

Lineær regression, kvadratsummer og forklaringsgrad

Når man udfører lineær regression, søger man som bekendt den lineære funktion på formen $f(x) = a \cdot x + b$, som minimerer følgende kvadratsum:

$$(1) \quad SS_{res} = \sum_{n=1}^n (y_i - f(x_i))^2$$

Opgaven er altså at finde a og b , så summen bliver mindst mulig. Forskellen mellem punkternes y -værdier og funktionsværdierne for den lineære funktion i de tilhørende x -værdier kaldes for *residualerne* og beregnes med r_i . Vi har altså:

$$(2) \quad r_i = y_i - f(x_i)$$

Med denne definition har vi altså, at

$$(3) \quad SS_{res} = \sum_{n=1}^n r_i^2$$

Sætning 1 (Lineær regression)

Givet en række datapunkter $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. Den lineære funktion på formen $f(x) = a \cdot x + b$, som minimerer kvadratsummen (1) er givet ved:

$$(4) \quad a = \frac{\sum_{n=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum_{n=1}^n (x_i - \bar{x})^2}$$

$$(5) \quad b = \bar{y} - a \cdot \bar{x}$$

hvor \bar{x} og \bar{y} er middelværdierne af henholdsvis x -værdierne og y -værdierne:

$$(6) \quad \bar{x} = \sum_{n=1}^n x_i$$

$$(7) \quad \bar{y} = \sum_{n=1}^n y_i$$

Bevis: Vi skal ikke give det her. Det kan findes mange steder. □

I dette tillæg skal vi kigge på to andre kvadratsummer end (1), nemlig:

$$(8) \quad SS_{tot} = \sum_{n=1}^n (y_i - \bar{y})^2$$

$$(9) \quad SS_{reg} = \sum_{n=1}^n (f(x_i) - \bar{y})^2$$

Før vi kommer til en interessant identitet involverende de tre kvadratsummer, skal vi først se på nogle andre interessante hjælpesætninger.

Sætning 2

Middelværdien af residualerne er 0:

$$(10) \quad \sum_{n=1}^n r_i = 0$$

Bevis:

$$\begin{aligned} \sum_{n=1}^n r_i &= \sum_{n=1}^n (y_i - f(x_i)) = \sum_{n=1}^n (y_i - a \cdot x_i - b) = \sum_{n=1}^n y_i - a \cdot \sum_{n=1}^n x_i - n \cdot b \\ &= n \cdot \bar{y} - a \cdot n \cdot \bar{x} - n \cdot (b - a \cdot \bar{x}) = 0 \end{aligned}$$

hvor vi i fjerde lighedstegn har benyttet (5). □

Sætning 3

$$(11) \quad \sum_{n=1}^n f(x_i) \cdot (x_i - \bar{x}) = \sum_{n=1}^n y_i \cdot (x_i - \bar{x})$$

Bevis: Vi regner på venstresiden:

$$\begin{aligned} \sum_{n=1}^n f(x_i) \cdot (x_i - \bar{x}) &= \sum_{n=1}^n (a \cdot x_i + b) \cdot (x_i - \bar{x}) \\ &= \sum_{n=1}^n (a \cdot x_i + \bar{y} - a \cdot \bar{x}) \cdot (x_i - \bar{x}) \\ &= \sum_{n=1}^n (\bar{y} + a \cdot (x_i - \bar{x})) \cdot (x_i - \bar{x}) \\ &= \bar{y} \cdot \sum_{n=1}^n (x_i - \bar{x}) + a \cdot \sum_{n=1}^n (x_i - \bar{x})^2 \\ &= \bar{y} \cdot \sum_{n=1}^n (x_i - \bar{x}) + \sum_{n=1}^n (y_i - \bar{y}) \cdot (x_i - \bar{x}) \\ &= \bar{y} \cdot \sum_{n=1}^n (x_i - \bar{x}) + \sum_{n=1}^n y_i \cdot (x_i - \bar{x}) - \bar{y} \cdot \sum_{n=1}^n (x_i - \bar{x}) \\ &= \sum_{n=1}^n y_i \cdot (x_i - \bar{x}) \end{aligned}$$

hvor vi har brugt følgende: 1. lighedstegn: (5). 4. lighedstegn: (4). □

Sætning 4

$$(12) \quad \sum_{i=1}^n r_i \cdot (x_i - \bar{x}) = 0$$

Bevis: Fås straks ved at trække venstresiden fra højresiden i (11) fra sætning 3:

$$0 = \sum_{i=1}^n (y_i - f(x_i)) \cdot (x_i - \bar{x}) = \sum_{i=1}^n r_i \cdot (x_i - \bar{x})$$

□

Sætning 5

$$(13) \quad \sum_{i=1}^n r_i \cdot x_i = 0$$

Bevis: Vi bruger resultatet i sætning 4:

$$0 = \sum_{i=1}^n r_i \cdot (x_i - \bar{x}) = \sum_{i=1}^n r_i \cdot x_i - \bar{x} \cdot \sum_{i=1}^n r_i = \sum_{i=1}^n r_i \cdot x_i$$

hvor vi for at få sidste lighedstegn har benyttet sætning 2.

□

Sætning 6 (Sætning med kvadratsummer)

$$(14) \quad SS_{tot} = SS_{reg} + SS_{res}$$

Bevis: Vi regner på venstresiden:

$$\begin{aligned} SS_{tot} &= \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \sum_{i=1}^n [y_i - f(x_i) + f(x_i) - \bar{y}]^2 \\ &= \sum_{i=1}^n (y_i - f(x_i))^2 + \sum_{i=1}^n (f(x_i) - \bar{y})^2 + 2 \cdot \sum_{i=1}^n (y_i - f(x_i)) \cdot (f(x_i) - \bar{y}) \\ &= SS_{res} + SS_{reg} + 2 \cdot \sum_{i=1}^n r_i \cdot (a \cdot x_i + b - \bar{y}) \\ &= SS_{res} + SS_{reg} + 2a \cdot \sum_{i=1}^n r_i \cdot x_i + 2 \cdot (b - \bar{y}) \cdot \sum_{i=1}^n r_i \\ &= SS_{res} + SS_{reg} \end{aligned}$$

hvor vi i 4. lighedstegn har udnyttet definitionerne (1), (9) samt (2). For at få sidste lighedstegn har vi udnyttet hjælpesætningerne 2 og 5.

□

Definition 7 (Forklaringsgraden)

Givet følgende sæt af datapunkter: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. Forklaringsgraden R^2 for den regressionslinje, som hører til datapunkterne, er defineret ved:

$$(15) \quad R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

Vi skal give en fortolkning af denne størrelse snart, men først skal vi give et par andre formler for forklaringsgraden.

Sætning 8

Der gælder følgende alternative formler for forklaringsgraden R^2 for regressionslinjen, som hører til sættet af datapunkter $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$:

$$(16) \quad R^2 = \frac{SS_{reg}}{SS_{tot}}$$

$$(17) \quad R^2 = a^2 \cdot \frac{\sum_{n=1}^n (x_i - \bar{x})^2}{\sum_{n=1}^n (y_i - \bar{y})^2} = a^2 \cdot \frac{\overline{x^2} - \bar{x}^2}{\overline{y^2} - \bar{y}^2}$$

hvor a er hældningskoefficienten for regressionslinjen.

Bevis: (16) fås direkte af sætning 6:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = \frac{SS_{tot} - SS_{res}}{SS_{tot}} = \frac{SS_{reg}}{SS_{tot}}$$

For at vise formlerne i (17), bemærker vi først følgende:

$$(18) \quad f(x_i) - \bar{y} = a \cdot x_i + b - \bar{y} = a \cdot x_i + (\bar{y} - a \cdot \bar{x}) - \bar{y} = a \cdot (x_i - \bar{x})$$

Vi regner videre på (16):

$$(19) \quad R^2 = \frac{SS_{reg}}{SS_{tot}} = \frac{\sum_{n=1}^n (f(x_i) - \bar{y})^2}{\sum_{n=1}^n (y_i - \bar{y})^2} = \frac{\sum_{n=1}^n a^2 \cdot (x_i - \bar{x})^2}{\sum_{n=1}^n (y_i - \bar{y})^2} = a^2 \cdot \frac{\sum_{n=1}^n (x_i - \bar{x})^2}{\sum_{n=1}^n (y_i - \bar{y})^2}$$

hvorved den første identitet er vist. For at vise den sidste i (17), omskriver vi først tælleren i brøken:

$$\begin{aligned}\sum_{n=1}^n (x_i - \bar{x})^2 &= \sum_{n=1}^n (x_i^2 - 2 \cdot \bar{x} \cdot x_i + \bar{x}^2) = \sum_{n=1}^n x_i^2 - 2 \cdot \bar{x} \cdot \sum_{n=1}^n x_i + \bar{x}^2 \cdot \sum_{n=1}^n 1 \\ &= \sum_{n=1}^n x_i^2 - 2 \cdot \bar{x} \cdot n \cdot \bar{x} + \bar{x}^2 \cdot n = \sum_{n=1}^n x_i^2 - n \cdot \bar{x}^2 = n \cdot (\overline{x^2} - \bar{x}^2)\end{aligned}$$

Nævneren omskrives helt på samme vis, hvorefter vi kan regne videre på (19):

$$R^2 = a^2 \cdot \frac{\sum_{n=1}^n (x_i - \bar{x})^2}{\sum_{n=1}^n (y_i - \bar{y})^2} = a^2 \cdot \frac{n \cdot (\overline{x^2} - \bar{x}^2)}{n \cdot (\overline{y^2} - \bar{y}^2)} = a^2 \cdot \frac{\overline{x^2} - \bar{x}^2}{\overline{y^2} - \bar{y}^2}$$

Herved er den sidste identitet vist. □

Vi er rede til at give en fortolkning af forklaringsgraden direkte ud fra definition 17. Den *totale kvadratsum* (Total Sum of Squares) SS_{tot} , den *residuale kvadratsum* SS_{res} (Residual Sum of Squares) og *regressions-kvadratsummen* SS_{reg} (Regression Sum of Squares) skal i spil her. SS_{tot} beskriver y -koordinaternes variation i forhold til middelværdien af disse. Den er faktisk n gange den empiriske varians af y -værdierne. SS_{res} beskriver y -koordinaternes variation i forhold til regressionslinjens y -værdier, altså kvadratsummen af residualerne. Bemærk, at det er den størrelse, som er minimeret ved regressionsprocessen! Den sidste størrelse SS_{reg} handler om variationen af de y -værdier, som regressionslinjen forudsiger (altså $f(x_i)$ 'erne) set i forhold til middelværdien af y -værdierne. Den kaldes også for *den forklarende kvadratsum* (Explained sum of Squares), fordi den handler om den naturlige variation, som regressionslinjen kan forklare.

Hvis vi kigger på definition 7 for forklaringsgraden, så kan brøken SS_{res}/SS_{tot} fortolkes som den brøkdel af den totale variation, som *ikke* kan forklares, stammende fra residualerne. Brøken skal trækkes fra 1 for at få forklaringsgraden. Det er derfor ikke mærkeligt, at man tolker forklaringsgraden som den brøkdel af den totale variation, som *kan forklares*. Det stemmer samtidigt fint med den alternative formel (16) for forklaringsgraden.

$$(20) \quad R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = \frac{SS_{reg}}{SS_{tot}}$$

Hvis forklaringsgraden i et eksempel fås til $R^2 = 0,85$, så siger vi, at regressionslinjen kan forklare 85% af variationen i y -værdierne. Der er altså 15% af variationen, som regressionslinjen *ikke* kan forklares. Vi ser samtidigt, at hvis $R^2 = 1$, så ligger alle punkterne på regressionslinjen, mens $R^2 = 0$ betyder, at $SS_{res} = SS_{tot}$ og $SS_{reg} = 0$. I øvrigt vil forklaringsgraden altid være et tal mellem 0 og 1. Kvadratsummerne er nemlig alle *ikke-negative*, hvorfor (14) i sætning 6 giver os, at SS_{tot} er mindst lige så stor som henholdsvis SS_{res} og SS_{reg} , hvorved $0 \leq R^2 \leq 1$ følger. Lige et forbehold: Hvis alle datapunkterne har samme y -værdier, så er forklaringsgraden slet ikke defineret!

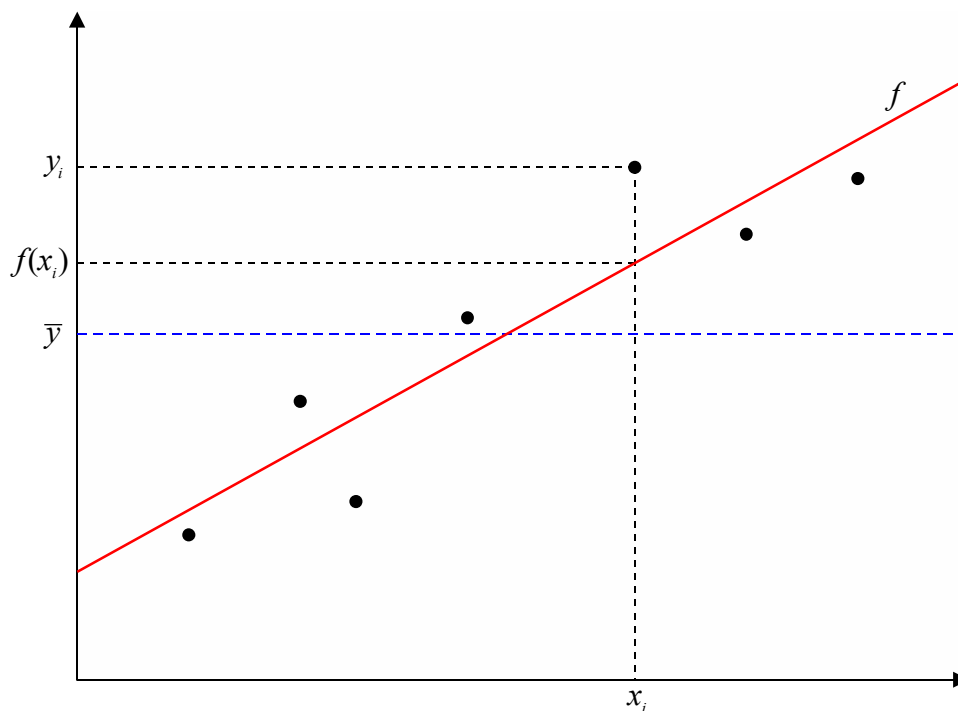
Lad os et øjeblik betragte de størrelser vi summer over i kvadratsummerne i et grafisk perspektiv: $y_i - \bar{y}$, $f(x_i) - \bar{y}$ og $y_i - f(x_i)$, hvor den sidste er det i 'te residual r_i . Der gælder klart at den første størrelse er summen af de to sidste:

$$(21) \quad y_i - \bar{y} = (f(x_i) - \bar{y}) + (y_i - f(x_i))$$

Det mest interessante er imidlertid, at hvis vi kvadrerer hvert af de tre led og summerer fra $i = 1$ til n , så gælder lighedstegnet også. Det var det, vi viste i sætning 6:

$$(22) \quad \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (f(x_i) - \bar{y})^2 + \sum_{i=1}^n (y_i - f(x_i))^2$$

der er det samme som: $SS_{tot} = SS_{reg} + SS_{res}$.



Korrelationskoefficienten

Definition 9 (Den empiriske korrelationskoefficient)

Den *empiriske korrelationskoefficient*, som beskriver samvariationen mellem x - og y -værdierne i sættet af datapunkter: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, er defineret ved:

$$(23) \quad \rho = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Sætning 10

Der gælder følgende formel for den empiriske korrelationskoefficient:

$$(24) \quad \rho = a \cdot \frac{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

hvor a er hældningskoefficienten for regressionslinjen.

Bevis: Vi omskriver udtrykket for definitionen af ρ fra definition 9 og udnytter udtrykket for hældningskoefficienten for regressionslinjen fra sætning 1:

$$\begin{aligned} \rho &= \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \frac{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= a \cdot \frac{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \end{aligned}$$

□

Sætning 11

Forklaringsgraden for datasættet $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ er kvadratet på den empiriske korrelationskoefficient:

$$(25) \quad R^2 = \rho^2$$

Bevis: Fremgår direkte ved at sammenligne udtrykkene i sætning 8 og 10.

□

Definition 12 (Empiriske spredninger)

De *empiriske spredninger*, af x - og y -værdierne i sættet $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ af datapunkter er defineret ved:

$$(26) \quad s_x = \sqrt{\frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{og} \quad s_y = \sqrt{\frac{1}{n-1} \cdot \sum_{i=1}^n (y_i - \bar{y})^2}$$

Sætning 13

For hældningskoefficienten a for linjen, som fremkommer ved lineær regression af datapunkterne $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ gælder følgende formel:

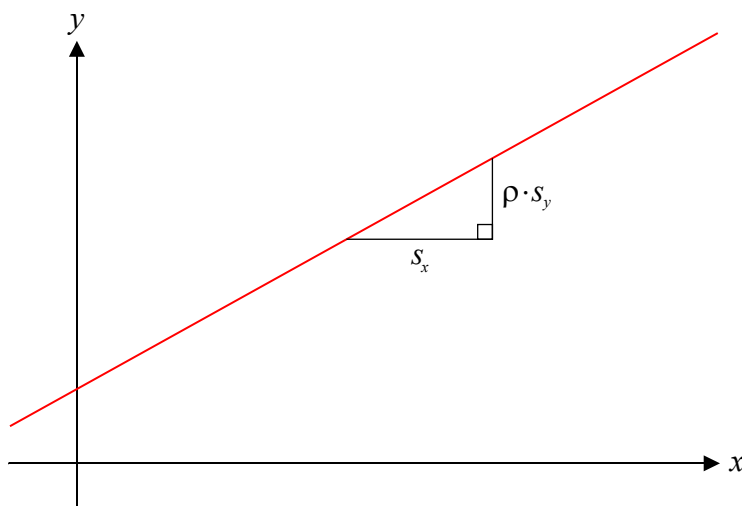
$$(27) \quad a = \rho \cdot \frac{s_y}{s_x}$$

Bevis: Det ses umiddelbart af sætning 10 og definition 12. □

Bemærkning 14

Sætning 13 ses også at gælde, selv om man ikke skulle anvende spredningerne, som er korrigeret for bias, altså selv om der i definition (26) benyttes $1/n$ i stedet for $1/(n-1)$. □

Sætning 13 kan grafisk illustreres således:



En ændring på én spredning (standardafvigelse) i x afstedkommer altså en ændring i y på korrelationskoefficienten ganget med en spredning (standardafvigelse) i y . Desuden skal det tilføjes, at den empiriske korrelationskoefficient er et tal mellem -1 og 1 : $-1 \leq \rho \leq 1$.