Design decisions and dilemmas in a new data science course

David Sterratt

9 May 2022

Outline

- Context
- Intended learning outcomes
- Course description
- Inspiration and courses
- Decisions

Context

- "Inf2 Foundations of Data Science" (Inf2-FDS, or just FDS) is a new 2nd year undergraduate course, taken by Informatics students on all degree programmes
- 330 students in 2020/21; 290 in 2021/22
- Delivered by Kobi Gal and David Sterratt
- Arises out of the pre-honours curriculum redesign
- Two semester, 20 credit course
- Designed to fit with "Discrete Maths and Probability" (DMP, in S1)
- Runs alongside "Inf2 Introduction to algorithms and Data Structures"
- Replaces part of "Inf2b Learning"
- Somewhat less overlap with "Introductory Applied Machine Learning" (year 3 course) than Inf2b had

(Intended) learning outcomes

At the end of the course students will be able to:

- 1. Describe and apply good practices for storing, manipulating, summarising, and visualising data. (Data storage, manipulation and visualisation)
- 2. Use standard packages and tools for data analysis and describing this analysis, such as Python and LaTeX. (**Tool use**)
- 3. Apply basic techniques from descriptive and inferential statistics and machine learning; interpret and describe the output from such analyses. (Stats and ML)
- Critically evaluate data-driven methods and claims from case studies, in order to identify and discuss a) potential ethical issues and b) the extent to which stated conclusions are warranted given evidence provided. (Critical evaluation)
- 5. Complete a data science project and write a report describing the question, methods, and results. (**Data science project**)

Course description - technical topics

1. Data wrangling and exploratory data analysis

- Working with tabular data
- Descriptive statistics and visualisation
- Linear regression and correlation
- Clustering
- 2. Supervised machine learning
 - Classification
 - More on linear regression; logistic regression
 - Generalization and regularization
- 3. Statistical inference
 - Randomness, simulation and sampling
 - Confidence intervals, law of large numbers
 - Randomized studies, hypothesis testing

Course description - real-world implications

A. Implications:

- Where does data come from? (Sample bias, data licensing and privacy issues)
- Visualisation: misleading plots, accessible design
- Machine learning: algorithmic bias and discrimination
- B. Thinking, working, and writing:
 - Claims and evidence: what can we conclude; analysis of errors
 - Reproducibility; programming "notebooks" vs modular code
 - Scientific communication; structure of a lab report
 - Reading and critique of data science articles

Design decision 1: Course outline

- Sketched out by Sharon Goldwater, Heather Yorston & Kobi Gal
- Rough order:
 - 1. Data wrangling and exploratory data analysis
 - Including linear regression, PCA & K-means
 - 2. Supervised machine learning
 - Including k-NN and evaluation
 - 3. Statistical inference (S2)
 - Randomness, sampling, statistical simulations, confidence intervals, hypothesis testing
 - Logistic regression and linear regression using max likelihood
 - 4. Project

Inspiration & sources

- Berkeley Data 8 a course introducing programming and statistical inference to students with no programming background
 - Statistical inference via statistical simulations rather than standard probability distributions, as per a standard stats course
 - Python, though using their own datascience module
- Devore & Berk Modern mathematical statistics with applications
 - Frequentist approach
 - Essentially the same book as used for the DMP course
- Gelman & Nolan Teaching statistics a bag of tricks
- Shannon Vallor's pack An Introduction to Data Ethics
- Mine Çetinkaya-Rundel's Intro to Datascience course (https://www.introds.org/)

Design decision 2







EOUNDATIONS OF DATA SCIENCE

Design decision 3: Assessment

- Originally 10% class test, 30% project, 60% exam
- In July (with other courses) made decision to move to:

What?		ILOs assessed
20% Automarked	S1	Stats and ML
20% Visualisation & data	52 S1	Data storage, manipulation and visualisation, Tool use
20% "Critical evaluation"	S2	Critical evaluation
40% Final project	S2	Data science project, Tool use, Stats and ML, Data storage, manipulation and visualisation

No participation points or microcredit – a mistake?

Design decision 4: Learn to live with Learn

€ THE UNIVERS	Learn		Learn Help 🕜		David Sterratt	<u>1</u> ₹ (
5 				My Learn	Self-Enrol H	elp
Informatics 2 - Foundations of	Data Science (2020-2021)[YR]	Announcements		0	Edit Mode Is: 🚺	N ?
Informatics 2 - Foundations of Data Science (2020-2021)(YR) Read Me first Welcome Course Information Announcements	Announcements	the university of edinburgh informatics	FOUN OF DATA SCIE	DAT NCE		3
Course Materials	repositionable bar to pin the Students do not see the bar o	anecty below the repositionate but, kearder by a digging and m to the top of the list and prevent new announcements from s and cannot reorder announcements.	uperseding them. The order shown	here is the order p	resented to students.	
Lab Instructions Workshops	Create Announcement					
Discussions (Piazza) Library Resources Wiki	New announcemen	nts appear below this line				
Miniproject examples Wikis 🛿 Assessment	Week 5 materials Posted on: Sunday, 7 Fe	and announcements abruary 2021 15:26:04 oʻclock GMT		5	Posted by: David Sterratt Posted to: Informatics 2 - Foundations of Data Scie 2020-2021)[YR]	nce

Why? Familiarity for students & supported tool
 But... Collaborate caused issues for some students -> move to Teams (not popular with all students and staff!)

Design decision 5: Weekly pattern

Statistical inference	
Semester 2, Week 1	Topic: Randomness, sampling and simulation (00:53:23) Topic: Estimation: point estimates and confidence intervals (01:31:47) Coursework: Essay citilquing data science study (academic paper or news article) (20%) [Released Monday] Reading: Computational and Inferential Timiting, Chapter 10 Reading: Computational and Inferential Timiting, Chapter 10 Reading: Computational and Inferential Timiting, Chapter 10 Reading: Computational Inferential Timiting, Chapter 10 Reading: Computational Inferential Timiting, Chapter 10 Reading: Computational Inferential Timiting, Chapter 13 Reading: Computational Inferential Timiting, Chapter 13 Reading: Computational Inferential Timiting, Chapter 13 Reading: Reading: Computational Inferential Timiting, Chapter 13 Reading: Reading: Reademic 13 Sections 7.1, 8.1-8.3 and 8.5 Python Lab: K-nearest neighbours Workshop: Reflecting on the mini-project Question and answer session (Moved to next week)
Semester 2, Week 2	Topic: Hypothesis testing and p-values (00:49:30) Topic: Logistic regression (01:14:40) Reading: <u>XSCD comic strip on multiple testing - funnyd</u> Reading: <u>Appointesis na Jubilitythought-provoking and amusing article</u> Task: Problem sheet for S2 Week 3 Workshop Python Lak: Randomness and sampling Question and answer session (Monday)
Semester 2, Week 3 25-29 January	Topic: A/B testing (00:35:17) Coursework Deadline: Essay [Monday] Reading: Computational and Inferential Thinking. Chapter 12 Reading: Modern Mathematical Statistics with Applications, Chapter 10 Extra: Extra: Bayesian Inference applied to A/B testing Python Lab: Estimation with the bootstrap (including confidence intervals) Workshop: Problem sheet on S2 Week 1 material Question and answer session (Monday)

- Important to give students a rhythm
- Took to it after a few weeks

A topic



Pzi determine what Ppi outo these loadings are

lookings Swe **(**))

14:12 / 21:41

More details

- Video lectures, grouped in 1 or 2 topics
- Most topics have lecture notes attached
- "Comprehension questions" on topic actually a bit more than comprehension.
 - Purely optional, but popular, especially before class test
- Lab
- Every other week: "Workshop"
- QA session every week one used for a guest lecture

Design decision 6: Write lecture notes before lecture recording

🔏 Menu	t i		fds-lecture-notes			Ab Review	Share Stare	🚱 Submit	History	🗩 Chat
1 in 2	/8	Sou	rce Rich Text	\sim	C Re	ompile -	😬 🕹			2
1	sample-means	1848	\xreviewedsec{DCS}		4 5	emester	, Week 3			
1	sampling-distri	1849 1850	% \begin{itemize}	Θ	4	.1 Тори	:: Ethics ii	i data analys	sis	
-	sampling-distri	1051	% \item Use PCA on dataset % \item Write PCA from scratch \todo(Check with ADS lecturers they	G		4.1.1	Video:	Data protec	tion and pri	vacy (17:
-	sampling-distri	1053	<pre>% don't seem to be biting on this one!}</pre>			4.1.2	Video:	Bias (12:31)	
1	sampling-distri	1054	> (end(itemize)			4.1.3	Video:	Case study:	AI classific	er exhibiti
-	squirrel-regres	1056 1057	\qa		4	.2 Work	shop: Eth	ical discussi	on	
-	squirrel-regres	1058 -	\part{Supervised machine learning}	1	4	.3 Pythe	on Lab: Pa	indas – Data	wrangling	
	Squirrel.ipynb	1059	\chapter{Semester 1, Week 9}		4	.4 Ques	tion and a	nswer sessio	n (Friday)	
-	swain-versus	1061 1062	\label{main:cha:week-9}		5 8	emester 1	. Week 4			
~ 📾 S2	2-01-2-estimatio	1863	<pre>\input{topics/9-1-intro-supervised-learning/9-1-intro-supervised-learning}</pre>		5	.1 Topic	: Data co	llection, bias	es in data .	
200	air-bootstran	1065	Clustering exercise (write from scratch?) Use K-means on	- 14		Read	ing: Com	outational ar	ud Inferentia	al Thinkir
 File out 	line	1865	dataset - possibly inises or similar)	>		5.1.1	Video:	Webscrapin	g and infere	ence (16:
Sem	nester 1, Week 6	1868 1869	\topic{Intro to mini-project (covered in QA session)}	- 1	5	2 Topic	: Statistic	al relationsh	ins	
Sem	nester 1, Week 8	1070	<pre>\video{Intro to mini-project}{00:00}</pre>			5.2.1	Video:	Statistical re	lationships	(20:55)
Supervision	sed machine le	1071	\lab[https://github.com/Inf2-FDS/week09-kmeans]{SKS-means}		4	3 Pythe	n Lah: D	ata Represer	tation I – N	Astrolotlib
Sem	nester 1, Week 9	1073	\dataset{Breast cancer again}		-	A Tack	Dranarati	on for Week	5 Worksho	n on Vier
Sem	nester 1, Week	1075	\xreviewedsec(DCS)		-	5 Oues	tion and a	newar sassio	n (Friday)	p on vist
Sem	nester 1, Revisi	1076	\workshop{Distance-based clustering}			Ques	don and a	iiswei sessie	(i (i i iday)	
 Statistic 	al inference	1078	Discuss hing mosts on other data science readings		6 8	emester 1	, Week 5			
Sem	nester 2, Week 1	1080			6	.1 Cour	sework: D	ata manipul	ation and vi	isualisatic
Sem	nester 2, Week	1081 1082	/da		6	.2 Topic	: Statistic	al Prelimina	ries	
Sem	nester 2, Week	1083	\xdidsec(DCS)			Read	ing: Mode	rn Mathema	tical Statis	tics with .
Sem	nester 2, week	1085 -	\chapter{Semester 1, Week 10}			6.2.1	Video:	Sample and	population	mean (09
 Project 		1086	<pre>\label{main:sec:week-10}</pre>			6.2.2	Video:	Sample and	population	median (
Sem	nester 2. Weeks	1888	\topic(Evaluation)		_	51212		- mpre uno	r-r-million	

Design decision 7: Lecture recording



My new best friends



Setup details

- OBS (recording)
- Wacom Tablet
- Xournal++ (whiteboard)
- Kdenlive (editing)
- Media hopper create (hosting and captioning)

Design decision 8: Ethics

- Concept is that ethics is embedded in the course, not just in the final lecture
- Started with group presentation and discussion based around scenarios in Shannon Vallor's pack:
 - Facebook's emotion manipulation experiment (informed consent)
 - OK Cupid data breach (legal, terms and conditions)
- Also mentioned algorithmic bias and transparency and interactions between law and outcomes in credit approval (invited lecture by Dr Galina Andreeva from the Business School)
- Questions on ethics and law in class tests and coursework
- Students seem to have been interested

Design decision 9: Collaboration groups

• Workshop \approx Tutorial

- Split workshop groups of 16 into 4 "collaboration groups" of 4, and set tasks
- Worked well for first 2 workshops...
- Worked well for some students: "Love the group work, I've gotten to know my group well and they're lovely. Its really refreshing to chat to new people"

But not others:
 "Managing the collaboration group can be tricky"
 "Maybe add a symbolic grade/engagement, to produce a more engaged discussion."

Design decision 10: The final project, Weeks 7-11

Three public datasets:

- 20-21: Edinburgh Just Eat Bike use, Higher Education stats, Scottish Munro features)
- 21-22: Scottish A&E waiting times, Spotify, EEDI education dataset
- Seed question + ideas for future questions
- Project can be undertaken by individuals, pairs or threes (page limits of 6, 8 and 10 respectively)
- No lectures, but workshops at which lightweight presentations are made. (Not marked, but prerequisite)
- ► Team formation & code submission via GH Classroom
- Form for assessment of contributions
- Lots of python scripts to administer

Feedback on the project

The process of writing my/our presentation and hearing other presentations helped with the final project report



Should groups be mandatory?



Being online

- Videos seem to work well & comprehension questions seem to work well
- Q&As appreciated by those who attend (20-70)
 I think the comprehension quizzes are a great way to see if
 I actually absorbed the knowledge from the slide and the
 q&as are great.
- Workshops and labs are harder...

How it started...



... how it's going



2021/22: Going hybrid

Labs in person

- Workshops moved to in person, using teaching studios
 - Worked well, allowing ad-hoc collaboration groups to form
- Re-used recorded lectures (with some corrections)
 - Not universally popular...
 - ... will move to in-person
- Assessment structure similar to last year
 - Gradescope used where appropriate
 - Modified rubrics

Rules for visualisation

Informatics 2 - Foundations of Data Science: Visualisation principles and guidance

Principle 1: Show the data

Aim to show as much of the data as possible without leading to a confusing visualisation. There are often multiple ways of representing the same dataset, and no "right" answer. The following guidance on arranging the the graphical elements of the plot should help you to show as much of the data as . Present many numbers in a small space possible:

- · Choose an appropriate plot type. Some basic types are:
- Bar charts: good for plotting numeric variables associated with categorical or ordinal variables, for example the mean weight (numeric variable) of male and female (categorical variable) squirrels
- Line charts: good for showing trends of numerical variables over time (a numerical variable).
- Scatterplots: Scatterplots show the relationship between two numeric variables.
- Boxplots: Boxplots are a way of representing the distribution of a numeric variable for multiple categories. For example, the median and inter-quartile range of the weights of male and female squirrels.
- good for showing the distribution of a single variable.
- · Show multiple variables by using length, shape, size and · Title or caption colour:
- The above plots are all bivariate, since they show the relationship between two variables. Use shape and colour to create extra dimensions for categorical variables. For example, in a scatterplot of squirrel weight versus length. we can indicate sex using colour. We could also indicate our age categories by changing the size or the shape of the markers. However, adding information using marker properties can detract from the plot.
- Barcharts can be extended to two categorical variables and one numerical variables by using colour.
- Use colour effectively. (Wexler et al., 2017, pp. 14–18)
- Choose an appropriate colour scale, depending on if the data is sequential, diverging or categorical.
- Colour can also be used to highlight features in the plot. e.g. the largest two bars in a bar plot.

- · Encourage the eve to compare several pieces of data. The following guidelines should help to avoid distorting the e.g. by using multiple plots with the same scale.
- Wexler et al. (2017), p. 31, is a nice example of how this can work better than using multiple symbols on a plot (p. 30).

- A boxplot takes up as much space as a barplot, but conveys more information. For example, a boxplot of the squirrel's weight versus sex shows information about the distribution of the weight as well as the mean weight.
- · Choose appropriate transforms
 - Sometimes it can make sense to transform data so that features of it are clearer. For example, a time series of Bitcoin over time will show very little detail about the early history of the currency, when it was not valuable. However, plotting the log of the value of Bitcoin on the u-axis allows this detail to be seen.

Principle 2: Make the meaning of the data clear

- Histograms and density plots: These univariate plots are A visualisation is meaningless if it's not labelled. Every plot should have

- · Axis labels as English words
- · Units given, where appropriate (e.g. "Length (mm)" not just "Length")
- · All variables labelled e.g. a legend indicating the colours used to represent squirrel sex
- Use graphical and textual annotation e.g. it can be helpful to highlight a time series with events that you know about

Principle 3: Avoid distorting what the data have to say

Choices in visualisation design can lead to the instant impression given by preattentive processing of the visualisation being quite different to the numbers in the dataset. Tufte (1982) measures the level of distortion in a visualisation by the "Lie factor":

$$Lie factor = \frac{size \text{ of effect shown in graphi}}{size \text{ of effect in data}}$$

data:

· Use appropriate scales and baselines

- A very common problem is that the baseline (i.e. the lowest point on the y-axis) in a barchart is not zero. This can lead to small differences appearing large.
- · Be aware of limitations of our perception of size
- Although marker area can be useful for indicating categories, humans are not very good at relating the area to a quantity - we are much better at comparing lengths.

Principle 4: Make the data accessible

A visualisation is meaningless if it's illegible and loses impact if it's difficult to read. To ensure data is accessible:

- · Make sure text is legible, i.e. font size of minimum 8 points in a PDF, or about 20 points in a presentation. (It is surprising how often talks are given in which it's impossible to read the labels on plots even from the front row.)
- · Use colours that work for people with colour-vision deficiency. Wexler et al. (2017). Chapter 1 has an excellent introduction to using colour in visualisations.

Principle 5: Focus on the content

Give the viewer's brain as little work to do as possible.

- · No chartjunk e.g. colours that don't have any meaning.
- · Reduce clutter
- · Consistent colours between plots in a study
- · Correct spelling

References

Tufte, E. (1982). The visual display of auantitative information. Graphics Press, Cheshire, Connecticut,

Wexler, S., Shaffer, J., and Cotgreave, A. (2017). The Big Book of Dashboards: Visualizing Your Data Using Real-World Business Scenarios. Wiley.

Not discussed

- Rules for visualisation?
- Application versus derivation?
- Lab development

Acknowledgements

Kobi Gal

- Sharon Goldwater, Heather Yorston
- Chris Williams
- The marvellous team of FDS tutors and lab demonstrators, including Jonathan Feldstein on Lab development
- ILTS
- Other colleagues
- My family