Data Analysis and Machine Learning 4 Week 1: Introduction, data modalities, variable types

Elliot J. Crowley, 16th January 2023







of EDINBURGH

What is data?

"information, especially facts or numbers, collected to be examined and considered and used to help decisionmaking, or information in an electronic form that can be stored and used by a computer"

Cambridge Dictionary

Data











Last Name	First Name	Age	Rank	Major	Gender	Current GPA	Photo
Adams	Grace	19	Sophomore	English	Female	3.78	DA
Bloomfield	Erika	21	Junior	Physics	Female	3.89	P
Chow	Kimmie	20	Senior	Political Science	Female	3.77	
Crutchfield	Seth	23	Senior	Psychology	Male	3.58	
Fitch	Fredrick	18	Freshman	Art	Male	4.0	
Grover	Oscar	26	Junior	Biology	Male	3.32	Alerty .



Wim MORRISON Supermarkets plc 803 7DL Noking Manager : Lee King Telephone : 01483 755552 Vat Number : 343475355

Savers Stamps Pick up a Card and Start Saving for Christmas Today

DATE: 19/06/2008 TIME: 17:54 TILL: 0019 NO: 01969232 You were served by: JENI

DESCRIPTION			£
M FRESH SEA M SIDE OF SA M KIPPER F M PORK LEG	A BREAD ALMON ILLETS STEAK	K	2.88 3.08 0.56 2.93
0.270kg 8 ET	ND HE	8	0.54
0.650kp 0 EC "M'BEST POTA" HORLICKS DUREX EXTRA S "M'TRIM BEANS "M'STRANBERS TETLEY TEA BA "M'VALUE ONIO "M'DOUBLE CRE "M'ENGLISH B "M'ENGLISH B "M'BUTTAR OF "M'BUTTAR OF MIVEA FOR MEN "M'LOOSE LEM	AFE SAFE SIES SOS SAFE SIES SOS SUTTER UTTER I Fer SI SI SONS		2.59 0.99 1.34 5.98 1.29 1.89 0.99 0.58 0.94 0.94 -0.08 1.98 1.98 0.28
Items Sold:	20	TOTAL	£35.54
		CASH	£40.00
	Chang	98	£4.45
VAT A 17.5% VAT 8 5.0% VAT D 0.0% VAT Total	-	£2.59): £5.98): £26.97):	£0.39 £0.28 £0.00 £0.67
ML	JLT:	ISAVE	

£0.08 SAVINGS AT MORRISONS

Thank you for shopping at Morrisons Please call again





Data Analysis

"The process of examining information, especially using a computer, in order to find something out, or to help with making decision"

Cambridge Dictionary

"The natural application of data; you use it to do cool stuff!" Elliot J. Crowley

Spotting patterns

Alcohol, total per capita (15+) consumption (in litres of pure alcohol) (SDG Indicator 3.5.2)



Year

Latest

Sex Both sexes



Source: https://www.who.int/data/gho/data/indicators/indicator-details/GHO/total-(recorded-unrecorded)-alcohol-per-capita-(15-)-consumption



Observing trends

Annual CO₂ emissions

not included.

LINEAR	LOG	Add country	Relative
35 billion	t		
30 billion	t		
25 billion	t		
20 billion	t		
15 billion	t		
10 billion	t		
5 billion	t		
01	t 		
	1750	1800	
Source: Glob	al Carbon	n Project	
1750	0		
CHAR	т	MAP	TA



Source: https://ourworldindata.org/co2-emissions



Telling (happy/sad?) stories





The Cartogram map shows the UK's 650 parliamentary seats as if they are hexagons of the same size. Hexagons by Esri

Source: https://www.bbc.co.uk/news/election/2019/results



Finding anomalies

Covid: Man offered vaccine after error lists him as 6.2cm tall

() 18 February 2021



Coronavirus pandemic



A man in his 30s with no underlying health conditions was offered a Covid vaccine after an NHS error mistakenly listed him as just 6.2cm in height.

Liam Thorp was wrongly classed as morbidly obese according to his height and weight



What is Machine Learning?

Machine Learning is robots and the colour blue



















Machine Learning is...

"the study of models that can learn from on new data."

Elliot J. Crowley

training data in order to make predictions

Machine Learning for a spring

- We want a model that given an arbitrary mass x can predict extension y
- We can attach some masses to the spring and record its extension
- These mass-spring measurements form our training data





Machine Learning for a spring

- Will will use a linear function y = mx + c as our model
- We can use the training data to find the *m*, *c* that give the best fit
- Given an arbitrary mass, we can input it to the function to predict extension



Is that it?





Face recognition









Detection and segmentation









Recommender systems

<

Books you may like



Mathematics for Machine Learning > Marc Peter Deisenroth Paperback \$46.99



Deep Learning (Adaptive Computation and Machine Learning series) > Ian Goodfellow 1,320 Hardcover \$39.00





Text to image



Text generation



This is a basic prompt for detecting sentiment.

Prompt

Decide whether a Tweet's sentiment is positive, neutral, or negative.

Tweet: "I loved the new Batman movie!" Sentiment:

Sample response

Positive



Corrects sentences into standard English.

Prompt

Correct this to standard English:

She no went to the market.

Sample response

She did not go to the market.

https://beta.openai.com/examples/



And more!



NEWS 22 July 2021

DeepMind's AI predicts structures for a vast trove of proteins

AlphaFold neural network produced a 'totally transformative' database of more than 350,000 structures from *Homo sapiens* and 20 model organisms.

Ewen Callaway



https://www.nature.com/articles/d41586-021-02025-4

The course

This is a brand new course!

- I hope you enjoy it
- There may be teething problems
- Feedback is very welcome

Course outline

- Data analysis (Weeks 1-3)
- Supervised learning and ethics (Week 4)
- Linear models (Weeks 5-7)
- Non-parametric and non-linear models (Weeks 8-10)

Course format

- In the **lecture** you are taught material
- In the **lab** session you will use this to solve problems using Python \bullet

these before the lab.

This is an applied course. Attending the labs is essential to success.

Each week's teaching consists of lecture (Monday AM) \rightarrow lab (Thursday PM)

There are **notes** that accompany each lecture that provide code. Go through

Assessment: Mini-tests (50%)

- There are three of these worth 16.66% each
 - Test 1 is at the end of the Week 3 lab and covers Weeks 1-3 material
 - Test 2 is at the end of the Week 7 lab and covers Weeks 5-7 material
 - Test 3 is at the end of the Week 10 lab and covers Weeks 8-10 material
- Each test consists of short-answer questions and some coding exercises
- They are open-book and must be taken live and in-person

Assessment: Coursework 1 (25%)

- This will be released in the Week 4 lecture on Monday 6th February
- You will create slides and record a short presentation using them
- This will be a case study on a real-world machine learning application
- You will critique this application from an ethical standpoint
- The deadline is Tuesday 21st February @ 1600 (Flexible learning week)

Assessment: Coursework 2 (25%)

- This will be released in the Week 8 lecture on Monday 13th March
- You will be given a dataset
- You will perform exploratory data analysis and apply machine learning to this dataset
- You will produce a short report on your findings supplemented with code
- The deadline is Tuesday 28th March @ 1600 (Week 10)

Data Modalities

Data exists in different modalities











· · · · · · · · · · · · · · · · · · ·		

Name	Name	Age	Rank	Major	Gender	Current GPA	Photo
Adams	Grace	19	Sophomore	English	Female	3.78	31
Bloomfield	Erika	21	Junior	Physics	Female	3.89	-
Chow	Kimmie	20	Senior	Political Science	Female	3.77	
Crutchfield	Seth	23	Senior	Psychology	Male	3.58	
Fitch	Fredrick	18	Freshman	Art	Male	4.0	
Grover	Oscar	26	Junior	Biology	Male	3.32	Attenue



Wm MORRISON Supermarkets plc 803 7DL Noking Manager : Lee King Telephone : 01463 755552

Vat Number : 343475355 Savers Stamps Pick up a Card and Start Saving for Christmas Today

DATE: 19/08/2008 TIME: 17:54 TILL: 0019 NO: 01969232 You were served by: JENI

DESCRIPTION			ž
M FRESH SE M SIDE OF S M KIPPER F M PORK LEG	A BREAK	M	2.88 3.08 0.56 2.93
0.270kg 8 E	1.99/k	8	0.54
0.650kp 0 EC "M'BEST POTA" HORLICKS DUREX EXTRA S "M'TRIM BEAKS "M'STRANBERS TETLEY TEA BA "M'VALUE ONIO "M'DOUBLE CRS "M'ENGLISH E "M'ENGLISH E "M'BUTTAR OT "M'RASPBERRIS NIVEA FOR MEN "M'LOOSE LEM	SAFE SAFE SEES NGS NGS SUTTER FUTTER FONS NONS	,	2.59 0.99 1.34 5.98 1.29 1.88 1.89 0.99 0.56 0.94 1.98 1.98 1.98 1.98 0.28 0.28
Items Sold:	20	TOTAL	£35.54
		CASH	£40.00
	Chang	98	£4.45
VAT A 17.5% VAT 8 5.0% VAT D 0.0% VAT Total		£2.59): £5.98): £26.97):	£0.39 £0.28 £0.00 £0.67
ML	JLT	ISAVE	

AT MORRISONS

Thank you for shopping at Morrisons Please call again





Time series data

- y axis is some quantity we care about
- x axis is time



Source: https://www.xe.com/currencycharts/?from=GBP&to=USD&view=10Y



Time series data

• For example, speech!



Source: https://towardsdatascience.com/beginners-guide-to-speech-analysis-4690ca7a7c05



Image data

- An image is a rectangular array of $H \times W$ pixels



Each pixel consists of three numbers: the amount of red, green, and blue

blue red green 255]

Image data

- This gives us a red, green, and blue 2D array
- These are stacked along the z axis to form a 3D array



Tabular data

- Looks like a table with rows and columns
- Rows are objects and columns are attributes of those objects
- An example is the iris dataset of 150 flowers

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal wid
0	5.1	3.5	1.4	
1	4.9	3.0	1.4	
2	4.7	3.2	1.3	
3	4.6	3.1	1.5	
4	5.0	3.6	1.4	
145	6.7	3.0	5.2	
146	6.3	2.5	5.0	
147	6.5	3.0	5.2	
148	6.2	3.4	5.4	
149	5.9	3.0	5.1	

n (cm)	species
0.2	setosa
2.3	virginica
1.9	virginica
2.0	virginica
2.3	virginica
1.8	virginica



Free-form data

- Largely unstructured and usually text
- Can (sometimes!) be hacked into e.g. tabular data

$\overrightarrow{}$ $\overrightarrow{}$

Stopped by on a Sunday afternoon, not so crowded and we got a table outside right away. Service was not attentive, we had to go in to get waitstaff including ordering and paying the bill. Food was meh. Ordered the prosciutto scramble, arugula and fennel salad, and Caesar salad. Don't think our scramble came with prosciutto, and arugula salad was extremely sour and quite plain. Fried cauliflower was quite tasty.

Overall a very mediocre place.

Dracula is a novel by Bram Stoker, published in 1897. As an epistolary novel, the narrative is related through letters, diary entries, and newspaper articles. It has no single protagonist, but opens with solicitor Jonathan Harker taking a business trip to stay at the castle of a Transylvanian noble, Count Dracula. Harker escapes the castle after discovering that Dracula is a vampire, and the Count moves to England and plagues the seaside town of Whitby. A small group, led by Abraham Van Helsing, hunt Dracula and, in the end, kill him.

Dracula was mostly written in the 1890s. Stoker produced over a hundred pages of notes for the novel, drawing extensively from Transylvanian folklore and history. Some scholars have suggested that the character of Dracula was inspired by historical figures like the Wallachian prince Vlad the Impaler or the countess Elizabeth Báthory, but there is widespread disagreement. Stoker's notes mention neither figure. He found the name Dracula in Whitby's public library while holidaying there, picking it because he thought it meant *devil* in Romanian.

Following its publication, Dracula was positively received by reviewers who pointed to its effective use of horror. In contrast, reviewers who wrote negatively of the novel regarded it as excessively frightening Comparisons to other works of Gothic fiction were common, including its structural similarity to Wilkie Collins' The Woman in White (1859). In the past century, Dracula has been situated as a piece of Gothic fiction. Modern scholars explore the novel within its historical context—the Victorian era—and discuss its depiction of gender roles, sexuality, and race.



Replying to @JakeBlueatSM

May be initiated not by the country leaders, but one of the AI's, if it decides that a prepemptive strike is most probable path to victory

11:36 PM - 3 Sep 2017





Follow

Nomenclature

- A dataset is a collection of data items (/ data points)
- A data item is a set of elements
- An **element** is a measurable or countable quantity









A dataset of **images**

An image is a set of pixels

A **pixel** measures the intensity of different colour(s)

Variable Types

Tabular data (again!)

- A table is a dataset and its rows are data items
- The measurements for a given attribute vary across the dataset

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal wid
0	5.1	3.5	1.4	
1	4.9	3.0	1.4	
2	4.7	3.2	1.3	
3	4.6	3.1	1.5	
4	5.0	3.6	1.4	
145	6.7	3.0	5.2	
146	6.3	2.5	5.0	
147	6.5	3.0	5.2	
148	6.2	3.4	5.4	
149	5.9	3.0	5.1	

• Each data item is a set of elements which are measurements of attributes

h (cm)	species
0.2	setosa
2.3	virginica
1.9	virginica
2.0	virginica
2.3	virginica
1.8	virginica



Variables

- The measurements for a given attribute vary across the dataset
- This means we can think of the attributes as variables
- There are different **types** of variables

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica

S

Categorical variables

Measurements of the variable correspond to descriptive categories

- For nominal variables the categories have no order
- For ordinal variables the categories are ordered

iris species (nominal)







0 setosa

versicolor

2 virginica

level of education (ordinal)



primary

secondary

2

3 university

Numerical variables

Measurements of the variable can be discrete or continuous

- For discrete variables they can only be whole numbers
- For continuous variables they can be any real number (within a given range)



The number of times this man tosses this coin is discrete The length of a tie is continuous

Continuous variables

These can be further divided into interval and ratio

- For interval variables there is no true zero measurement; It's relative
- For ratio variables zero has a clear meaning i.e. the absence of something

Temperature in Celcius is interval

Temperature in Kelvin is absolute



Summary

- We have considered different modalities of data e.g. tabular, image, freeform
- We have established the nomenclature for talking about data
- We have seen how attributes in tabular data can be treated as variables
- We have considered different variable types