

## Projekt 8.4 Geogebra-eksperiment med residualspredningen: Bedste estimator for varians for residualer ved lineær regression

Eksperimentet i dette papir ligger også som video på <https://youtu.be/9xDcGeWaP7E>.

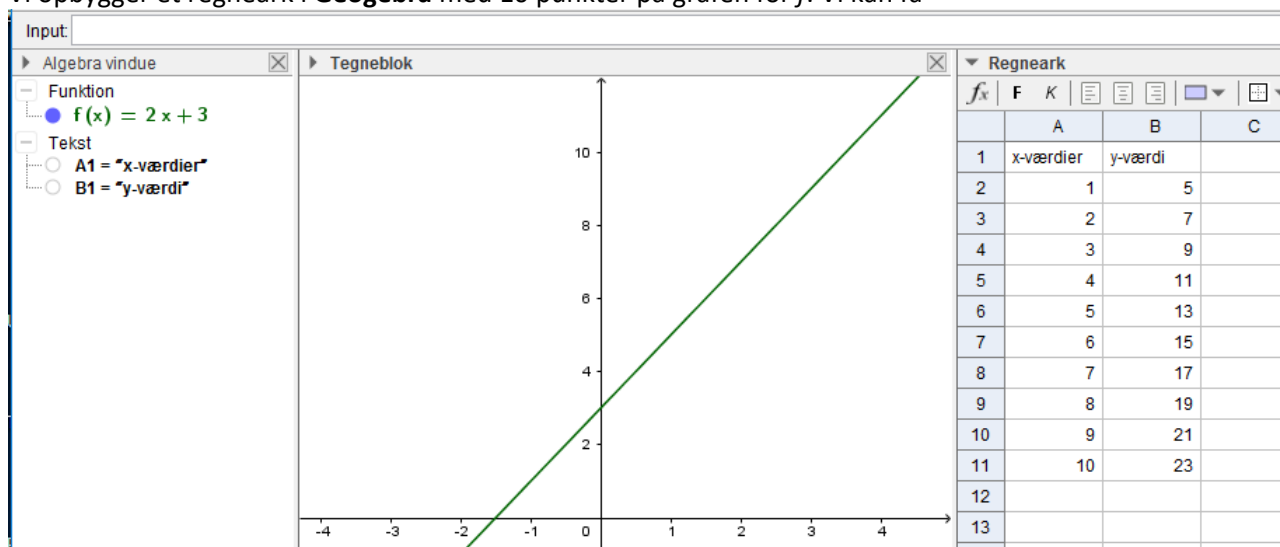
Problemstilling: Varians for residualer ved lineær regression – hvilken værdi er den bedste estimator?

Dette projekt gennemføres med værktøjet Geogebra. Du kan [her](#) finde en værktøjsuafhængig beskrivelse af, hvordan projektet kan gennemføres med andre værktøjer.

Vi regner på et datasæt med 10 talpar.

Del 1: Vi vælger her, at udgangspunktet er en lineær sammenhæng  $f(x) = 2x + 3$ .

Vi opbygger et regneark i Geogebra med 10 punkter på grafen for  $f$ . Vi kan få

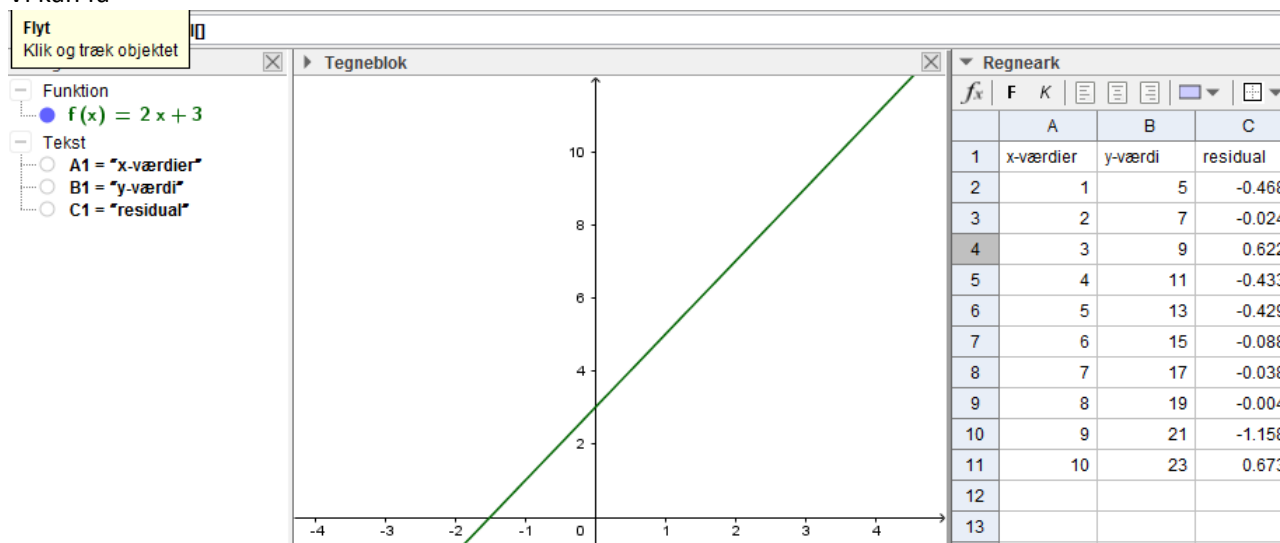


Desuden antager vi, at vi har residualer  $e_i$ , der alle er normalfordelte med  $N\left(0, \frac{1}{2\theta}\right)$ .

Dvs. den "sande" varians er  $\frac{1}{4}$ . Det er denne vi i det følgende vil antage, vi ikke kender, og som vi vil estimere.

I kolonnen "C" simulerer vi 10 residualer med kommandoen TilfældigNormal(<middelværdi>,<spredning>).

Vi kan få

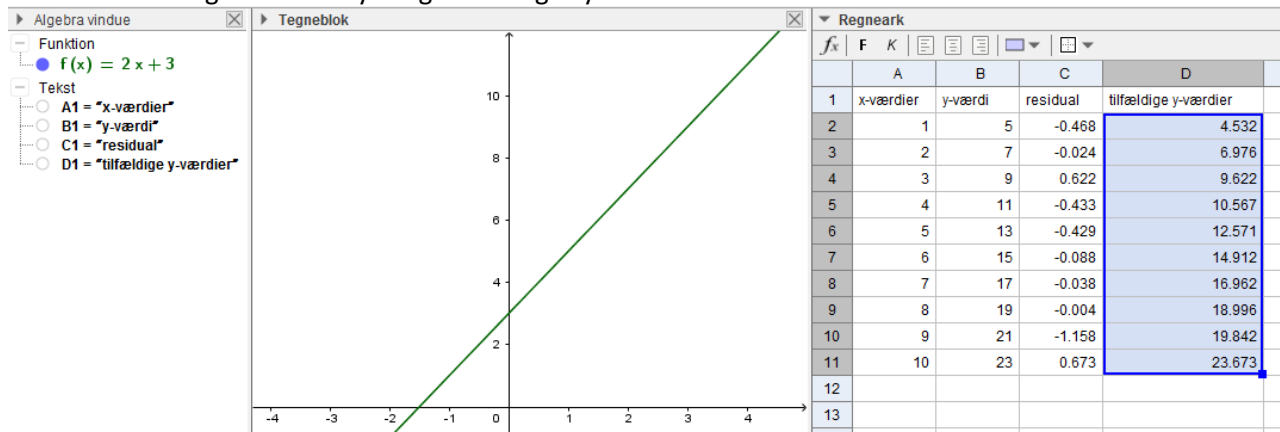


Når vi taster Ctrl R, så får vi 10 nye residualer.

Del 2: Vi vil nu simulere "nye" og "tilfældige" y-værdier, hvor den bag ved liggende lineære model er  $f$ .

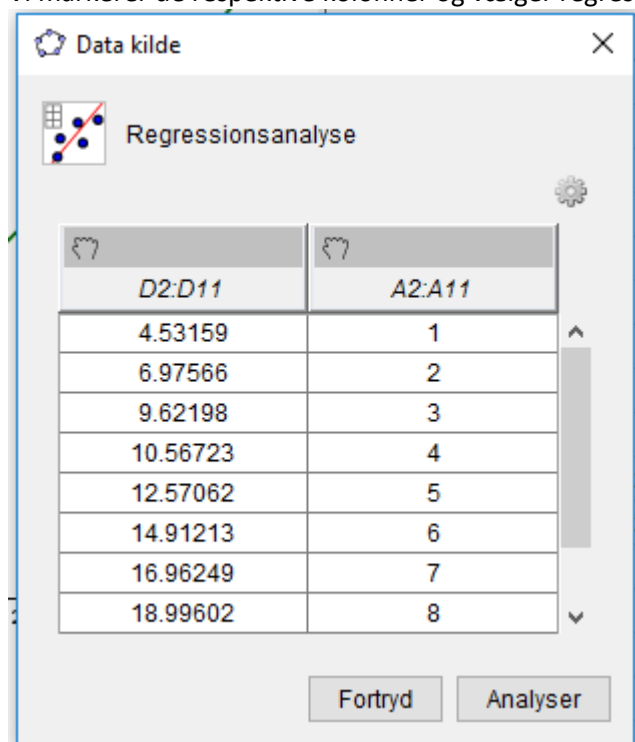
Projekter: fra kapitel 8 Projekt 8.4b Geogebra-eksperiment med residualsprejningen

I kolonne D udregner vi de "nye" og "tilfældige" y-værdier.



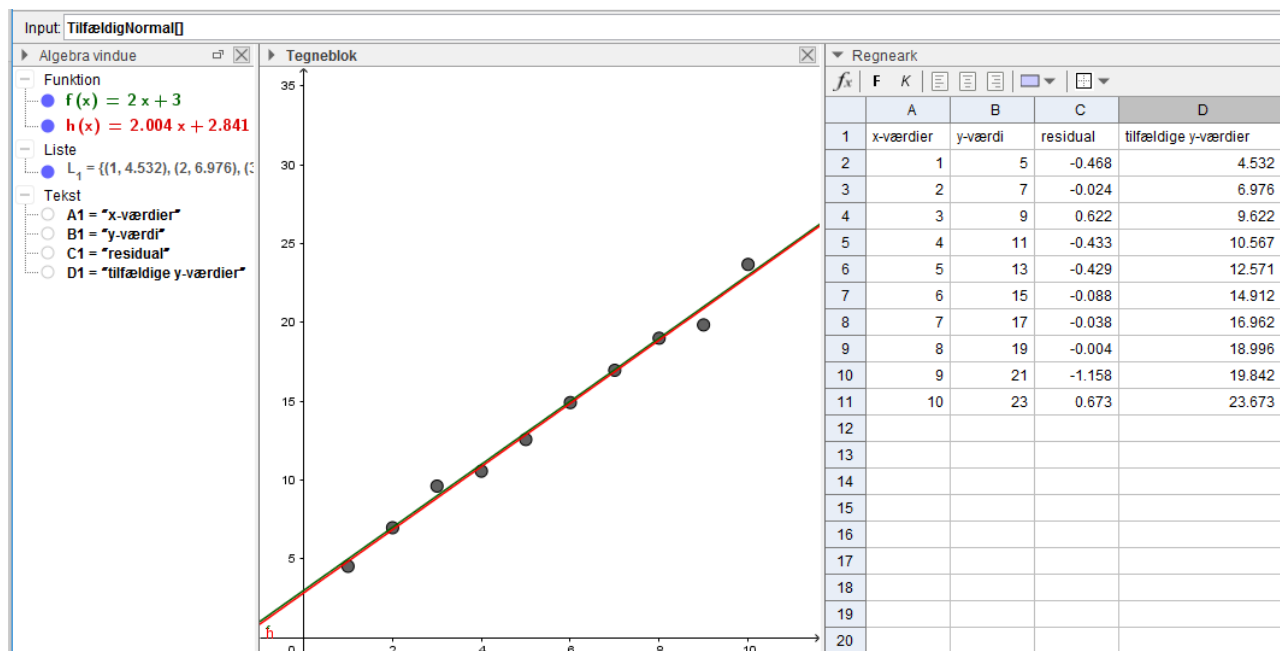
**Del 3:** Vi udfører nu lineær regression på punkterne (<x-værdier>,<tilfældige y-værdier>).

Vi markerer de respektive kolonner og vælger regressionsværktøjet. Måske får vi



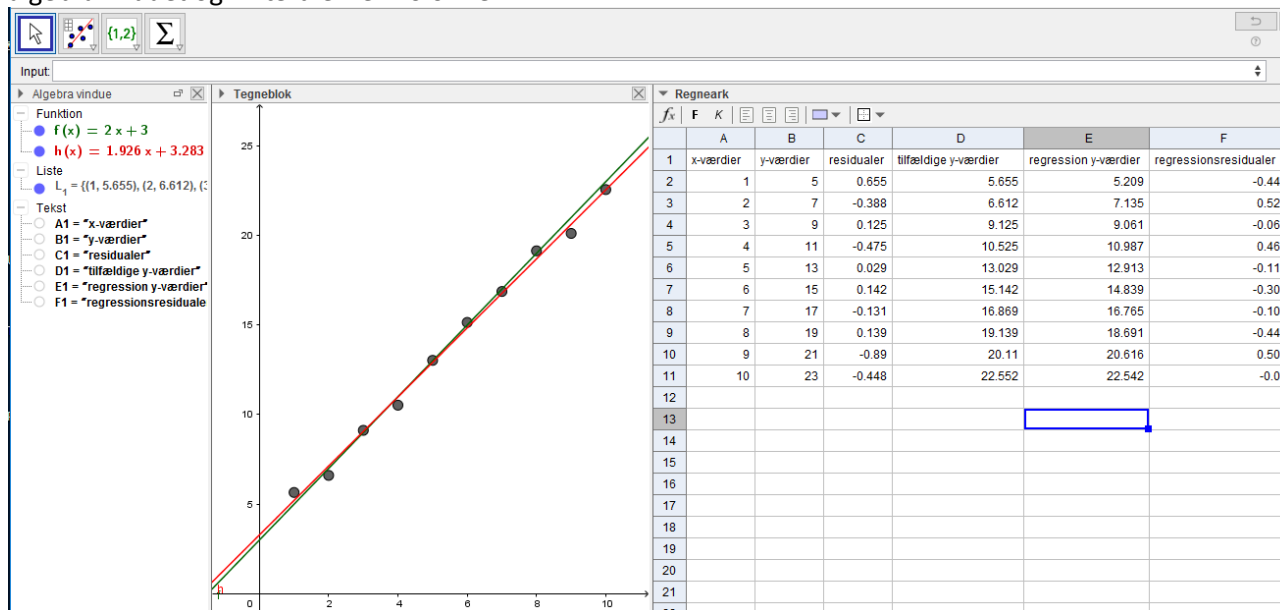
Marker kolonne A og tryk på "hånden", derefter marker kolonne D og tryk på hånden. Klik derefter på Analyser.

Hvis vi vælger "Lineær", og derefter "Kopier til tegneblok", så kan vi få:



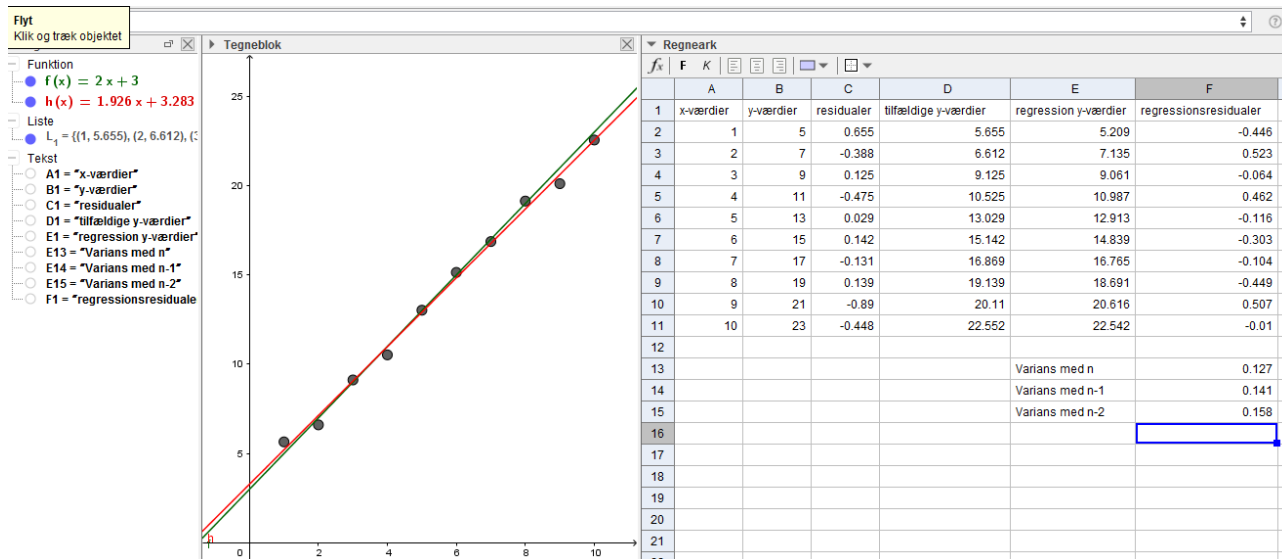
Klikker vi på Ctrl R, så laver vi hver gang en ny simulering.

**Del 4:** For hver simulation vil vi udregne regressionsmodellens y-værdier ud fra regressionsforskriften i algebravinduet og x-værdierne i kolonne A.

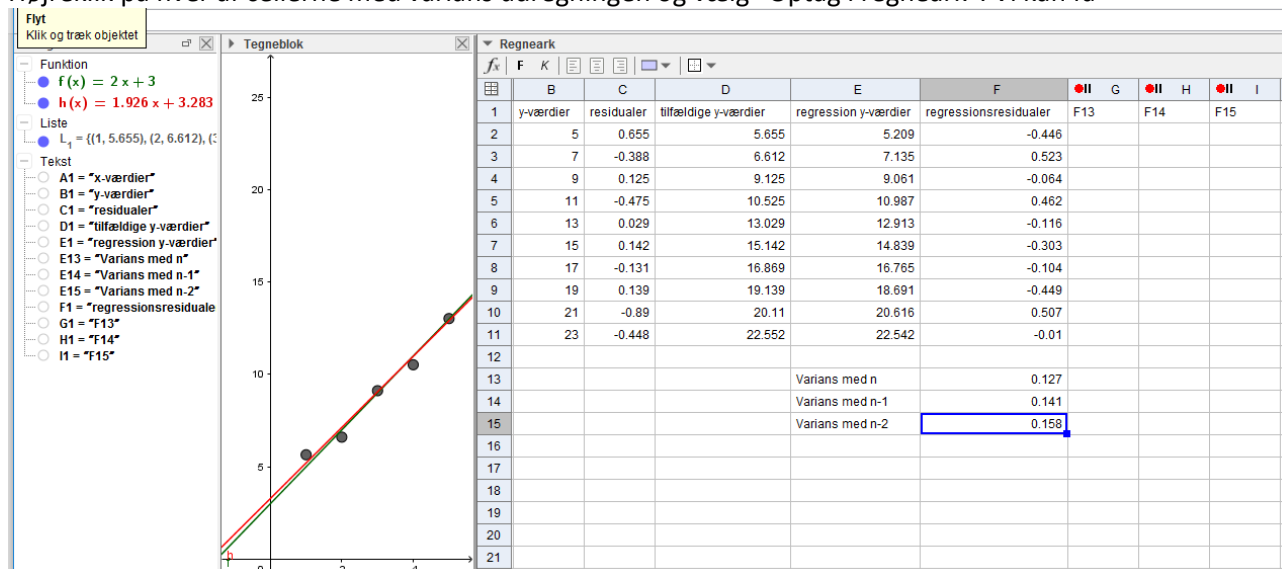


Vi kan nu udregne residualerne for regressionsmodellen i kolonne F, og derefter residualernes varians i tre udgaver. En version, hvor vi dividerer kvadratsummen med  $n$ ,  $n-1$  og  $n-2$ . Vi kan bruge kommandoen Varians(F2:F11).

Projekter: fra kapitel 8 Projekt 8.4b Geogebra-eksperiment med residualsprejningen

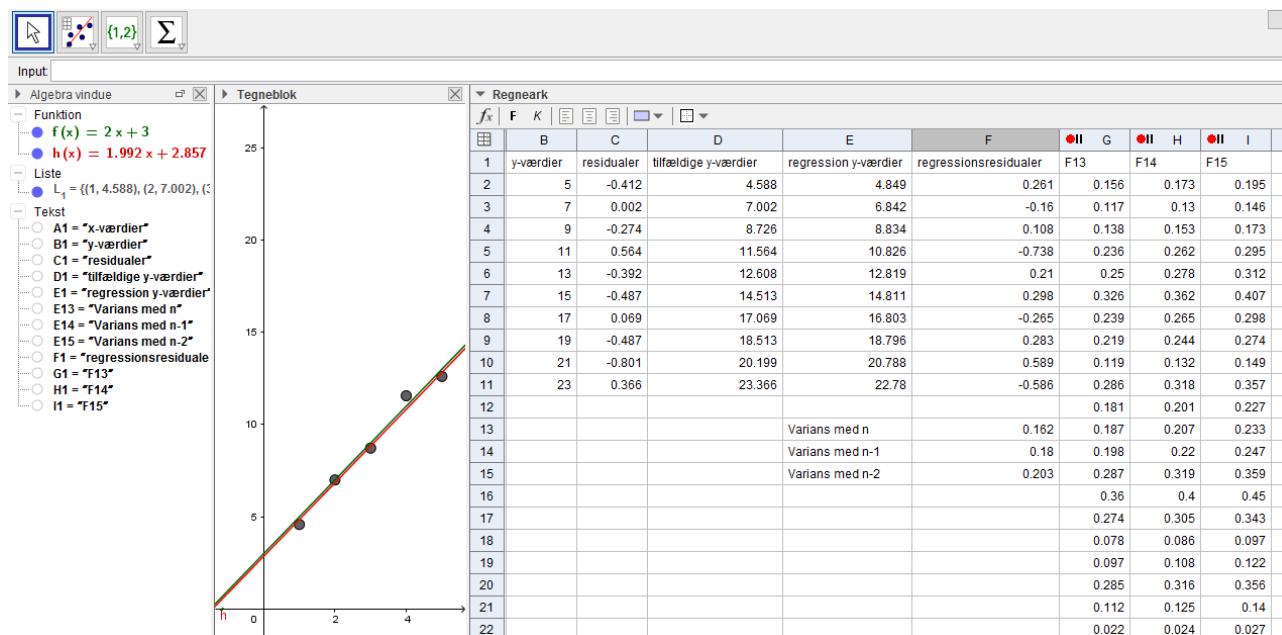


Del 5: Vi vil nu lave 1000 simuleringer af de regressionen, så vi får 1000 værdier for hver af varianserne. Højreklik på hver af cellerne med varians udregningen og vælg "Optag i regneark". Vi kan få



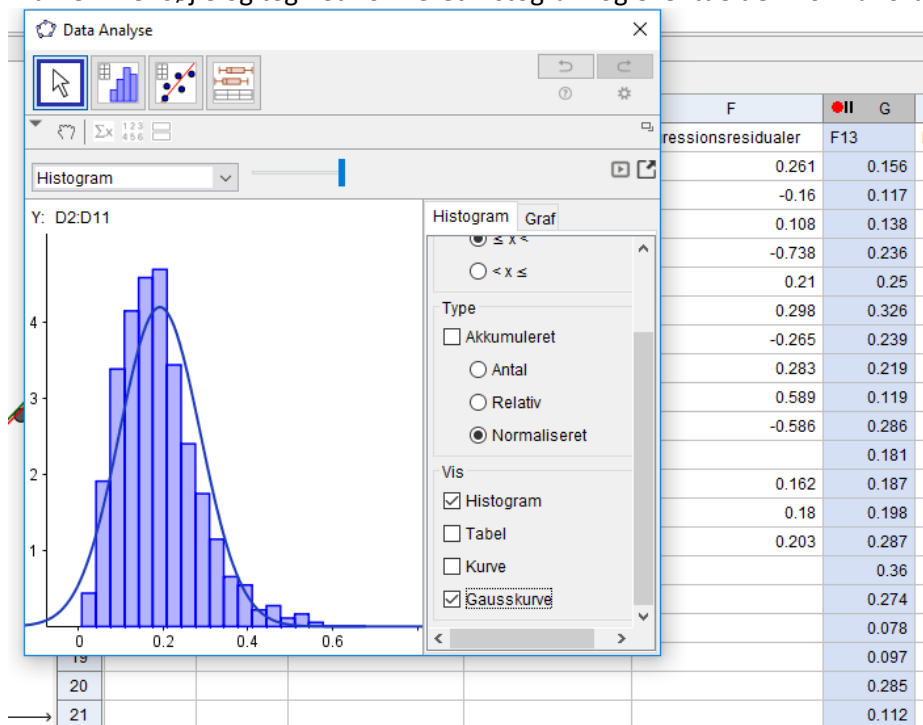
Tryk på Ctrl R mange gange, og vi kan få

Projekter: fra kapitel 8 Projekt 8.4b Geogebra-eksperiment med residualsprejningen

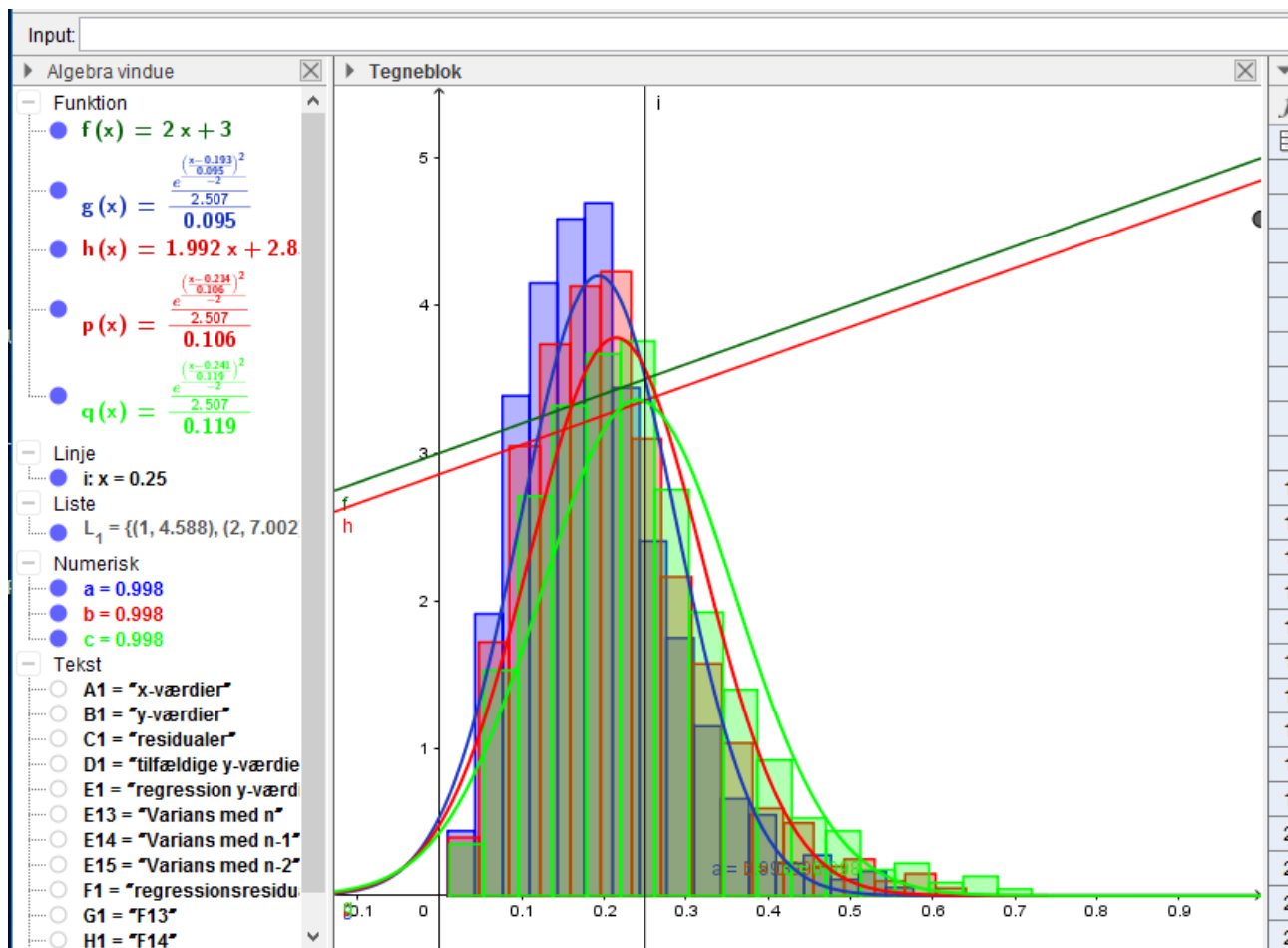


Del 6: Vi undersøge fordelingen af hver af varianserne.

Marker hver søjle og tegn et normeret histogram og eventuelt en normalfordelingskurve. Vi kan få



Vi vælger at kopiere alle histogrammer til tegneblokken, sammen med den "sande" varians. Samlet kan vi få:



Konklusion: Vores simuleringer viser, at residualvariansen, hvor vi dividerer med  $n - 2$  er den bedste estimator for den "sande" residualvariens.