

Obligatorisk læringsaktivitet i introuger Flow 1, Dataanalyse

Afleveringsfrist:

fredag den 13. september 2024, kl. 16 Udarbejdet af:

Thorbjørn Baum, adjunkt Cphbusiness Thorbjørn Wulf, adjunkt Cphbusiness



Indholdsfortegnelse

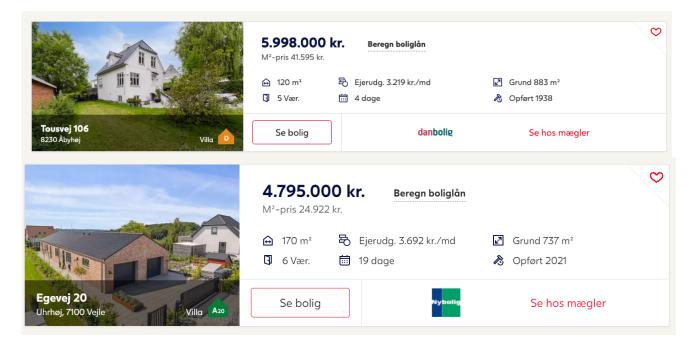
Opgave 1 – Data Science Modellen	3
Opgave 1.1 – find data	3
Opgave 1.2 - Vælg	
Opgave 1.3 – Beskriv variabler	3
Opgave 1.4 – Research goal	3
Opgave 2 – Korrelation og simple lineær regression	4
Opgave 2.1 – Beskrivende statistik	
Opgave 2.2 - Korrelation	
Opgave 2.3 – Simple regressioner	
Opgave 2.4 – Korrelation og lineær regression	
Opgave 3 – Tilfældigheder og terninger	5
Opgave 3.1 – Funktion til terninger	5
Opgave 3.2 – Plot I	5
Opgave 3.3 – Plot II	5
Opgave 3.4 – Lav dine egne data	5
Opgave 4 – Danskernes forhold til alkohol	6
Opgave 4.1 – Hent data	
Opgave 4.2 – Korrelation	
Opgave 4.3 – Kritisk tænkning	
Opgave 5 – Dataframes	7
Opgave 5.1 – Månedlige observationer	7
Opgave 5.2 – Kvartalsvise observationer	7
Opgave 5.3 - Pivot	



Opgave 1 – Data Science Modellen

Opgave 1.1 – find data

Data på ejendomme til salg i opgave 1 er fremkommet via webscapping af boligsiden.dk. Der ligger en csv-fil i mappen med OLA-opgaven. Her er to billeder fra boligsiden.dk. Find de to rækker i csv-filen, som matcher de to huse.



Opgave 1.2 - Vælg

Udvælg 2 ejendomme fra csv-filen og find dem på boligsiden.dk.

Opgave 1.3 – Beskriv variabler

Forklar NA-værdierne i csv-filen ud fra, hvad I har observeret i opgave 1.1 og 1.2. Derudover gør rede for variable I mener mangler i csv-filen sammenlignet med boligsiden.dk. (Hint: Hvordan vil I unikt identificere en bolig til salg via boligsiden?)

Opgave 1.4 – Research goal

Med udgangspunkt i Data Science Modellen skal I gøre rede for de skridt, der er blevet taget for at nå frem til csv-filen. I skal komme med et bud på et "research goal" som kunne have optimeret processen med at fremskaffe data fra boligsiden.



Opgave 2 – Korrelation og simple lineær regression

Opgave 2.1 – Beskrivende statistik

Lav beskrivende statistik for data på boliger til salg via boligsiden.dk fra opgave 1. Vær opmærksom på, hvilke variable I mener kan betragtes som x'er og y i en række af simpel lineære regressioner.

Opgave 2.2 - Korrelation

Hvad er korrelationen mellem m² og prisen for boliger lagt på boligsiden.dk? Giv en forklaring på begrebet korrelation.

Opgave 2.3 – Simple regressioner

Lav minimum 5 simple regressioner mellem pris pr. m² og 5 andre variable i csv-filen fra opgave 1. Giv en forklaring på, hvilken af de fem modeller, der bedst forklarer pris pr. m². Derudover skal I understøtte jeres 5 modeller med et bud på om, der er sammenhæng (korrelation) mellem jeres 5 udvalgte variable (hint: lav en korrelationsmatrix).

Opgave 2.4 – Korrelation og lineær regression

Forklar den teoretiske sammenhæng mellem korrelation og simple lineær regression.



Opgave 3 – Tilfældigheder og terninger

Opgave 3.1 – Funktion til terninger

Lav et script (en funktion) i R-studio, der kan slå med 25.000 terninger. Hvor mange 5'ere har jeres terning slået? Hvad er sandsynligheden, givet jeres resultatet med de 25.000 terninger, for, at jeres script slår en 5'er?

Opgave 3.2 – Plot I

Lav et script (en funktion) i R-studio, der kan slå med 6 terninger og vise summen. Slå nu 10.000 gange med de 6 terninger og lav et barplot af jeres resultat. (hint: barplot kræver fx pakken "ggplot2"). Forklar om jeres resultat giver mening i forhold til jeres funktion (hint: I bruger sample funktionen).

Opgave 3.3 – Plot II

Brug jeres script fra 3.2 og slå nu 1.000.000 gange med de 6 terninger. Lav igen et barplot og sammenlign med jeres plot fra 3.2.

Opgave 3.4 – Lav dine egne data

Lav et script i R-studio, der viser en tilfældigt opstillet række af tallene 1, 2, 3, 5, 6. Lav en matrix med to kolonner og fem rækker, hvor den første kolonne skal være tallene 2 til 6 og den anden kolonne skal være jeres tilfældige række af tallene 1, 2, 3, 5, 6. (hint: Google funktionen cbind(), der sætter lige lange kolonner sammen i en matrix).



Opgave 4 – Danskernes forhold til alkohol

Opgave 4.1 – Hent data

Hent data fra tabel FU02, alle forbrugsgrupper under 02.1 (alkoholiske drikkevarer) i faste priser for perioden 2000 til 2022 og indlæs i R. Illustrer udviklingen i de enkelte grupper.

Opgave 4.2 – Korrelation

Lav en korrelationsmatrix over forbrugsgrupperne under 02.1 og konkludér på resultaterne

Opgave 4.3 – Kritisk tænkning

Forhold jer kritisk til de indlæste data når I reflekter over jeres resultater.



Opgave 5 – Dataframes

Opgave 5.1 – Månedlige observationer

Lav en 36 x 3 dataframe med kolonnenavne "Klasse", "Uge", "Score". Første kolonne skal fyldes med A,B,C,D så der startes med 9 A'er, derpå 9 B'er osv. Anden kolonne skal fyldes med tallene 1 til 9, der gentages for hvert bogstav. Sidste kolonne skal fyldes med observationer. Det er op til jer, hvilke værdier I vil putte i framen. (hint: benyt R-funktionen *seq()*)

Opgave 5.2 – Kvartalsvise observationer

I skal lave en ny dataframe, der er 9x3 og bygger på den dataframe I lavede opgave 5.1. I skal tage udgangspunkt i framen fra opgave 5.1 og loope igennem. I loopet skal I hver tredje gang lave en dataframe 1x3 dataframe med samme navne som i 5.1. (Hint: brug modulo-operatoren til at ramme hver tredje). Indholdet skal være som følger: 1 element og 2 element henter I fra 5.1-framen. Det sidste element skal være gennemsnittet af de forrige tre observationer.

Fx I går fra denne (vær opmærksom på det er et eksempel på kvartal):

-	INDIKATOR =	TID [‡]	value ‡	Over til:	_	INDIKATOR *	TID ‡	value 🗦
1	Α	2000-01-01	5.079502		1	Α	2000-04-01	4.861433
	Α	2000-02-01			2	Α	2000-08-01	5.415208
					3	Α	2000-12-01	4.530170
3	Α	2000-03-01	3.543534		4	Α	2001-04-01	4.738251
4	Α	2000-04-01	4.218260		5	Α	2001-08-01	4.790563
5	Α	2000-05-01	5.320402		6	Α	2001-12-01	3.779561
6	Α	2000-06-01	4.555218		7	В	2000-04-01	4.687455
7	Α	2000-07-01	6.370004		8	В	2000-08-01	5.366524
	Α	2000-08-01	5 673254		9	В	2000-12-01	4.474418
	<i>A</i>	2000 00 01	3.073231		10	В	2001-04-01	4.240166
24	Α	2001-12-01	5.840540		11	В	2001-08-01	4.178769
25	В	2000-01-01	4.714155		12	В	2001-12-01	4.485748
26	В	2000-02-01	5.504126		13	С	2000-04-01	4.882734
27	R	2000-03-01	3 844083		14	С	2000-08-01	5.447672
					15	С	2000-12-01	5.016085
28		2000-04-01	4.0/2851		16	С	2001-04-01	4.850881
					17	С	2001-08-01	4.004480
					18	С	2001-12-01	5.427528

Opgave 5.3 - Pivot

I skal nu konvertere denne nye dataframe til en ny dataframe som har følgende navne på kolonnerne:

"Uge","A","B","C",D" og rækkerne indeholder de gennemsnit som I har beregnet. (Hint: Brug funktionen pivot-wider fra pakken *tidyr*).

•	TID ‡	A =	₿ ‡	c ÷
1	2000-04-01	4.861433	4.687455	4.882734
2	2000-08-01	5.415208	5.366524	5.447672
3	2000-12-01	4.530170	4.474418	5.016085
4	2001-04-01	4.738251	4.240166	4.850881
5	2001-08-01	4.790563	4.178769	4.004480
6	2001-12-01	3.779561	4.485748	5.427528