

# DS 1000B

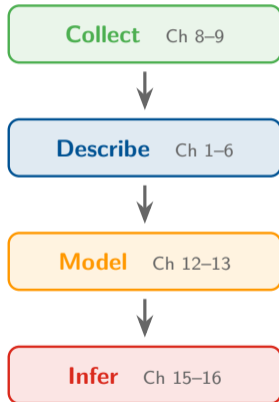
## Data Science Concepts

---

Last Lecture — Winter 2026

# DS 1000B in One Picture

---



- Gather data with sound methods
- Summarise and visualise.  
Explore the data prior to modelling it.
- Quantify uncertainty with probability
- Infer what the data suggests about the world beyond the data

## One Example, Four Stages

---

A researcher asks: how much sleep do university students get?

**Collect** Survey 200 students. SRS from the registrar's list.

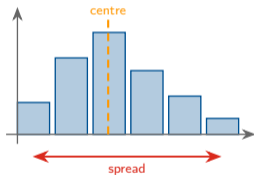
**Describe**  $\bar{x} = 6.8$  hours,  $\sigma = 1.4$  hours

**Model** Sampling distribution of  $\bar{x}$ : approximately Normal with  $SE = 1.4/\sqrt{200} \approx 0.10$ .

**Infer** 95% CI:  $6.8 \pm (1.96)(0.10) = 6.8 \pm 0.20$ , so between 6.6 and 7.0 hours.

## Chapter 1 Picturing Distributions with Graphs

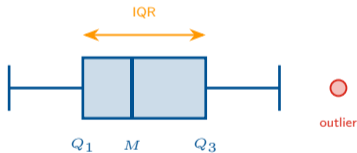
---



- individuals, variables, variable types
- histograms, bar charts, pie charts, stemplots
- describing distributions
- timeseries plots

## Chapter 2 Describing Distributions with Numbers

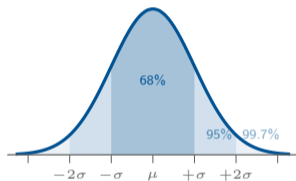
---



- mean, median, mode
- variance, standard deviation, IQR
- boxplots and modified boxplots
- summary measures

## Chapter 3 The Normal Distribution

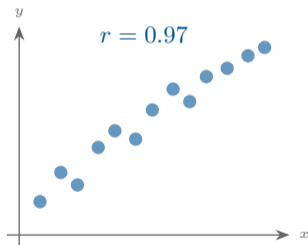
---



- density curves
- Normal distribution  $N(\mu, \sigma)$
- empirical rule
- calculate and interpret a z-score
- transforming to the standard Normal
- percentiles, probability calculations

## Chapter 4 Scatterplots and Correlation

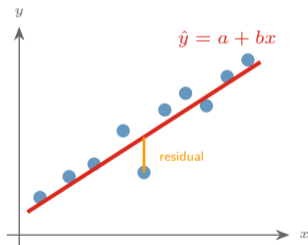
---



- scatterplots
- correlation coefficient

## Chapter 5 Regression

---



- slope and  $y$ -intercept interpretation
- predictions
- coefficient of determination  $R^2$
- residuals
- residual plots for assessing model fit
- influence and leverage

## Chapter 6 Contingency Tables

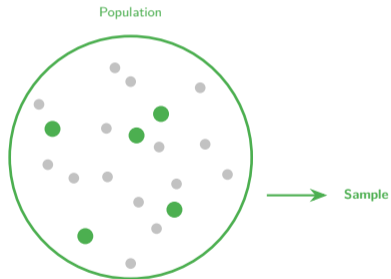
---

	Dog	Cat	Total
Male	1	4	5
Female	3	1	4
Total	4	5	9

- contingency tables
- marginal and conditional distributions
- independence of two categorical variables
- relative risk, odds ratio
- Simpson's paradox

## Chapter 8 Sampling

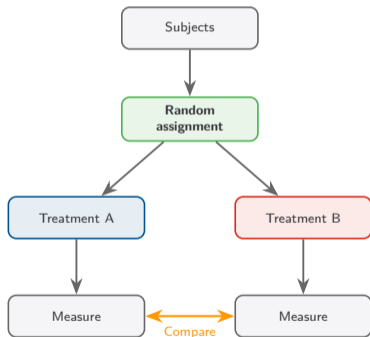
---



- population vs. sample; parameter vs. statistic
- problems with non-random sampling methods
- SRS, stratified, systematic, cluster, and multistage sampling
- practical issues with good designs

## Chapter 9 Designing Experiments

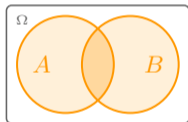
---



- observational studies vs. experiments
- explanatory and response variables
- factors, levels, treatments, and subjects
- principles of control, blocking, randomisation, and replication
- design types: CRD, block design, matched pairs
- blinding and placebo effects

## Chapter 12 Introduction to Probability

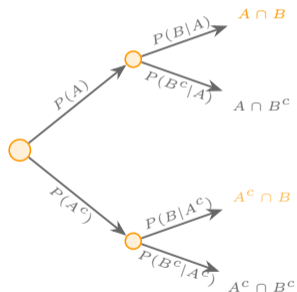
---



- outcomes, events, sample space, and probability
- algebra of events
- addition and complement rules
- Venn diagrams
- random variables

## Chapter 13 Conditional Probability

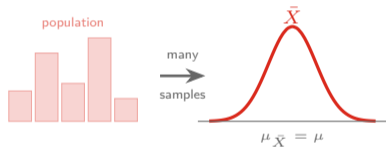
---



- conditional probability
- multiplication rule
- independent events
- tree diagrams
- Bayes' theorem
- Law of Total Probability

## Chapter 15 Sampling Distributions

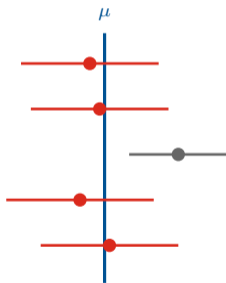
---



- sampling distribution of a statistic
- the behaviour of the sampling distribution of  $\bar{X}$

## Chapter 16 Confidence Intervals: The Basics

---



- CI interpretation
- CI construction
- relationship between CI and sample size, confidence level and population variability
- CI for a transformed parameter



## DS 1000B Final Exam Self-Assessment & Study Plan

Winter 2026 Final Exam: Sunday, April 19, 2026 10:00 AM - 1:00 PM

The goal is to help you gather an accurate assessment of your strengths and weaknesses across the course material, and to use that information to create a focused study plan for the final exam.

### Unit 1 Exploring Data Chapters 1 – 6

#### Chapter 1 Picturing Distributions with Graphs

I can ...	Not yet	Developing	Proficient
Distinguish categorical variables from quantitative variables	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Read and interpret a histogram	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Read and interpret a bar chart or pie chart for categorical data	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Interpret the features of a time series-plot: trends, seasonality and outliers	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

#### Chapter 2 Describing Distributions with Numbers

I can ...	Not yet	Developing	Proficient
Calculate and interpret the mean, median and quartiles of a dataset	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Calculate and interpret the IQR and standard deviation	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
State and apply the properties of standard deviation	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Identify outliers from a modified boxplot	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Read and construct a boxplot or a modified boxplot	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Choose appropriate summary measures for a distribution	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

## DS1000B – Notation and Terminology Glossary

Winter 2026

### Contents

Ch.1 – Picturing Distributions with Graphs	2
Ch.2 – Describing Distributions with Numbers	3
Ch.4 – Scatterplots and Correlation	5
Ch.5 – Regression	6
Ch.6 – Contingency Tables	7
Ch.8 – Sampling	8
Ch.9 – Designing Experiments	9
Ch.12 – Introduction to Probability	10
Ch.13 – Conditional Probability	12
Ch.3 – The Normal Distribution	13
Ch.15 – Sampling Distributions	14
Ch.16 – Confidence Intervals: The Basics	15
Common Points of Confusion	17
Alphabetical Index	19

#### Quick Reference: Common Notation

Symbol	Name	Meaning	Symbol	Name	Meaning
$\bar{x}$	x-bar	Sample mean	$P(A)$	Probability of event $A$	
$s$		Sample standard deviation	$P(A   B)$	Prob. of $A$ given $B$	
$\hat{p}$	p-hat	Sample proportion	$A \cap B$	Intersection ( $A$ and $B$ )	
$n$		Sample size	$A \cup B$	Union ( $A$ or $B$ )	
$\mu$	mu	Population mean	$\bar{X}, A^c, A'$	Complement (not $A$ )	
$\sigma$	sigma	Population standard deviation	$N(\mu, \sigma)$	Normal distribution	
$p$		Population proportion	$Z$	Standard Normal variable	
$Q_1, Q_3$		First, third quartile	$\Phi(z)$	Cumul. Normal prob. $P(Z \leq z)$	
$\hat{y}$	y-hat	Predicted value	$z^*$	Critical value	
$e_i$		Residual	$\bar{X}_n$	Sample mean (random var.)	
$a, b$		Slope, intercept	SE	Standard error	
$r$		Correlation coefficient	MOE	Margin of error	
$R^2$		Coefficient of determination			

### Ch.1 – Picturing Distributions with Graphs

#### Bar graph

Displays distribution of a categorical variable using rectangular bars.



#### Categorical variable

Places individuals into groups or categories.

#### Centre

A typical or representative value for the distribution.

#### Dataset

Structured collection of data on individuals and their variables.

#### Distribution

The values a variable can take and how often it takes them.



#### Exploratory Data Analysis (EDA)

Using graphs and summaries to describe variables and relationships.

#### Histogram

Visualisation for quantitative data:  $x$ -axis divided into bins, bar height shows count.



#### Individuals

The objects described by a set of data.

#### Observation

A single row in a dataset; all variable values for one individual.

#### Outlier

A value that lies far from the rest of the data.  $\circ$  Outlier, Ch.4



#### Pie chart

Circle divided into slices proportional to category frequencies.



#### Quantitative variable

Takes numerical values where arithmetic makes sense.

#### Shape

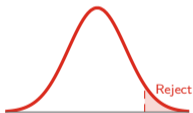
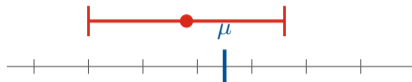
Overall pattern of a distribution: symmetric, left skewed or right skewed.

# Where It Leads

---

# Inference

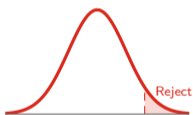
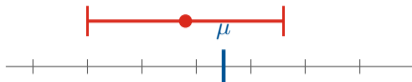
---



## Hypothesis testing

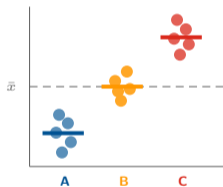
Is the observed effect real,  
or just noise?

# Inference



## Hypothesis testing

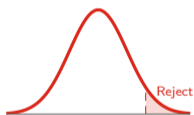
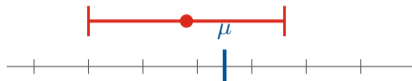
Is the observed effect real,  
or just noise?



## ANOVA

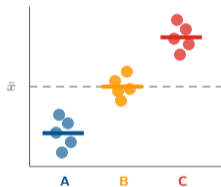
Is there a difference in centre between  
the groups?

# Inference



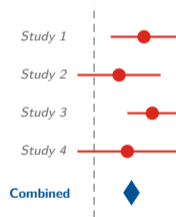
## Hypothesis testing

Is the observed effect real,  
or just noise?



## ANOVA

Is there a difference in centre between  
the groups?

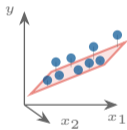
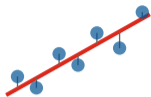


## Meta-analysis

Several studies have conflicting  
conclusions.  
How can we combine their evidence,  
and what does it say?

# Regression

---



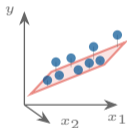
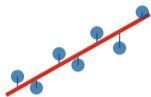
## Multiple regression

What if we have more than one independent variable?

SS 3859 Regression  
SS 4860 Advanced Regression

# Regression

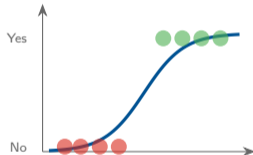
---



## Multiple regression

What if we have more than one independent variable?

SS 3859 Regression  
SS 4860 Advanced Regression



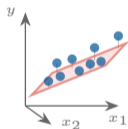
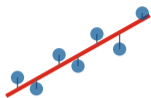
## Logistic regression

How can we predict other types of dependent variables?

SS 3860 Generalized Linear Models  
SS 3869 Applied Linear Models  
BIOSTATS 3110 Multivariable Methods

# Regression

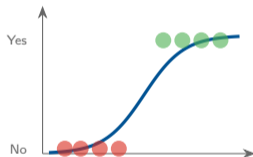
---



## Multiple regression

What if we have more than one independent variable?

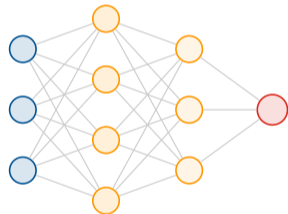
SS 3859 Regression  
SS 4860 Advanced Regression



## Logistic regression

How can we predict other types of dependent variables?

SS 3860 Generalized Linear Models  
SS 3869 Applied Linear Models  
BIOSTATS 3110 Multivariable Methods



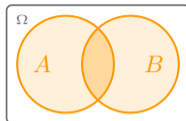
## Neural networks

How does AI translate languages or generate images?

DS 3000 Intro to ML  
CS 4452 Deep Learning  
AM 4264 Intro to Neural Networks

# Probability

---



$$\mathbb{E}[X]$$

Expectation

$$\text{Var}(X)$$

Variance

$$\mathbb{E}[\mathbb{E}[X | Y]]$$

Law of Total Expectation

$$M_Z(t) = e^{\frac{t^2}{2}}$$

Moment generating functions

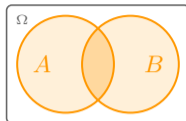
## Probability theory

Why is the Central Limit Theorem true?

How to make decisions under  
uncertainty? What are good betting  
strategies?

# Probability

---



$$\mathbb{E}[X]$$

Expectation

$$\text{Var}(X)$$

Variance

$$\mathbb{E}[\mathbb{E}[X | Y]]$$

Law of Total Expectation

$$M_Z(t) = e^{\frac{t^2}{2}}$$

Moment generating functions

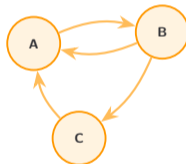
## Probability theory

Why is the Central Limit Theorem true?

How to make decisions under uncertainty? What are good betting strategies?

SS 2857 Probability I

SS 3657 Intermediate Probability

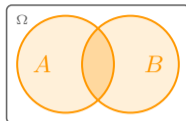


## Markov chains

How does Google rank billions of webpages?

SS 4654 Markov Chains

# Probability



$$\mathbb{E}[X]$$

Expectation

$$\text{Var}(X)$$

Variance

$$\mathbb{E}[\mathbb{E}[X | Y]]$$

Law of Total Expectation

$$M_Z(t) = e^{\frac{t^2}{2}}$$

Moment generating functions

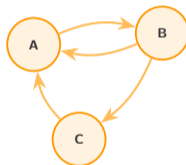
## Probability theory

Why is the Central Limit Theorem true?

How to make decisions under uncertainty? What are good betting strategies?

SS 2857 Probability I

SS 3657 Intermediate Probability



## Markov chains

How does Google rank billions of webpages?

SS 4654 Markov Chains

The cat sat on the ...



## Large Language Models

Conditional probability:

$$P(\text{next word} | \text{context})$$

AISE 4010 Deep Learning (Transformers)

CS 3346 Intro to AI

# Before You Go

---

---

## Show Up

- Attend your classes
- Ask questions in lecture, in office hours, and beyond

---

## Build Experience

- Pursue internships and research opportunities
- Develop skills beyond your coursework

---

## Adapt

- AI
- Struggle is part of learning

# Thank You

---

It has been a genuine pleasure teaching you this course.

Thank you for your curiosity, your effort, and your questions. They made the lectures a lot more enjoyable.

Good luck on the final and in everything that comes next. Don't hesitate to reach out.

## One Last Ask

---

Your feedback shapes how this course is taught in the future.

Please take a few minutes to complete the **SQCT course evaluations** they are anonymous and I go through them all carefully.

[go.blueja.io/Uu165VRnRUC1AZaZnmGHhA](https://go.blueja.io/Uu165VRnRUC1AZaZnmGHhA)



<https://go.blueja.io/Uu165VRnRUC1AZaZnmGHhA>