional way of distributing documentation about programs. The term “man page” itself is short for “manual page”, as they correspond to the pages of the printed manual; the man pages “sections”¹⁰ correspond to sections in the full UNIX manual. Support is still there if you want to print a man page to paper, although this is rarely done these days, and the sheer number of man pages make it just impossible to bind them all into a single book.

In the early '90s, the GNU project decided that “man” documentation system was outdated, and wrote the `info` command to replace it: `info` has basic hyperlinking features and a simpler markup language to use (compared to the `troff`¹¹ system used for man pages). In addition, GNU advocates against the use of man pages at all and contends that complex software systems should have complete and comprehensive documentation rather than just a set of short man pages.

There are actually other documentation systems in use, besides man and `info`: GNOME and KDE have their own, HTML-based system, etc.

In the end, the form in which you get documentation depends on the internal policies of the project that provided the software in the first place – there is no globally accepted standard.

6.5.3 The `whatis` command

This command is used with another command just to show a one liner usage of the latter command from its manual. It's a quick way of knowing the usage of a command without going through the whole manual.

`whatis` command in Linux is used to get a one-line manual page description. In Linux, each manual page has some sort of description within it. So, this command search for the manual pages names and show the manual page description of the specified filename or argument.

To use this command, launch the terminal by pressing `Ctrl + Alt + T` or just click on the terminal icon in the task bar. Just type `whatis` in the terminal and after a space, type the name of the command whose one liner description you want (for example `ls`) and then press `\return`.

1 whatis ls | head -10 C.R. 9
bash

```
1 dcmcjpls(1)          - Encode DICOM file to JPEG-LS transfer syntax      text
2 dcmdjpls(1)           - Decode JPEG-LS compressed DICOM file
3 gdircolors(1), dircolors(1) - color setup for ls
4 gls(1), ls(1)         - list directory contents
5 gdircolors(1), dircolors(1) - color setup for ls
6 git-ls-files(1)        - Show information about files in the index and the working tree
7 git-ls-remote(1)        - List references in a remote repository
8 git-ls-tree(1)          - List the contents of a tree object
9 git-mktree(1)           - Build a tree-object from ls-tree formatted text
10 gls(1), ls(1)          - list directory contents
```


Part IV

Robot Operating System

Chapter 7

Installation

Table of Contents

7.1	Introduction	155
7.2	Installing ROS Humble Hawksbill	156

7.1 Introduction

Setting up the Locale

Make sure you have a locale which supports UTF-8¹. If you are in a minimal environment (such as a docker container), the locale may be something minimal like POSIX. We test with the following settings. However, it should be fine if you're using a different UTF-8 supported locale.

```
1  locale # check for UTF-8
2
3  sudo apt update && sudo apt install locales
4  sudo locale-gen en_US en_US.UTF-8
5  sudo update-locale LC_ALL=en_US.UTF-8 LANG=en_US.UTF-8
6  export LANG=en_US.UTF-8
7
8  locale # verify settings
```

C.R. 1
bash

¹UTF-8 is a character encoding standard used for electronic communication. Defined by the Unicode Standard, the name is derived from Unicode Transformation Format - 8-bit. Almost every webpage is stored in UTF-8.

7.2 Installing ROS Humble Hawksbill

Before we can begin working with ROS, we must install all the necessary files and dependencies required. For these lectures we will install ROS 2 Humble Hawksbill which is currently available for Ubuntu Jammy (22.04 LTS).

While currently there are more up-to-date version of ROS 2 available, due to its long term support and established compatibility, we will be using ROS 2 Humble.

It is **heavily** recommended to be on Ubuntu 22.04 LTS as ROS 2 Humble is only officially supported on this version. It is possible to compile ROS 2 on other Ubuntu or Linux distributions, however, you need to compile the binaries yourself.

7.2.1 Set locale

Make sure you have a locale which supports [UTF-8](#). If you are in a minimal environment (i.e., docker container), the locale may be something minimal like [POSIX](#). We test with the following settings. However, it should be fine if you're using a different [UTF-8](#) supported locale. To start with installation, open up your terminal (+ + T).

```
1 sudo apt install software-properties-common  
2 sudo add-apt-repository universe
```

C.R. 2

bash

UTF-8 A variable-length character encoding standard used for electronic communication.

POSIX A family of standards specified for maintaining compatibility between operating systems.

Docker A set of platform as a service products that use OS-level virtualization to deliver software in packages.

7.2.2 Setup Sources

You will need to add the ROS 2 [apt²](#) repository to your system.

First ensure that the Ubuntu Universe repository is enabled.

```
1 sudo apt install software-properties-common  
2 sudo add-apt-repository universe
```

C.R. 3

bash

Now add the ROS 2 GPG key with apt.

```

1 sudo apt update && sudo apt install curl -y
2 sudo curl -sSL \
3     https://raw.githubusercontent.com/ros/rosdistro/master/ros.key \
4     -o /usr/share/keyrings/ros-archive-keyring.gpg

```

C.R. 4

bash

Then add the repository to your sources list.

```

1 echo "deb [arch=$(dpkg --print-architecture) \
2     signed-by=/usr/share/keyrings/ros-archive-keyring.gpg] \
3     http://packages.ros.org/ros2/ubuntu \
4     $(. /etc/os-release && echo $UBUNTU_CODENAME) main" | \
5     sudo tee /etc/apt/sources.list.d/ros2.list > /dev/null

```

C.R. 5

bash

echo A command that outputs the strings that are passed to it as arguments.

7.2.3 Install ROS 2 packages

Update your apt repository caches after setting up the repositories. ROS 2 packages are built on frequently updated Ubuntu systems. It is always recommended that you ensure your system is up to date before installing new packages.

```
1 sudo apt update && sudo apt upgrade
```

C.R. 6

bash

Desktop Install (Recommended): ROS, RViz, demos, tutorials. Development tools: Compilers and other tools to build ROS packages

```
1 sudo apt install ros-foxy-desktop python3-argcomplete
```

C.R. 7

bash

If you have read the document before doing copy and pasting, you can use the rosinstall.sh provided to you to automatically install everything required for this course.

```
1 sudo apt install ros-foxy-desktop python3-argcomplete
```

C.R. 8

bash

7.2.4 Setting up the Environment

Set up your environment by sourcing the following file in your terminal (+ + T) .

```

1 # Replace ".bash" with your shell if you're not using bash
2 # Possible values are: setup.bash, setup.sh, setup.zsh
3 source /opt/ros/humble/setup.bash

```

C.R. 9

bash

Alternatively, if you prefer to automate this action, simply type the following code in your terminal:

The aforementioned command simply writes the command to a document called bashrc which is a script running when the OS boots up.

The .bashrc file is a script file that's executed when a user logs in. The file itself contains a series of configurations for the terminal session. This includes setting up or enabling: coloring, completion, shell history, command aliases, and more.

It is a hidden file and simple `ls` command won't show the file.

```
1 #!/bin/bash
2
3 # First check your Ubuntu Version
4 # For maximum compatibility with ROS it needs to be 22.04 LTS
5
6 # Creating log for troubleshooting
7 echo "##### BEGIN ATTEMPT #####">>install_log.txt
8
9 echo "Welcome to ROS 2 Automated Installation"
10 echo ""
11 echo ""
12
13 # Accessing the Ubuntu version using AWK and piping it to grep for Regex
14 version=$(
15     awk "/VERSION_ID/" IGNORECASE=1 /etc/*release |
16         grep -Eo "[[:digit:]]+([.][[:digit:]])?\""
17 )
18
19 # Checks version for ROS Compliance
20 if [[ "${version}" == *"22.04"* ]]; then
21     echo "Version is supported."
22     sleep 1
23     echo "Continuing installation..."
24     sleep 1
25 else
26     echo "Your version: ${version}, What is needed: 22.04"
27     echo "I am sorry but your version is not supported."
28     echo "This install script will terminate"
29     exit
30
31 fi
32
33 echo ""
34 echo "Installing UTF-8 Compliance ..."
35
36 {
37     locale # check for UTF-8
38
39     sudo apt update
40     sudo apt install locales
```

C.R. 10

bash

```

C.R. 11
bash

41     sudo locale-gen en_US en_US.UTF-8
42     sudo update-locale LC_ALL=en_US.UTF-8 LANG=en_US.UTF-8
43     export LANG=en_US.UTF-8
44
45     locale # verify settings
46 } &>install_log.txt
47
48 echo ""
49 echo "Enabling Ubuntu Universe Repositories..."
50
51 {
52     sudo apt install software-properties-common
53     echo | sudo add-apt-repository universe
54 } &>install_log.txt
55
56 echo ""
57 echo "Adding ROS 2 GPG Keys ..."
58
59 {
60     sudo apt update
61     sudo apt install curl -y
62     sudo curl -sSL \
63         https://raw.githubusercontent.com/ros/rosdistro/master/ros.key \
64         -o /usr/share/keyrings/ros-archive-keyring.gpg
65 } &>install_log.txt
66
67 echo ""
68 echo "Adding ROS 2 to repository ..."
69
70 {
71     echo "deb [arch=$(dpkg --print-architecture) \
72 signed-by=/usr/share/keyrings/ros-archive-keyring.gpg] \
73 http://packages.ros.org/ros2/ubuntu \
74 $(. /etc/os-release && echo $UBUNTU_CODENAME) main" \
75     | sudo tee /etc/apt/sources.list.d/ros2.list >/dev/null
76 } &>install_log.txt
77
78 echo ""
79 echo "Getting Updates ..."
80
81 {
82     sudo apt update
83     echo yes | sudo apt upgrade
84 } &>install_log.txt
85
86 echo ""
87 echo "Installing ROS ..."
88
89 {
90     echo yes | sudo apt install ros-humble-desktop
91     yes | sudo apt install ros-dev-tools
92 } &>install_log.txt
93

```

```
94 {  
95     sudo apt install dbus-x11  
96 } &>install_log.txt  
97  
98 echo ""  
99 echo "Sourcing ROS file ..."  
100 sleep 1  
101  
102 echo "source /opt/ros/humble/setup.bash" >~/.bashrc  
103  
104 echo ""  
105 echo "Removing unnecessary files ..."  
106 sleep 1  
107 {  
108     yes | sudo apt autoremove  
109 } &>install_log.txt
```

Part V

Appendix

Appendix A

Tables

Table of Contents

A.1	Introduction	163
A.2	Student-t Distribution	164
A.3	Chi-Square Distribution	166

A.1 Introduction

The following are tables used in solving the exercises present in this book.

A.2 Student-t Distribution

Table A.1: Values of z for given values of the distribution function $F(z)$ with $m = 1 - 10$.

$F(z)$	Degrees of Freedom									
	1	2	3	4	5	6	7	8	9	10
0.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.6	0.32	0.29	0.28	0.27	0.27	0.26	0.26	0.26	0.26	0.26
0.7	0.73	0.62	0.58	0.57	0.56	0.55	0.55	0.55	0.54	0.54
0.8	1.38	1.06	0.98	0.94	0.92	0.91	0.9	0.89	0.88	0.88
0.9	3.08	1.89	1.64	1.53	1.48	1.44	1.41	1.4	1.38	1.37
0.95	6.31	2.92	2.35	2.13	2.02	1.94	1.89	1.86	1.83	1.81
0.975	12.71	4.3	3.18	2.78	2.57	2.45	2.36	2.31	2.26	2.23
0.99	31.82	6.96	4.54	3.75	3.36	3.14	3.0	2.9	2.82	2.76
0.995	63.66	9.92	5.84	4.6	4.03	3.71	3.5	3.36	3.25	3.17
0.995	63.66	9.92	5.84	4.6	4.03	3.71	3.5	3.36	3.25	3.17
0.999	318.31	22.33	10.21	7.17	5.89	5.21	4.79	4.5	4.3	4.14

Table A.2: Values of z for given values of the distribution function $F(z)$ with $m = 11 - 20$.

$F(z)$	Degrees of Freedom									
	11	12	13	14	15	16	17	18	19	20
0.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.6	0.26	0.26	0.26	0.26	0.26	0.26	0.26	0.26	0.26	0.26
0.7	0.54	0.54	0.54	0.54	0.54	0.54	0.53	0.53	0.53	0.53
0.8	0.88	0.87	0.87	0.87	0.87	0.86	0.86	0.86	0.86	0.86
0.9	1.36	1.36	1.35	1.35	1.34	1.34	1.33	1.33	1.33	1.33
0.95	1.8	1.78	1.77	1.76	1.75	1.75	1.74	1.73	1.73	1.72
0.975	2.2	2.18	2.16	2.14	2.13	2.12	2.11	2.1	2.09	2.09
0.99	2.72	2.68	2.65	2.62	2.6	2.58	2.57	2.55	2.54	2.53
0.995	3.11	3.05	3.01	2.98	2.95	2.92	2.9	2.88	2.86	2.85
0.995	3.11	3.05	3.01	2.98	2.95	2.92	2.9	2.88	2.86	2.85
0.999	4.02	3.93	3.85	3.79	3.73	3.69	3.65	3.61	3.58	3.55

Table A.3: Values of z for given values of the distribution function $F(z)$ with $m = 21 - 30$.

$F(z)$	Degrees of Freedom (m)									
	21	22	23	24	25	26	27	28	29	30
0.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.6	0.26	0.26	0.26	0.26	0.26	0.26	0.26	0.26	0.26	0.26
0.7	0.53	0.53	0.53	0.53	0.53	0.53	0.53	0.53	0.53	0.53
0.8	0.86	0.86	0.86	0.86	0.86	0.86	0.86	0.85	0.85	0.85
0.9	1.32	1.32	1.32	1.32	1.32	1.31	1.31	1.31	1.31	1.31
0.95	1.72	1.72	1.71	1.71	1.71	1.71	1.7	1.7	1.7	1.7
0.975	2.08	2.07	2.07	2.06	2.06	2.06	2.05	2.05	2.05	2.04
0.99	2.52	2.51	2.5	2.49	2.49	2.48	2.47	2.47	2.46	2.46
0.995	2.83	2.82	2.81	2.8	2.79	2.78	2.77	2.76	2.76	2.75
0.995	2.83	2.82	2.81	2.8	2.79	2.78	2.77	2.76	2.76	2.75
0.999	3.53	3.5	3.48	3.47	3.45	3.43	3.42	3.41	3.4	3.39

A.3 Chi-Square Distribution

Table A.4: Values of z for given values of the distribution function $F(z)$ with $m = 1 - 10$.

$F(z)$	Degrees of Freedom (m)									
	1	2	3	4	5	6	7	8	9	0
0.005	0.0	0.01	0.07	0.21	0.41	0.68	0.99	1.34	1.73	2.16
0.01	0.0	0.02	0.11	0.3	0.55	0.87	1.24	1.65	2.09	2.56
0.025	0.0	0.05	0.22	0.48	0.83	1.24	1.69	2.18	2.7	3.25
0.05	0.0	0.1	0.35	0.71	1.15	1.64	2.17	2.73	3.33	3.94
0.95	3.84	5.99	7.81	9.49	11.07	12.59	14.07	15.51	16.92	18.31
0.975	5.02	7.38	9.35	11.14	12.83	14.45	16.01	17.53	19.02	20.48
0.99	6.63	9.21	11.34	13.28	15.09	16.81	18.48	20.09	21.67	23.21
0.995	7.88	10.6	12.84	14.86	16.75	18.55	20.28	21.95	23.59	25.19

Table A.5: Values of z for given values of the distribution function $F(z)$ with $m = 11 - 20$.

$F(z)$	Degrees of Freedom (m)									
	11	12	13	14	15	16	17	18	19	20
0.005	2.6	3.07	3.57	4.07	4.6	5.14	5.7	6.26	6.84	7.43
0.01	3.05	3.57	4.11	4.66	5.23	5.81	6.41	7.01	7.63	8.26
0.025	3.82	4.4	5.01	5.63	6.26	6.91	7.56	8.23	8.91	9.59
0.05	4.57	5.23	5.89	6.57	7.26	7.96	8.67	9.39	10.12	10.85
0.95	19.68	21.03	22.36	23.68	25.0	26.3	27.59	28.87	30.14	31.41
0.975	21.92	23.34	24.74	26.12	27.49	28.85	30.19	31.53	32.85	34.17
0.99	24.72	26.22	27.69	29.14	30.58	32.0	33.41	34.81	36.19	37.57
0.995	26.76	28.3	29.82	31.32	32.8	34.27	35.72	37.16	38.58	40.0

Table A.6: Values of z for given values of the distribution function $F(z)$ with $m = 21 - 30$.

$F(z)$	Degrees of Freedom (m)									
	21	22	23	24	25	26	27	28	29	30
0.005	8.03	8.64	9.26	9.89	10.52	11.16	11.81	12.46	13.12	13.79
0.01	8.9	9.54	10.2	10.86	11.52	12.2	12.88	13.56	14.26	14.95

Continued on next page

Table A.6: Values of z for given values of the distribution function $F(z)$ with $m = 21-30$. (Continued)

$F(z)$	Degrees of Freedom (m)									
	21	22	23	24	25	26	27	28	29	30
0.025	10.28	10.98	11.69	12.4	13.12	13.84	14.57	15.31	16.05	16.79
0.05	11.59	12.34	13.09	13.85	14.61	15.38	16.15	16.93	17.71	18.49
0.95	32.67	33.92	35.17	36.42	37.65	38.89	40.11	41.34	42.56	43.77
0.975	35.48	36.78	38.08	39.36	40.65	41.92	43.19	44.46	45.72	46.98
0.99	38.93	40.29	41.64	42.98	44.31	45.64	46.96	48.28	49.59	50.89
0.995	41.4	42.8	44.18	45.56	46.93	48.29	49.64	50.99	52.34	53.67

Bibliography

- [1] Roland Siegwart, Illah Reza Nourbakhsh, and Davide Scaramuzza. *Introduction to autonomous mobile robots*. MIT press, 2011.
- [2] Michael LaBarbera. "Why the wheels won't go". In: *The American Naturalist* 121.3 (1983), pp. 395–408.
- [3] Julian FV Vincent et al. "Biomimetics: its practice and theory". In: *Journal of the Royal Society Interface* 3.9 (2006), pp. 471–482.
- [4] Fran ccois Druelle et al. "Convergence of bipedal locomotion: why walk or run on only two legs". In: *Convergent Evolution: Animal Form and Function*. Springer, 2023, pp. 431–476.
- [5] Damian M Lyons and Kiran Pamnany. "Rotational legged locomotion". In: *ICAR'05. Proceedings., 12th International Conference on Advanced Robotics, 2005*. IEEE. 2005, pp. 223–228.
- [6] G Schweitzer. "ROBOTRAC-a Mobile Manipulator Platform for Rough Terrain". In: *Proc. of Int. Symp. on Advanced Robot Technology*. 1991, pp. 411–416.
- [7] Mathias Thor et al. "A dung beetle-inspired robotic model and its distributed sensor-driven control for walking and ball rolling". In: *Artificial Life and Robotics* 23 (2018), pp. 435–443.
- [8] Z. P. Square R. Jones. *All Praise The Humble Dung Beetle*. 2018. URL: <https://www.smithsonianmag.com/science-nature/the-humble-dung-beetle-180967781/>.
- [9] Sharp Photography. *Common ostrich (Struthio camelus australis) male running (composite image), Damaraland, Namibia*. 2018. URL: [https://commons.wikimedia.org/wiki/File:Common_ostrich_\(Struthio_camelus_australis\)_male_running_composite.jpg](https://commons.wikimedia.org/wiki/File:Common_ostrich_(Struthio_camelus_australis)_male_running_composite.jpg).
- [10] Porges. *Zebra in Wellington Zoo*. 2025. URL: https://commons.wikimedia.org/wiki/File:Zebra_sideview.jpg.
- [11] Encyclopaedia Britannica. *Ant*. 2025. URL: <https://cdn.britannica.com/42/223142-050-7033F421/Red-ant-on-a-green-branch.jpg>.
- [12] Zhanbing Song et al. "The Impact of Exercise Play on the Biomechanical Characteristics of Single-Leg Jumping in 5-to 6-Year-Old Preschool Children". In: *Sensors* 25.2 (2025), p. 422.
- [13] Sven Böttcher. "Principles of robot locomotion". In: *Proceedings of human robot interaction seminar*. 2006.
- [14] Andrzej Krzywinski, Anna Niedbalska, and L Twardowski. "Growth and development of hand reared fallow deer fawns." In: *Acta theriologica* 29.29 (1984), pp. 349–356.
- [15] John Brackenbury. "Caterpillar kinematics". In: *Nature* 390.6659 (1997), pp. 453–453.

- [16] José L Pons. *Wearable robots: biomechatronic exoskeletons*. John Wiley & Sons, 2008.
- [17] William P Zyhowski, Sasha N Zill, and Nicholas S Szczechinski. "Adaptive load feedback robustly signals force dynamics in robotic model of *Carausius morosus* stepping". In: *Frontiers in Neurorobotics* 17 (2023), p. 1125171.
- [18] David A Winter. *Biomechanics and motor control of human gait: normal, elderly and pathological*. 1991.
- [19] Francesco Lacquaniti, Yuri P Ivanenko, and Myrka Zago. "Patterned control of human locomotion". In: *The Journal of physiology* 590.10 (2012), pp. 2189–2199.
- [20] Hyunglae Lee and Neville Hogan. "Investigation of human ankle mechanical impedance during locomotion using a wearable ankle robot". In: *2013 IEEE International Conference on Robotics and Automation*. IEEE. 2013, pp. 2651–2656.
- [21] R McN Alexander. "Optimization and gaits in the locomotion of vertebrates". In: *Physiological reviews* 69.4 (1989), pp. 1199–1227.
- [22] MIT. *The Raibert Hopper*. 1984. URL: http://www.ai.mit.edu/projects/leglab/robots/3D_hopper/3D_hopper.html.
- [23] Citizendum. *Asimo*. 2005. URL: <https://citizendum.org/wiki/ASIMO>.
- [24] Seshashayee S Murthy and Marc H Raibert. "3D balance in legged locomotion: modeling and simulation for the one-legged case". In: *ACM SIGGRAPH Computer Graphics* 18.1 (1984), pp. 27–27.
- [25] Marc H Raibert, H Benjamin Brown Jr, and Michael Chepponis. "Experiments in balance with a 3D one-legged hopping machine". In: *The International Journal of Robotics Research* 3.2 (1984), pp. 75–92.
- [26] Ben Brown and Garth Zeglin. "The bow leg hopping robot". In: *Proceedings. 1998 IEEE International Conference on Robotics and Automation (Cat. No. 98CH36146)*. Vol. 1. IEEE. 1998, pp. 781–786.
- [27] Robert Ringrose. "Self-stabilizing running". In: *Proceedings of International Conference on Robotics and Automation*. Vol. 1. IEEE. 1997, pp. 487–493.
- [28] Flamingo. *Spring Flamingo*. 2000. URL: http://www.ai.mit.edu/projects/leglab/robots/Spring_Flamingo/Spring_Flamingo.html.
- [29] Ilon Bengt Erland. "Rad fuer ein laufstabiles, selbstfahrendes fahrzeug". In: *German Patent No. DE2354404A1* (1974).
- [30] Kevin Dowling et al. "NAVLAB An Autonomous Navigation Testbed". In: *Vision and Navigation: The Carnegie Mellon Navlab*. Springer, 1990, pp. 259–282.
- [31] ackerman. *Ackerman Steering Mechanism*. 2025. URL: <https://www.dubizzle.com/blog/cars/ackerman-steering-mechanism/>.
- [32] Farnell. *Rotary Encoder, Module, Optical, Incremental, 500 PPR, 0 Detents, Vertical, Without Push Switch*. 2025. URL: <https://at.farnell.com/en-AT/broadcom-limited/aedb-9140-a13/encoder-3channel-500cpr-8mm/dp/1161087>.

- [33] Flyrobo. *GY-26 Digital Electronic Compass Sensor Module*. 2025. URL: <https://www.flyrobo.in/gy-26-digital-electronic-compass-sensor-module>.
- [34] FindLight. *Optical Gyroscopes: Measuring Rotational Changes With Sagnac Effect*. 2025. URL: <https://www.findlight.net/blog/optical-gyroscopes-measuring-rotations/>.
- [35] PiHut. *HC-SR04 Ultrasonic Range Sensor on the Raspberry Pi*. 2025. URL: <https://the-phut.com/blogs/raspberry-pi-tutorials/hc-sr04-ultrasonic-range-sensor-on-the-raspberry-pi>.
- [36] Reinhold. *Structured light sources on display at the 2014 Machine Vision Show in Boston*. 2014. URL: https://commons.wikimedia.org/wiki/File:Structured_light_sources.agr.jpg.
- [37] Andrzej. *CCD image sensor SONY ICX493AQA 10,14 (Gross 10,75) M pixels APS-C 1.8" 28.328mm (23.4 x 15.6 mm) from module IS-026 from digital camera SONY DSLR-A200 or DSLR-A300 sensor side*. 2014. URL: https://commons.wikimedia.org/wiki/File:CCD_SONY_ICX493AQA_sensor_side.jpg.
- [38] TeledyneCCD. *How a Charge Coupled Device (CCD) Image Sensor Works*. 2020. URL: https://www.teledyneimaging.com/media/1300/2020-01-22_e2v_how-a-charge-coupled-device-works_web.pdf.
- [39] K. Hirakawa and T.W. Parks. "Chromatic adaptation and white-balance problem". In: *IEEE International Conference on Image Processing 2005*. Vol. 3. 2005, pp. III-984. DOI: 10.1109/ICIP.2005.1530559.
- [40] Fstoppers. *Is There a Difference Between Color Temperature and White Balance?* 2022. URL: <https://fstoppers.com/natural-light/there-difference-between-color-temperature-and-white-balance-596031>.
- [41] Adrian Ilie and Greg Welch. "Ensuring color consistency across multiple cameras". In: *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*. Vol. 2. IEEE. 2005, pp. 1268–1275.
- [42] Teledyne. *Saturation and Blooming*. 2025. URL: <https://www.photometrics.com/learn/imaging-topics/saturation-and-blooming>.
- [43] Zhen Zhang et al. "Analysis and simulation to excessive saturation effect of CCD". In: *2nd International Symposium on Laser Interaction with Matter (LIMIS 2012)*. Vol. 8796. SPIE. 2013, pp. 96–102.
- [44] Baumer. *Operating principle and features of CMOS sensors*. 2025. URL: <https://www.baumer.com/int/en/service-support/function-principle/operating-principle-and-features-of-cmos-sensors/a/EMVA1288>.
- [45] econ Systems. *CMOS camera module*. <https://www.directindustry.com/prod/e-con-systems/product-168594-2365044.html>. URL: <https://www.directindustry.com/prod/e-con-systems/product-168594-2365044.html>.
- [46] TS Holst and GC Lomheim. *CMOS/CCD Sensors*. JCD publishing, 2007.
- [47] Mdf. *A photon noise simulation*. 2010. URL: <https://commons.wikimedia.org/wiki/File:Photon-noise.jpg>.

- [48] Mark-j. *Open Camera Blog*. 2018. URL: <https://sourceforge.net/p/opencamera/blog/2018/09/focus-bracketing-with-open-camera/>.
- [49] David C Hoaglin, Frederick Mosteller, and John W Tukey. *Understanding robust and exploratory data analysis*. Vol. 76. John Wiley & Sons, 2000.
- [50] Oded Koren. "A study of the Linux kernel evolution". In: *ACM SIGOPS Operating Systems Review* 40.2 (2006), pp. 110–112.
- [51] Maurice J Bach. "The Design of the UNIX". In: *RTM. Operating system Prentice Hall* (1986), pp. 312–329.
- [52] Steven C Johnson and Dennis M Ritchie. "UNIX time-sharing system: Portability of C programs and the UNIX system". In: *The Bell System Technical Journal* 57.6 (1978), pp. 2021–2048.
- [53] A. Lundqvist. *Timeline of Debian Linux and related projects*. 2012. URL: <https://commons.wikimedia.org/wiki/File:DebianFamilyTree1210.svg>.

