

KONINKLIJK NEDERLANDS
METEOROLOGISCH INSTITUUT

Wetenschappelijk Rapport W.R. 57-001 (III-197)

Dr. C. Levert

Het normaliseren van een niet-normale verdeling.
Normalization of a non-normal distribution.

De Bilt, 1957



All Rights Reserved.

Nadruk zonder toestemming van het K.N.M.I. is verboden.

Dr. C. Levert

Het normaliseren van een niet-normale verdeling.
Normalization of a non-normal distribution.

blz.

INHOUD

2	0	Inleiding; waarom "normaliseren" ?
3	1	Het normaliseren heeft twee aspecten.
	1.1	Eerste geval; $\int f(x)dx$ is exact bekend.
	1.2	Tweede geval; $\int f(x)dx$ moet geschat worden uit de steekproef.
3	2	Eerste geval; $\int f(x)dx$ is exact bekend.
5	3	Tweede geval; $\int f(x)dx$ te schatten uit steekproef.
11	4	Numerieke voorbeelden.
15	5	Literatuur
16	6	Summary

0. Inleiding; waarom "normaliseren"?

Het komt herhaaldelijk voor, dat men "statistisch gehandicapt" is door het feit, dat de onderhavige variable x (niet een normale verdeling volgt, ook niet in voldoende benadering. In vele gevallen moet men weten, dat aan voorwaarde van normaliteit voldaan is om bepaalde statistische "methodes" (toetsingen enz.) te mogen toepassen. Enige voorbeelden:

- 0.1 Men wil een grootheid x met een andere grootheid y correleren. Niets betel natuurlijk een berekening van de klassieke correlatiecoëfficiënt r (c.c.), berustende op de sommen $\sum x, \sum y, \sum x^2, \sum y^2, \sum x.y$, doch wil men weten welke de "betekenis" deze r heeft, dan moet men de significantie-tabel raadplegen. Deze onderstelt echter, dat men reeds weet, dat zowel x - als het y - universum een normale verdeling heeft.
- 0.2 Men wil onderzoeken of een tijdreeks x vrij van persistentie, d.i. auto- of interne correlatie, is en berekent daartoe de c.c- tussen de reeksen x_1, x_2, x_3, \dots en x_2, x_3, x_4, \dots en die tussen de reeksen x_1, x_2, x_3, \dots en x_3, x_4, x_5, \dots enz. (d.w.z. autocorrelatiecoëfficiënten van orden 1, 2, enz.) Weer eist het raadplegen van de significantie-tabel, dat x normaal verdeeld is.
- 0.3 Men wil berekenen tussen welke grenzen de overschrijdingskans van zekere x ligt in een x -populatie, waarover wij slechts op basis van N metingen informatie bezitten. Nodig is dan dat de x -verdeling normaal is.
- 0.4 Men wil de t-toets op het verschil van de gemiddelden van twee reeksen metingen toepassen. Elk der populaties, waaruit de twee reeksen komen, moet dan een normale verdeling hebben.
- 0.5 Men wil de variantie-analyse toepassen. Weer is de normaliteit der verdeling een eis.
- 0.6 Men wil de toets van Bartlett toepassen (: is $\sigma_1 = \sigma_2 = \dots = \sigma_n$?), d.w.z. men wil onderzoeken of de n populaties, waaruit de n reeksen van metingen, t.w. $x_{11}, x_{12} \dots x_{1,a_1}$ (a_1 stuks); $x_{21}, x_{22} \dots x_{2,a_2}$ (a_2 stuks) ... $x_{n1}, x_{n2} \dots x_{n,a_n}$ (a_n stuks) stammen, ongelijke varianties bezitten. De toets onderstelt, dat men weet dat deze populaties normale verdelingen hebben.

*) De grootheid x is onderstreept als ze een waarschijnlijkheidsverdeling volgt (een stochastische of statistische variable, een "stochastiek") Een speciale waarde ervan is niet onderstreept.

- 0.7 Men wil de z - of F-toets toepassen. (Zie "Statistische Toetsingsmethoden" R III 120 1953)
- 0.8 De toepassing van vele "uitwerpskriteria" eist, dat het universum normaal verdeeld is (de parameter vrije criteria eisen dit niet). Voor een uitschieterskriterium en een daarop gebaseerd nomogram, zie (1), waarin tevens een lijst met artikelen over uitschieterskriteria opgenomen is.
- 0.9 Men wil aan de vorm van de transformatiefunctie $z = \psi(x)$, zie verderop, zekere fysische beschouwingen vastknopen.

In sommige gevallen, althans bij "toetsingen", is het mogelijk om parameter-vrije (d.i. verdelingsvrije) toetsen te gebruiken, waarbij geen (of zeer weinige en weinig strenge) onderstellingen gemaakt worden, maar men dient steeds te bedenken, dat deze toetsingen dientengevolge minder onderscheidend zijn dan de niet-parameter vrije toetsen, die men toepassen mag wanneer men op grond van neveninformatie of van afzonderlijke onderzoeken wel weet dat aan deze voorwaarden gehoorzaamd wordt.

1. Het normaliseren heeft twee aspecten.

Het normaliseren heeft twee verschillende aspecten:

- 1. de kansverdeling in het universum is exact bekend en
- 2. de kansverdeling is niet exact bekend.

Wij verkeren in de praktijk gewoonlijk in het tweede geval. Omdat de behandeling van het tweede geval het duidelijkst is na die van het eerste, beginnen wij met aan het eerste aandacht te schenken, ook al maakt dit wellicht een zeer theoretische, academische indruk.

1.1 Het eerste geval; $f(x)dx$ is exact bekend.

Denk gegeven de exacte distributieve kansverdeling $f(x)dx$ (de z.g. kansdichtheid) of de (cumulatieve) kansverdeling (onderschrijdingskans) $P(x \leq x) = \int_a^x f(x)dx = F(x)$ in het universum van de x -variable, beide gedefinieerd voor een gegeven traject: $a \leq x \leq b$. Met $-\infty < a < \infty$ en $-\infty < b < \infty$. Dit betekent: $\int_a^b f(x)dx = 1$, d.i. $F(b) = 1$ en $F(a) = 0$.; $F(x)$ zal een monotoon stijgende functie zijn en natuurlijk is $f(x) \geq 0$ voor elke $a \leq x \leq b$. Beide verdelingen kunnen gegeven zijn op twee manieren:

- 1. d.m.v. een analytische uitdrukking of voor de kansdichtheid of voor de kansverdeling. Enkele voorbeelden:

a) $f(x)dx = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-8)^2}{18}\right\} dx$, d.i. een normale verdeling $N(8;3)$, waarvoor $-\infty < x < \infty$ en $\int_{-\infty}^{\infty} f(x)dx = 1$ $\mu = \int_{-\infty}^{\infty} x f(x)dx = 8$ en variantie = $\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x)dx = 3^2$

b) $f(x)dx = \frac{x^{\frac{1}{2}n-1} e^{-\frac{1}{2}x}}{(\sqrt{2})^n \Gamma(\frac{1}{2}n)} dx$, d.i. de χ^2 -verdeling met n graden van vrijheid,

waarvoor $0 \leq x < \infty$ en $\int_0^{\infty} f(x)dx = 1$; $\mu = \int_0^{\infty} x f(x)dx = n$ en $\sigma^2 = \int_0^{\infty} (x-\mu)^2 f(x)dx = 2n$

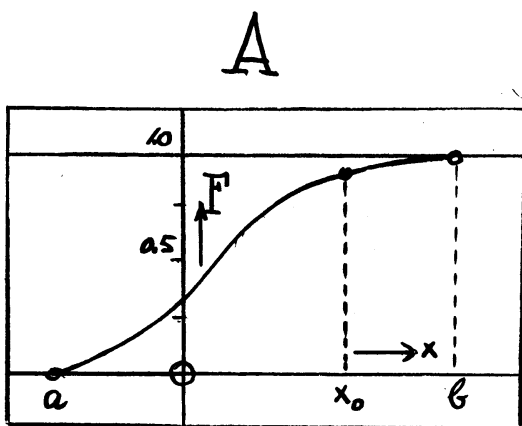
c) $f(x)dx = \frac{\exp\{-(x^2-3x+2)\}}{\int_1^{10} \exp\{-(x^2-3x+2)\} dx}$ voor $1 \leq x \leq 10$

d) $F(x)dx = 1 - e^{-(x-3)^2/4}$, d.i. een der typen goodrich-verdelingen, gedefinieerd in het traject $3 \leq x < \infty$, zodat $\int_3^x f(x)dx = 1 - \exp\{-(x-3)^2/4\}$

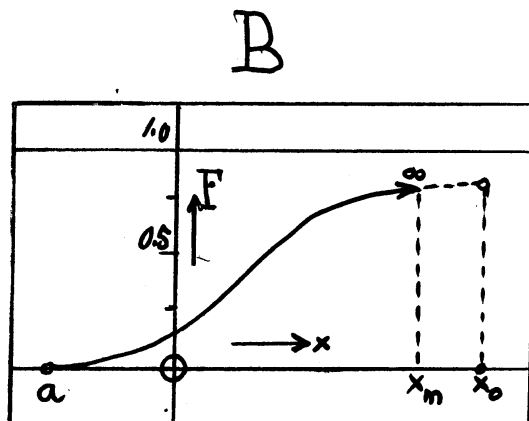
2. als een kromme, getekend in een dubbellijnair net (langs de ene as de x in een lineaire schaal, langs de andere as de F eveneens in lineaire schaal) of in een lineair-waarschijnlijkheidsnet (langs de ene as de x in een lineaire schaal, langs de andere as de F in een niet-lineaire, z.g. gausz-schaal, waarover verderop meer).

In geval 1ste kunnen wij bij iedere gegeven x_0 , waarvoor $a \leq x_0 \leq b$, de onderschrijdingskans berekenen of als $\int_a^{x_0} f(x)dx$, indien $f(x)$ gegeven is, of als $F(x_0)$, indien $F(x)$ gegeven is, zo nauwkeurig als men maar wil. In geval 2de lezen wij de $F(x_0)$ af van het prentje (zie hieronder, A). Daarbij denken wij ons de kromme $F(x)$ tegen x over het gehele traject $a \rightarrow b$ getekend (zeker als a en b eindig is). Het kan echter zijn, bijv. als a eindig en b positief oneindig is, dat de kromme $F(x)$ tegen x niet over het gehele traject (dat dan immers oneindig uitgestrekt is) getekend gegeven is (zie hieronder, B). Beschouwt men nu een x_0 -waarde, groter dan de grootste met de kromme corresponderen x -waarde (x_m), dan wordt er niet een scherp bepaald punt op de kromme gefixeerd. In dit geval moet de getekende curve uit de hand geëxtropoleerd worden, hetgeen over kleine afstanden wellicht in voldoende mate betrouwbaar is, over grote echter niet.

De volgende schetsjes mogen een en ander verduidelijken.



a en b eindig



a eindig ; $b = +\infty$

Definitie van normale verdeling.

Onderstel x is zelf niet normaal verdeeld. Wij spreken af een variabele y normaal verdeeld te noemen, als er tenminste één stel waarden α en β bestaat, zodanig dat voor iedere specifieke waarde y_1 uit het traject $-\infty \leq y < \infty$ de kans op y gelegen tussen y_1 en $y_1 + dy_1$ geschreven kan worden als

$$g(y_1)dy_1 = \frac{1}{\beta\sqrt{2\pi}} \cdot \frac{-(y_1-\alpha)^2}{2\beta^2} dy_1$$

Dat daarbij deze constanten α en β^2 de speciale, eenvoudige betekenis hebben, t.w. gemiddelde = $\xi_y = \int_{-\infty}^{\infty} yg(y) dy$ resp. variantie = $\xi_y - \xi_y^2 = \sigma_y^2 = \int_{-\infty}^{\infty} (y-\alpha)^2 g(y) dy$, is eigenlijk min of meer bijkomstig.

Wat is normaliseren? ¹⁾

Gesteld, x is zelf niet normaal verdeeld, is er dan tenminste één functie van x , $z = \psi(x)$, zodanig, dat z wel normaal verdeeld is? Het transformeren d.m.v. $z = \psi(x)$ heet dan normaliseren; deze $\psi(x)$ heet de normalisatie-transformatie. Kapteyn en Van Uven, zie [4], spreken van "normale functie", welke term wij liever vermijden.

$$T = g(z) dz$$

N.B. Als $z = \psi(x)$ en $x = h(z) =$ inversie functie, dan is $\int f(x) dx = \int f\{h(z)\} \frac{dh(z)}{dz} dz$. De onderschrijdingskans $\int_a^{x_1} f(x) dx$ van een gegeven x_1 in de gegeven x -verdeling moet gelijk zijn aan die van $z_1 = \psi(x_1)$ in de normale z -verdeling $\int_{-\infty}^{z_1} g(z) dz$, en wel voor iedere x_1 uit het traject $a \leq x \leq b$. Door dit voorschrift is aan iedere x -waarde één z -waarde toegevoegd en omgekeerd.

1.2 Het tweede geval; $\int f(x) dx$ te schatten uit steekproef

Er zijn N "metingen" $x_1, x_2 \dots x_n$ verricht; de kansverdeling $f(x) dx$ van x in het universum is onbekend. Wij hebben misschien een vermoeden, dat ze van een bekend type is (met onbekende parameters), of zelfs dit vermoeden is er niet. Laten wij onderstellen, dat deze $f(x) dx$ geen normale verdeling voorstelt.

1) Men zou, zo men wil (maar dit is zeer ongebruikelijk), ook kunnen "Pearson-niseren", d.w.z. zoeken naar zulk een transformatie $z = \phi(x)$, dat de z verdeeld is volgens de Pearson-verdeling van type III, gedefinieerd door

$$p(z) dz = \frac{2/\beta}{\beta\sqrt{2\pi}} z^{\alpha-1} e^{-z^{2\alpha}/2\beta^2} dz \quad \text{met } 0 \leq z < \infty, \text{ en } \alpha \text{ en } \beta \text{ reële constanten.}$$

($\beta > 0$) De $z = \phi(x)$ zou dan vastgelegd zijn door het voorschrift:

$$\int_a^x f(x) dx = \int_{\phi(a)}^{\phi(x)} \frac{2/\alpha}{\beta\sqrt{2\pi}} z^{\alpha-1} e^{-z^{2\alpha}/2\beta^2} dz, \quad \text{voor iedere } a \leq x \leq b.$$

(Voortzetting volgende blz.)

De vraag of $f(\underline{x})$ zulk een vorm hebben kan, dat normaliseren überhaupt mogelijk is behoeft niet gesteld te worden, aangezien iedere kansverdeling $f(\underline{x})d\underline{x}$ genormaliseerd kan worden, wat reeds onder 1.1. bleek en in 2 duidelijker worden zal. Wel heeft het zin te vragen hoe $f(\underline{x})$ eruit moet zien opdat de normalisatie-transformatie $z = \psi(\underline{x})$ van eenvoudige gedaante (dit nader te definiëren) zal zijn. De beantwoording eist de ontwikkeling van een kriterium om uit te maken of $f(\underline{x})$ zodanig kan zijn, dat van een gegeven $z = \psi(\underline{x})$ sprake kan zijn. Zulk een kriterium zal gebaseerd moeten zijn op bepaalde functies der N metingen en relaties daartussen.

2. Eerste geval; $f(\underline{x}) d\underline{x}$ is exact bekend.

$f(\underline{x})d\underline{x}$ is gegeven over $a \leq x \leq b$, zodat $\int_a^b f(\underline{x})d\underline{x} = 1$. Bij ieder paar waarden μ_z, σ_z (met $-\infty < \mu_z < +\infty$; $0 < \sigma_z < \infty$) is de normalisatie-transformatie $z = \psi(\underline{x})$ vastgelegd (en hiermee impliciet bekend) door het voorschrift:

$$1. \quad \int_a^x f(\underline{x})d\underline{x} = \int_{\psi(a)}^{\psi(x)} \frac{1}{\sigma_z \sqrt{2\pi}} \exp. \left[-\frac{(z - \mu_z)^2}{2\sigma_z^2} \right] dz$$

Gezien het feit, dat deze betrekking moet gelden voor iedere $a \leq x \leq b$, moet $\psi(a) = -\infty$ en $\psi(b) = +\infty$ zijn, waardoor reeds twee voorwaarden aan $z = \psi(\underline{x})$ opgelegd zijn. Tegelijk zijn μ_z en σ_z het gemiddelde en de standaarddeviatie van de normale z -verdeling. Door ieder ander stel σ, μ ¹⁾ is een andere $z = \psi(\underline{x})$ vastgelegd, tussen welke transformaties een eenvoudige relatie bestaat, waarover straks meer.

zodat ook gelden moet $\phi(a) = 0$ en $\phi(b) = \infty$. Nu is $\phi(x)$ bekend, zodra α en β gegeven zijn. Ieder ander stel α en β bepaalt een andere $\phi(x)$ en het is de vraag of ook nu tussen al deze $\phi(x)$'s en een eenvoudige relatie bestaat, zo simpel als bij het normaliseren.

1) Misschien lijkt het logisch te eisen, dat de gemiddelde waarden van x en z (μ_x en μ_z) gelijk zullen zijn, evenals hun varianties (σ_x^2 en σ_z^2). Men kan immers μ_z en σ_z kiezen zoals men wil; d.w.z. $\mu_z = \mu \equiv \int_a^b x f(\underline{x})d\underline{x}$ en $\sigma_z^2 = \sigma_x^2 \equiv \int_a^b (x - \mu)^2 f(\underline{x})d\underline{x}$. Dit is een kwestie van smaak. Men kan ook prefereren $\mu_z = 0$ en $\sigma_x = 1$ of $\mu_z = \text{modus}$ of mediaan in $f(\underline{x})d\underline{x}$; enz. enz.

Met dit te zeggen is natuurlijk $z = \psi(x)$ als functie van x expliciet nog niet bekend. Het kan zijn, dat wij juist deze functie willen kennen.

2. Hoe vinden wij z als functie van x , exact of in voldoende goede benadering?

Men bedenke:

$$3. \int_a^x f(\underline{x}) d\underline{x} = \int_{-\infty}^{\psi(x)} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(z-\mu)^2}{2\sigma^2}} dz = \int_{-\infty}^{\tau(x)} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\tau^2} d\underline{\tau}, \text{ waarin}$$

$$\phi(\underline{\tau}) d\underline{\tau} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\tau^2} d\underline{\tau} = \mathcal{N}(0;1) = \text{de gestandaardiseerde}$$

normale verdeling op de variabele τ (gemiddelde = 0; standaarddeviatie: 1) voorstelt. Hierbij is $\tau = (z - \mu) : \sigma$ (de indices z aan μ_z en σ_z laten we maar weg.)

Men berekent voor elk van een groot aantal ($=N$) waarden van x : $a \leq x_1 < x_2 \dots < x_n \leq b$, de onderschrijdingskans $F(x \leq x_1) = \int_a^{x_1} f(\underline{x}) d\underline{x}$ en zet deze op lineair-gauszisch papier uit. Dit papier heeft twee assen: een lineaire x -as en een niet lineaire kansen-as G . Langs de laatste kan men zich ook een lineaire τ -as voorstellen, met $\tau=0$ in het punt $G = 50\%$, zodat bijv. $\tau = \pm 1$ overeenkomst met de kansen 0.841 resp. 0.159 en $\tau = \pm 2$ met de kansen 0.977 resp. 0,023 etc. (zie de tabel van de gestandaardiseerde normale verdeling, volgens welke $\int_{-\infty}^1 \phi(\tau) d\tau = 0.841$ en $\int_{-\infty}^{-1} \phi(\tau) d\tau = 1 - 0.841 = 0.159$ enz.)

Vindt men bijv. bij $x=x_3$ een $F_3=0.0107$, dan komt daarmee overeen het punt $G_3=0.0107$; $x=x_3$, doch ook het punt $\tau_3=-2.3$; $x=x_3$, aangezien $\int_{-\infty}^{-2.3} \phi(\tau) d\tau = 0.0107$. Geeft een andere waarde $x=x_8$ een $F_8=0.9332$, dan wordt hiermede het punt $G_8=0.9332$; $x=x_8$ vastgelegd, d.i. tevens $\tau_8=1.5$; $x=x_8$, omdat $\int_{-\infty}^{1.5} \phi(\tau) d\tau = 0.9332$. Aldus ontstaan er N punten $\tau_i; x_i$. Het verband tussen τ en x zal kunnen worden beschreven door $\tau = \psi(x)$, zo nauwkeurig als men wil, aangezien men N onbeperkt kan laten toenemen. (het "hoe" bespreken wij niet; dit is een afzonderlijk probleem). Men kan natuurlijk met zekere graad van benadering tevreden zijn teneinde $\psi(x)$ niet te ingewikkeld te maken.

Voor zekere \tilde{x} (de z.g. mediaan) zal $F(x \leq \tilde{x}) = \frac{1}{2}$ zijn. Gesteld dat men deze \tilde{x} voldoende nauwkeurig berekenen kan (de berekening is soms zeer moeilijk), dan moet in elk geval $\psi(x)$ een dusdanige vorm hebben, dat $\psi(\tilde{x})=0$, of ook, de normalisatie-transformatie $z=z(x)$ geeft voor $x=\tilde{x}$ het gemiddelde μ in de normale z -verdeling, d.w.z. $\mu=z(\tilde{x})$. D.w.z. de kromme $\tau = \psi(x)$ in het dubbellineaire τ - x -vlak moet door het punt $\tau=0$, $x=\tilde{x}$ gaan.

Men moet bedacht zijn, in het algemeen, op de ongelijkheid der schaal-eenheden langs de τ - en x -as op het gauszisch papier. (e_τ resp. e_x). Op het gebruikelijke waarschijnlijkheidspapier is de τ -eenheid ca. 31.5 mm papier, terwijl men voor de x -eenheid nemen kan, naar believen, 42, 21, 10½, 2, 1 enz. mm. Als men gelijke eenheden wenst, moet men schrijven $\tau^{\bar{x}} = q \cdot \psi(x)$ als $q = e_\tau : e_x$. Deze $\psi(x)$ bepaalt de gezochte normalisatie-transformatie, echter niet éénduidig. Dit blijkt a.v.: zodra men $\psi(x)$ schrijven kan in de vorm $\psi(x) = \frac{\psi(x) - \mu}{\sigma}$ is deze $z = \psi(x)$ de gezochte normalisatie-transformatie, waarbij de normaal verdeelde z het gemiddelde μ en de standaarddeviatie σ heeft. Maar dan is er ook een andere $z^{\bar{x}} = \psi_1(x)$, die normaal verdeeld is rondom een ander gemiddelde μ_1 met een andere σ_1 , die zich a.v. laten bepalen. Het is altijd mogelijk te schrijven $\mu = A \mu_1 + B$ en $\sigma = A \sigma_1$.

Nodig is $F(x) \int_a^x f(\underline{x}) d\underline{x} = \int_{-\infty}^{z(x)} (\sigma \sqrt{2\pi})^{-1} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\} dz = \int_{-\infty}^{z^{\bar{x}}(\underline{x})} (\sigma_1 \sqrt{2\pi})^{-1}$

$\exp\left\{-\frac{(z-\mu_1)^2}{2\sigma_1^2}\right\} dz$. Maar

$$\int_{-\infty}^{z(x)} (\sigma \sqrt{2\pi})^{-1} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\} dz = \int_{-\infty}^{z(x)} (A \sigma_1 \sqrt{2\pi})^{-1} \exp\left[-\frac{\{z-(A\mu_1+B)\}^2}{2 A^2 \sigma_1^2}\right] dz =$$

$$\int_{-\infty}^t (\sigma_1 \sqrt{2\pi})^{-1} \exp\left\{-\frac{(t-\mu_1)^2}{2\sigma_1^2}\right\} dt, \text{ als } t = \frac{z-B}{A}.$$

De nieuwe bovengrens ? wordt

dan $\{z(x)-B\} : A$. Opdat dit waar zal zijn voor elke x uit $a \leq x \leq b$ moet gelden

$$4. \boxed{z^{\bar{x}}(x) \equiv \left\{z(x)-B\right\} / A = \frac{\sigma_1}{\sigma} \cdot z(x) + \mu_1 - \frac{\sigma_1}{\sigma} \mu.}$$

Zodra dus één normaliserende transformatie $z = \psi(x)$ bekend is, bij een bepaald paar μ en σ , is een ∞^2 verzameling van andere normalisatie transformaties bekend, want geeft men μ_1 en σ_1 , dan is $z^{\bar{x}}$ volgens (4) in z uit te drukken.

Getallenvoorbeelden

a) $f(\underline{x}) d\underline{x} = \frac{1}{(\underline{x} + 25) \sqrt{2\pi}} \exp\left[-\frac{\{\lg(\underline{x} + 25)\}^2}{2}\right] d\underline{x}$, voor $-25 \leq x \leq \infty$.

Nu is zeker $z = \lg(x+25)$ een normalisatie-transformatie, waarbij z normaal verdeeld is met gemiddelde 0 en standaarddeviatie 1. ($-\infty \leq z \leq \infty$, als $-25 \leq x < \infty$). Wil men een normale verdeling met gemiddeld 4 en standaarddeviatie 5 dan is (volgens (4)) de nieuwe norm. transformatie $z^{\bar{x}}(x) = 5 z(x) + 4$, d.i. $z^{\bar{x}}(x) = 5 \lg(x+25) + 4$ ($-\infty \leq z^{\bar{x}} \leq \infty$ als $-25 \leq x < \infty$) ; kennelijk is de mediaan $\tilde{x} = -24$, waarvoor $z = \lg(x+25) = 0$.

$$b) \int_{\underline{x}} f(\underline{x}) d\underline{x} = \frac{e^{\frac{x}{4}} + \frac{1}{x^2}}{\sqrt{2\pi}} e^{-\left(e^{\frac{x}{4}} - \frac{1}{x} - 4\right)^2 / 2} d\underline{x}, \text{ voor } 0 \leq \underline{x} \leq \infty.$$

Nu is zeker $z = e^{\frac{x}{4}} - \frac{1}{x}$ één der ∞^2 normalisatie-transformaties, waarvoor $-\infty \leq z \leq \infty$, als $0 \leq x \leq \infty$; de normale z -verdeling heeft een gemiddelde 4 en een standaarddeviatie 1. Wil men liever een gemiddelde 0 en standaarddeviatie 10, dan is volgens (4) de nieuwe norm. transformatie $z^*(x) = 10\left(e^{\frac{x}{4}} - \frac{1}{x}\right) - 40$. De mediaan $\hat{x} \cong 1.54$, want hiervoor is $e^{\frac{x}{4}} - \frac{1}{x} = 4$

Deze twee voorbeelden lijken sterk van academische aard, omdat zij voldoende doorzichtig zijn. Het kan echter voorkomen, dat $f(\underline{x})$ zelf zo ingewikkeld is, dat men de behoefte gevoelt een betrekkelijk eenvoudige normalisatie-transformatie $z = \psi(x)$ te vinden, zodanig, dat z voldoende goed normaal verdeeld is, terwijl \underline{x} zelf het beslist niet is. Is eenmaal $z = \psi(x)$ gevonden, dan kan het zijn, dat het minder rekenwerk vergt om voortaan voor iedere andere x - waarde van de gebruikte N stuks niet direct $\int_a^x f(\underline{x}) d\underline{x}$ te berekenen, maar d.w.z. de z , d.w.z. als $\int_{-\infty}^{z(x)} \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(z - \mu)^2}{2\sigma^2}\right\} dz = \int_{-\infty}^{\tau(x)} \phi(\tau) d\tau$ met $z(x) = \psi(x)$, en $\tau(x) = \{z(x) - \mu\} / \sigma$, waarbij $z = \psi(x)$ en daarmee $\tau(x)$ berekend moet worden en verder de tabel van de gestandaardiseerde normale verdeling geraadpleegd moet worden.

Bovendien: zoals reeds gezegd werd, is het tweede aspect der normalisatie veel belangrijker dan het eerste, aangezien wij nooit met universa, doch met steekproeven te maken hebben en vrijwel nooit de $\int f(x) dx$ exact bekend is. Het kwam ons systematischer voor toch met de behandeling van dit eerste aspect te beginnen.

3. Tweede geval; $\int f(\underline{x}) d\underline{x}$ moet worden geschat uit de steekproef.

Met dit geval hebben we gewoonlijk te maken. Immers van $\int f(\underline{x}) d\underline{x}$ kennen wij veelal niet eens het type, en zo ja, dan toch niet de numerieke waarden der bepalende parameters. Er zijn N (voorlopig verschillend gedachte) metingen $x_1, x_2 \dots x_n$ verricht, welke waarden tot een populatie met een onbekende waarschijnlijkheidsverdeling $f(\underline{x}) d\underline{x}$ behoren. Gewoonlijk kent men wel a priori de kleinst en de grootst mogelijke x - waarde (onder geval 1 resp. a en b genoemd) Dagsommen neerslag : $a = 0, b = \infty$; zonnenschijnpercentages : $a = 0, b = 100$; maximum-temperatuur : $a = -273^\circ C, b = \infty^\circ C$; aantal ijsdagen in december: $a = 0, b = 31$, enz. Wij kennen derhalve ook niet de ware waarde van $F' = \int_a^{x_i} f(\underline{x}) d\underline{x}$, voor $x_i = x_2, \dots, x_n$. Wij kennen slechts schattingen \hat{F} . Over de beste schattingen verschillen de meningen.

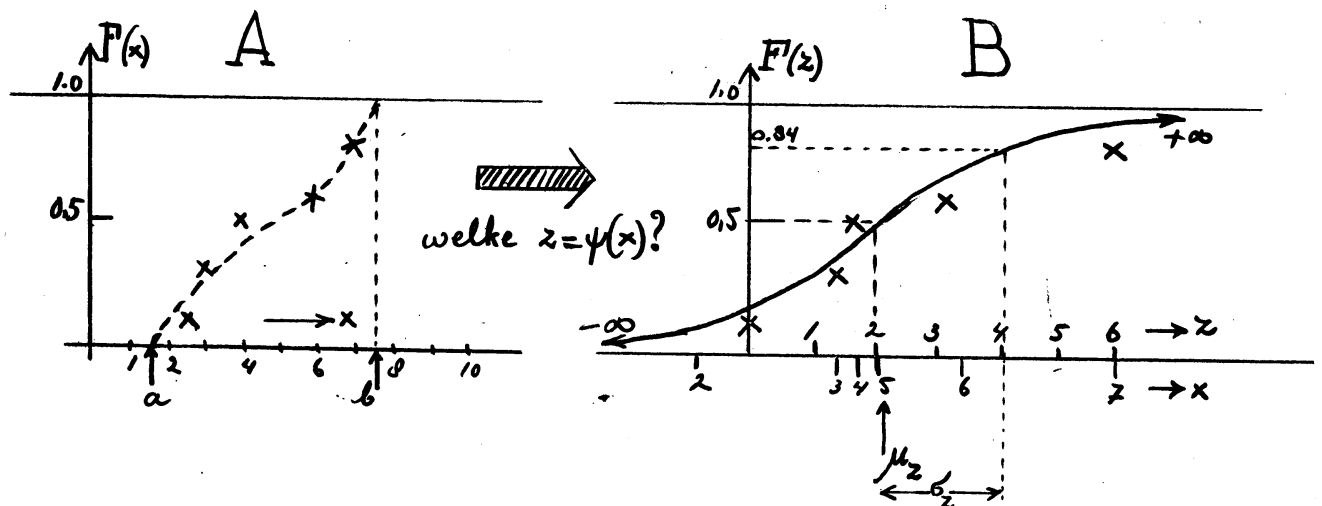
Er is veel over geschreven; dikwijls ligt het knellende punt in de definitie van "beste". Laat de x -waarden reeds gerangschikt zijn: $x_1 \leq x_2 \leq \dots \leq x_n$. (waaronder eventueel gelijke waarden). Men neemt wel: $\hat{F}(x_1) = \frac{1}{N}$, waarbij met het eventuele meervoudig optreden van zekere x -waarden rekening gehouden is. Komt bijv. x_1 n_1 keren voor; x_2 n_2 keren; ... x_n n_n keren, zodat $\sum_{i=1}^n n_i = N$, dan is $\hat{F}(x_i) = \frac{\sum_{j=1}^i n_j}{N}$. Anderen prefereren $\frac{i - \frac{1}{2}}{N}$ (d.i. $\frac{2i - 1}{2N}$); $\frac{i - 1}{N}$; $\frac{i + 1}{N}$ enz.

In [2] wordt uiteengezet, dat, als men te voren weet, dat $f(x)dx$ van het type der normale verdeling is, het best $\hat{F}(x_1) = \frac{1 - 0.3}{N + 0.4}$ genomen kan worden. Verondersteld, men heeft op zekere gronden redenen om de normaliteit van de x -verdeling te poneren. Nadat dan de n punten in het lineair-gauszisch net aangebracht zijn, kan men, bijv. op het oog, de beste rechte tussen deze punten door trekken, waardoor weer de beste μ en σ vastgelegd worden. Doch dit geval interesseert ons hier niet.

Wij onderstellen nl., dat we niets over de gedaante van $f(x)$ weten (althans niet meer dan wat de steekproef daarover suggereert) en dat wij willen onderzoeken of $f(x)$ kan behoren tot die verzameling van functies, die zich met simpele vormen van transformaties $z = \psi(x)$ laten normaliseren.

Ten behoeve van dit onderzoek kan het nuttig zijn criteria (toetsen) te ontwikkelen. Zulk een criterium en wel t.a.v. $z = \psi(x) = x^a$ werd door Neumann en Kotz in [6] ontwikkeld.

De volgende schetsjes illustreren een en ander. Thans is in het x - F -vlak niet een ononderbroken kromme gegeven, maar een (klein) aantal punten (bij N metingen: N of $N-1$ stuks, al naar de opvatting omtrent de beste $\hat{F}(x)$). Tussen die punten door loopt de bij het universum behorende, onbekende, onderschrijdingscurve, tussen de eindpunten a en b van het variatiegebied $a \leq x \leq b$. Kan deze curve zulk een gedaante hebben, dat er een simpele transformatie $z = \psi(x)$ bestaat, zodanig, dat z wel normaal verdeeld is, rondom μ en met standaarddeviatie σ ?



De verdere procedure is als onder geval 1, doch nu liggen de punten τ , x in het τ - x - veld "meer verspreid", d.w.z. ze liggen niet precies op, doch door toevalseffecten (steekproefeffecten) rondom de onbekende kromme $\tau = \varphi(x)$, hetgeen het vinden van deze $\varphi(x)$ bemoeilijkt. Nadat dan een bevredigende $\varphi(x)$ gevonden is, hebben wij tegelijk de beste schatting van $f(\underline{x})d\underline{x}$ gevonden, nl. $\hat{f}(\underline{x})d\underline{x} = \frac{dz/dx}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\} d\underline{x}$, waarin $z = \varphi(x)$ normaal verdeeld is, rondom $Ez = \mu$ en met variantie σ^2 , waarna met de χ^2 -toets kan worden onderzocht hoe goed empirische en theoretische distributie frequentieverdeling "overeenkomen". Hierbij is $\varphi(x)$ geschreven gedacht als $(\varphi(x) - \mu) : \sigma$

Verschillende auteurs gaan van de andere kant uit, d.w.z. zij denken $z = \varphi(x)$ gegeven en vragen zich dan af welke vorm $f(\underline{x})d\underline{x}$ bezit en welke klimatologische grootheden een distributie hebben die goed met deze typen $f(\underline{x})d\underline{x}$ te benaderen zijn. Men neemt allicht eenvoudige functies $\varphi(x)$, zoals $a+bx$; $(a+bx)^c$; $a_0 + a_1x + a_2x^2 + \dots$; $d + \lg(x+c)$ enz.

Helaas ontwikkelen ze geen criterium, door middel waarvan men met behulp der N metingen kan uitmaken (met voorgeschreven betrouwbaarheid) of de ene normalisatie-transformatie wel, de andere niet in aanmerking komt.

4 Numerieke voorbeelden.

4.1 Uit de praktijk.

Wij citeren een voorbeeld uit het rapport van Rijkooft [7]. Hij beschouwde in elk van 35 winters de grootste diepte x (beneden aardniveau), waartoe de vorst doordrong. Er waren winters, waarin de vorst niet in de grond drong. Voor deze werd x nul genoemd. De 35 x -waarden waren 8 keer 0 ($x_1 = \dots x_8 = 0$); $x_9 = 3$; $x_{10} = x_{11} = 4$; $x_{12} = 5$; $x_{13} = x_{14} = 6$; $x_{15} = 7$; $x_{16} = 10$; $x_{17} = 1$; $x_{18} = 12$; $x_{19} = x_{20} = 13$; $x_{21} = 15$; $x_{22} = 16$; $x_{23} = 17$; $x_{24} = 19$; $x_{25} = 21$; $x_{26} = 22$; $x_{27} = x_{28} = 31$; $x_{29} = 34$; $x_{30} = 35$; $x_{31} = 36$; $x_{32} = 38$; $x_{33} = 53$; $x_{34} = 54$; $x_{35} = 60$ (24 ongelijke x - waarden)

Als men log-gausz-papier gebruikt en met de x ' en zelf werkt, komen de punten allesbehalve lineair te liggen, zie in [7]. Het ligt dus voor de hand het met $x + b$ (b te zoeken) te proberen, wat impliceert, dat men wil onderstellen, dat $x + b$ normaal verdeeld is (van welke normale verdeling men verder de μ en σ weten wil). Aldus ging Rijkooft te werk. Naar wat waarden geprobeerd te hebben kwam hij tot $20 < b < 30$, terwijl μ en σ van het papier afgelezen kunnen worden. Deze uitkomst houdt in, dat men voor de onbekende kansverdeling van x meent gevonden te hebben in voldoende goede benadering?

$$f(\underline{x})d\underline{x} = \frac{1}{\sigma(\underline{x}+b)\sqrt{2\pi}} \exp. \left[-\lg(\underline{x}+b) - \mu^2/2\sigma^2 \right] d\underline{x}$$

Toch kan dit niet exact juist zijn, want alleen over het traject $-b \leq x < \infty$ is $\int_{-b}^{\infty} f(\underline{x})d\underline{x} = 1$, terwijl \underline{x} alleen waarden $0 \leq x < \infty$ aannemen kan. Men kan aan de moeilijkheid ontkomen door uit te gaan van $\int_{-b}^{\infty} f^{\#}(\underline{x})d\underline{x} = \left\{ \int_{-b}^{\infty} f(\underline{x})d\underline{x} \right\} / K$, als $K = \int_0^{\infty} f(\underline{x})d\underline{x} < 1$, want nu is $\int_0^{\infty} f^{\#}(\underline{x})d\underline{x} = 1$. Vermoedelijk ligt K zeer dicht bij één. Doet men dit, dan zal de aanpassing allicht iets minder goed zijn en zullen wij op iets andere waarden van b , μ en σ moeten overgaan om haar te herstellen ¹⁾.

Deze complicaties kunnen ook vermeden worden door allereerst alleen die winters (derhalve minder dan 35 stuks) te beschouwen, waarin x of precies nul of positief is en de andere winters, waarin zelfs niet het aardniveau op 0°C-temperatuur kwam, uit te zonderen. Alleen voor de dan overblijvende x -waarden (een aantal ≤ 35) zoeken we de kansverdeling in het universum. Bijgevolg moeten wij dan eigenlijk de logaritmisch-normale verdeling met $b > 0$ uitsluiten en direct naar andere uitzien ²⁾. Het is niet de bedoeling hier verder op door te gaan; wel willen wij, louter bij wijze van numeriek voorbeeld, het probleempje van Rijkoort hier uitwerken:

Er geldt $\int_{-b}^{\infty} f(\underline{x})d\underline{x} = \int_{-\infty}^{\infty} \psi(z)dz = 1$; zet $t = (z-\mu)/\sigma$ en $\phi(t) = \frac{e^{-\frac{1}{2}t^2}}{\sqrt{2\pi}}$ dt

Er zijn 24 verschillende x -waarden x_i . Voor elk geldt:

$$\int_{-b}^{x_i} f(\underline{x})d\underline{x} = \int_{-\infty}^{z_i} \psi(z)dz = \int_{-\infty}^{t_i} \phi(t)dt \quad \text{met } z_i = \lg(x_i+a) \text{ en } t_i = (z_i-\mu) / \sigma$$

Door dit voorschrift is aan elke x_i één t_i toegevoegd. We moeten de beste a , μ , σ berekenen; de beste schatting van $\int_{-b}^{x_i} f(\underline{x})d\underline{x}$ (onbekend, aangezien $f(\underline{x})d\underline{x}$ onbekend is) is

$\hat{F}(\underline{x} \leq 0) = \frac{8}{35} = 0.228$, zodat $t = -0.75$;
 die van $\int_{-b}^3 f(\underline{x})d\underline{x}$ is $\hat{F}(\underline{x} \leq 3) = \frac{8+1}{35} = 0.257$, zodat $t = -0.65$;
 die van $\int_{-b}^4 f(\underline{x})d\underline{x}$ is $\hat{F}(\underline{x} \leq 4) = \frac{8+1+2}{35} = 0.314$, zodat $t = -0.485$ enz.

Aldus ontstaan er 23 vergelijkingen van de vorm $t = \frac{\lg(x+b)-\mu}{\sigma}$ of $\lg(x+b) = t\sigma + \mu$, met onbekenden b , σ en μ . Hoe deze oplossen? Neem enkel-log-papier. Kies een b , bijv. 10 en zet de punten bij $\lg(10+x)$ tegen x uit (fig. 1). Zie of ze lineair liggen.

1) Maar deze finesse is van geen betekenis als het ons alleen gaat om een schatting van de overschrijdingskans $\int_m^{\infty} f(\underline{x})d\underline{x}$ voor grote m en hier ging het Rijkoort inderdaad in hoofdzaak om.

2) Zie mijn ander rapport [5], waarin ik normaliseer d.m.v. $z = x^a$.

Zo niet, neem een wat groter b , bijv. 20. Misschien liggen de punten nu op een wat minder sterk gebogen curve. Vergroot b weer wat en ga zo voort tot de punten zo goed mogelijk lineair liggen. Dit op het oog te beoordelen (want hoe objectief?) Wij vonden $\hat{b}=25$; $\hat{\mu}=3.6$ en $\hat{\sigma}=0.5$ cm. De normalisatie transformatie is dus $z=\lg(x+25)$ en de kansverdeling

$$f(\underline{x})d\underline{x} = \frac{e^{-\frac{[\lg(\underline{x}+25)-3.6]^2/2(\frac{1}{2})^2}}{\frac{1}{2}(\underline{x}+25)\sqrt{2\pi}} d\underline{x}, \text{ waarvoor } \int_{-25}^{\infty} f(\underline{x})d\underline{x}=1$$

Hoe goed is de benadering?

Zie het volgende tabelletje.

interval voor x	aantallen	
	geteld	met formule
≤ 10 cm	16	16.2
11-20	8	7.2
21-30	2	4.6
31-40	6	2.9
41-50	0	1.6
51-60	3	1.4
≥ 61	0	1.1
Som	35	35.0

De χ^2 -toets is helaas niet toe te passen omdat (afgezien van de reden, dat hier de beste schattingen grafisch gebeurden en niet volgens de Methode der Maximum Likelihood) het aantal graden van vrijheid tenminste 1 moet zijn; d.w.z. het aantal intervallen moet tenminste $1+4=5$ zijn (3 parameters en sommen gelijk), terwijl in elk interval het gemiddeld te verwachten aantal liefst tenminste 5 moet zijn. Voor dit alles is het totale aantal 35 te klein. Maar wij zien ook wel op het oog hoe (goed) de aanpassing is.

N.B. In wezen is de hier behandelde "methode" gelijk aan die van Rijkoort. In beide gevallen ontbreekt het "lineariteitskriterium"

4.2 Gefantaseerd voorbeeld.

Er werden 10 metingen verricht; met behulp van $\hat{P}_i = 1/10$ ($i=1.2 \dots 9$) werden op lineair-gauszisch papier de 9 kruisjes verkregen. (fig. 2) Na wat proberen lijkt de kromme $\tau = e^x - \frac{1}{x} - 4$ (tussen $-\infty$ en $+\infty$) een goede benadering van de ware cumulatieve kanscurve.

Is men met een aanpassing op het oog niet tevreden en wenst men het begrip "voldoende goede benadering" objectief mathematisch te definiëren, dan is dit een allesbehalve eenvoudig vraagstuk, ondermeer vanwege het feit, dat de punten min of meer afhankelijk zijn, d.w.z. met toenemende x een groter gewicht hebben. Uit genoemde $\tau(x)$ volgt weer, dat de beste schatting van de onbekende kansverdeling in het x -universum is

$$\hat{f}(x) dx = \frac{e^x + \frac{1}{x^2}}{\sqrt{2\pi}} \exp. \left[-\frac{(e^x - \frac{1}{x} - 4)^2}{2} \right] dx, \text{ voor } 0 \leq x \leq \infty.$$

De normalisatie-transformatie is dus $z = \psi(x) = e^x - \frac{1}{x}$.

De normale z -verdeling is $\frac{1}{\sqrt{2\pi}} \exp. \left\{ -\frac{(z-4)^2}{2} \right\} dz$, met $\mu_z = 4$ en $\sigma_z = 1$.

5. Literatuur.

- [1] Levert, C. Beschouwing over betrouwbaarheidsintervallen en een uitschieterscriterium K.N.M.I. W.R. 56-005, 1956.
- [2] Benard, A. en Bes-Levenbach, E.C. Het uitzetten van waarnemingen op waarschijnlijkheidspapier. Statistica 7 163 1953.
- [3] Iterson, G. van Het nut van waarschijnlijkheidspapier voor de variantiestatistiek. Statistica 4 129 1950.
- [4] Exalto, L.J.H. Gebruik en toepassing van waarschijnlijkheidspapier. Statistica 4 121 1950.
- [5] Levert, C. Normaliseren door middel van de transformatie $z=x^a$. K.N.M.I. W.R. 57-002, 1957.
- [6] Neumann, J. and Kotz, S. Some Pearsonian-like frequency functions capable of a modified normalization. In "The scientific basis of weather modification studies", Arizona 1956.
- [7] Rijkooert, P. Bijdrage tot het bepalen van de meest gunstige diepte voor het leggen van waterleidingbuizen in verband met het bevroeringsrisico. K.N.M.I. W.R. 56-003, 1957.

6. Summary

0. In this chapter 8 examples for the necessity of normalization are mentioned:
- 0.1 Both \underline{x} and \underline{y} must be distributed normally if one wants to verify the statistical reality of the common correlation coefficient by means of the common correlation table.
- 0.2 The members of a time series $x_1, x_2, x_3 \dots$ must all be distributed normally if one wants to investigate the significance of the serial correlation coefficients (autocorrelation) of lag 1, 2 and so on, as in 0.1.
- 0.3 Having at our disposal only N measurements $x_1 \leq x_2, \dots \leq x_n$ it is possible to compute the confidence limits for the exceedance probability of a given value x_g (even if $x_g < x_1$ or $> x_n$) almost only if the \underline{x} -distribution is normal.
- 0.4 Application of the t -test needs a normal distribution in both populations.
- 0.5 Application of the Variance Analysis needs a normal distribution.
- 0.6 Application of the Bartlett-test (is $\sigma_1 = \sigma_2 = \dots = \sigma_n$?) supposes that the n universa all have normal distributions.
- 0.7 Applications of the z - or F - test needs the same fact.
- 0.8 Several criteria for rejection of "outliers" require that the distribution in the population is normal [1].
- 0.9 In special theories the form of the transformation $z = \psi(x)$ is seen in relation to special theoretical considerations.

1.1. $f(\underline{x})d\underline{x}$ is known exactly.

Either $f(\underline{x})d\underline{x}$ or $F(\underline{x} \leq x) = \int_a^x f(\underline{x})d\underline{x}$ is given analytically, or $F(\underline{x} \leq x)$ is given as a curve on double linear paper (\underline{x} and F in linear scales) or e.g. on linear-probability paper (\underline{x} in linear scale and F in gausz-scale). The variation range of x should be $a \leq \underline{x} \leq b$.

Definition: The variable \underline{y} is distributed normally if at least one couple of values η and σ exists so that for each value y_1 out of the range $-\infty < y < \infty$ the probability of a value situated between y_1 and $y_1 + dy_1$ can be written as

$$(1) \quad g(y_1)dy_1 = (\sigma\sqrt{2\pi})^{-1} \exp\left[-(y_1 - \eta)^2 / 2\sigma^2\right] dy_1$$

Definition of normalizing: suppose \underline{x} itself is not distributed normally. We ask: Is there at least one normalizing transformation $z = \psi(x)$ so that \underline{z} is distributed normally? This involves:

$$(2) \quad \int_a^{x_1} f(\underline{x})d\underline{x} = \int_{-\infty}^{z_1} g(\underline{z})d\underline{z} \quad \text{with } z_1 = \psi(x_1)$$

As soon as the mean and standard deviation μ_z and σ_z are chosen the function $z = \psi(x)$ is fixed uniquely by the relation (2) implicitly.

1.2 $f(\underline{x})d\underline{x}$ is not known, but N values $x_1, x_2 \dots x_n$ have been "measured". Suppose the \underline{x} population is not distributed normally. Stress is laid on the following question: is it possible that the unknown $f(\underline{x})$ has such a (simple) form that it can be normalized by a simple normalization transformation $z = \psi(x)$ (from a given class)? Of course it is desirable and useful to derive a criterion with the aid of which it is possible to investigate wholly statistically whether the first question can be answered in the affirmative. Such a criterion will be based on special functions of the N measurements and relations between these functions.

2. Suppose $f(\underline{x})d\underline{x}$ is known, defined for $a \leq x \leq b$, with $\int_a^b f(\underline{x})d\underline{x} = 1$. Then $z = \psi(x)$ must be found (and is defined) so that

$$(3) \int_a^{x_1} f(\underline{x})d\underline{x} = \int_{\psi(a)}^{\psi(x_1)} \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\} dz \quad \text{for each } a \leq x \leq b \text{ and given values of } \mu \text{ and } \sigma.$$

(4) Hence $\psi(a) = -\infty$; $\psi(b) = \infty$ and $f(x) = \frac{dz/dx}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(z-\mu)^2}{2\sigma^2}}$ Further:

$$(5) \int_a^{x_1} f(\underline{x})d\underline{x} = \int_{-\infty}^{\psi(x_1)} \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(z-\mu)^2}{2\sigma^2}} dz = \int_{-\infty}^{\psi(x_1)} \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}\tau^2} d\tau$$

(6) where $\phi(\tau)d\tau = \frac{e^{-\frac{1}{2}\tau^2}}{\sqrt{2\pi}} d\tau$

is the normal probability distribution in standard form ($\mu=0$; $\sigma=1$). The solution is as follows: Choose a large number N of values $x_1, x_2 \dots x_N$, compute the values $F(x \leq x_1)$ and plot these pairs x, F on linear-probability paper. The non linear F -scale of this paper is at the same time the linear τ scale defined by the probability integral $F = \int_{-\infty}^{\tau} \phi(\tau)d\tau$. Suppose the relation between x and τ can be described sufficiently well by a certain function $\tau = \psi(x)$, with $\psi(\tilde{x}) = 0$ if \tilde{x} = median, defined by $F(x \leq \tilde{x}) = \frac{1}{2}$. Now write

(7) $\psi(x) = \left\{ \psi(x) - \mu \right\} : \sigma$, in which expression μ and σ are not functions of x . Then $z = \psi(x)$ is distributed normally in good approximation with mean μ and standard deviation σ . Each other function

(8) $z^* = \frac{\sigma^*}{\sigma} z + \mu^* - \frac{\sigma^*}{\sigma} \mu$ is also distributed normally with mean μ^* and standard deviation σ^* . Hence the functions $z^*(x)$ represent all normalizing transformations functions, approximative solutions of (4).

(9) Then $f(\underline{x})d\underline{x} = (\sigma \sqrt{2\pi})^{-1} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\} dz$

3. Suppose $f(\underline{x})d\underline{x}$ is not known, but we only have the disposal of N "measurements" $x_1 \leq x_2 \leq x_3 \dots \leq x_n$. Mostly we do know lower and upper limits of the variable \underline{x} , e.g. daily rainfall amounts $0 \leq x < \infty$; daily sunshine percentages $0 \leq x \leq 100$; number of freezing days per december $0 \leq x \leq 31$.

Now we must make estimations \hat{F} for the n unknown values of $F = \int_a^{x_1} f(x) dx$. Let the x values be ordered as follows $x_1 \leq x_2 \dots \leq x_n$; n_1 times x_1 ; n_2 times x_2 ; ... n_n times x_n , with $N = \sum_{j=1}^n n_j$. Some authors then take $\hat{F}(x \leq x_1) = \sum_{j=1}^1 n_j / N$. Other authors take $(\sum_{j=1}^1 n_j - \frac{1}{2}) : N$ or $(\sum_{j=1}^1 n_j - 1) : N$ or $(\sum_{j=1}^1 n_j + 1) : N$ or especially $(\sum_{j=1}^1 n_j - 0.3) : (N + 0.4)$ if the distribution may be supposed to be normal [2]. Plot then points x_i, \hat{F} on linear-probability paper (the non linear probability axis contains the linear τ scale) Try to find the best fitting function $\tau = \varphi(x)$. With this function corresponds a curve. The curve should be situated in such a way that the points scatter around this curve as narrowly as possible. Write $\varphi(x)$ as $\{\psi(x) - \mu\} : \sigma$. Hence the best estimation $\hat{f}(x) dx$ of $f(x) dx$ becomes:

$$(9) \quad \hat{f}(x) dx = \frac{dz/dx}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{(z-\mu)^2}{2\sigma^2} \right\} dx$$

with $z = \varphi(x) = \mu + \sigma \varphi(x)$

4. In this chapter two numerical examples are treated numerically. In the first one in each winter the maximum frost depth^x is considered. (de Bilt), see [7]. On the basis of the results obtained in 35 winters the following distribution is found:

$$f(x) dx = \frac{e^{-\left\{ \lg(x+25) - 3.6 \right\}^2 / 2 \left(\frac{1}{2}\right)^2}}{\frac{1}{2}(x+25) \sqrt{2\pi}} dx$$

with $\int_{-25}^{\infty} f(x) dx = 1$.

The agreement between the calculated frequency distribution (in the sample of 35 x values) and the theoretical one, based on the distribution mentioned above, is very good.

Vraag : Voor welke b liggen de 24 punten zo „lineair mogelijk“?
 Antw: $b \cong 25$

Dus $\lg(X+25) = t\sigma + \mu$

hier $\lg(X+25) = \frac{1}{2}t + 3.6$ of $t = 2\lg(X+25) - 3.6$

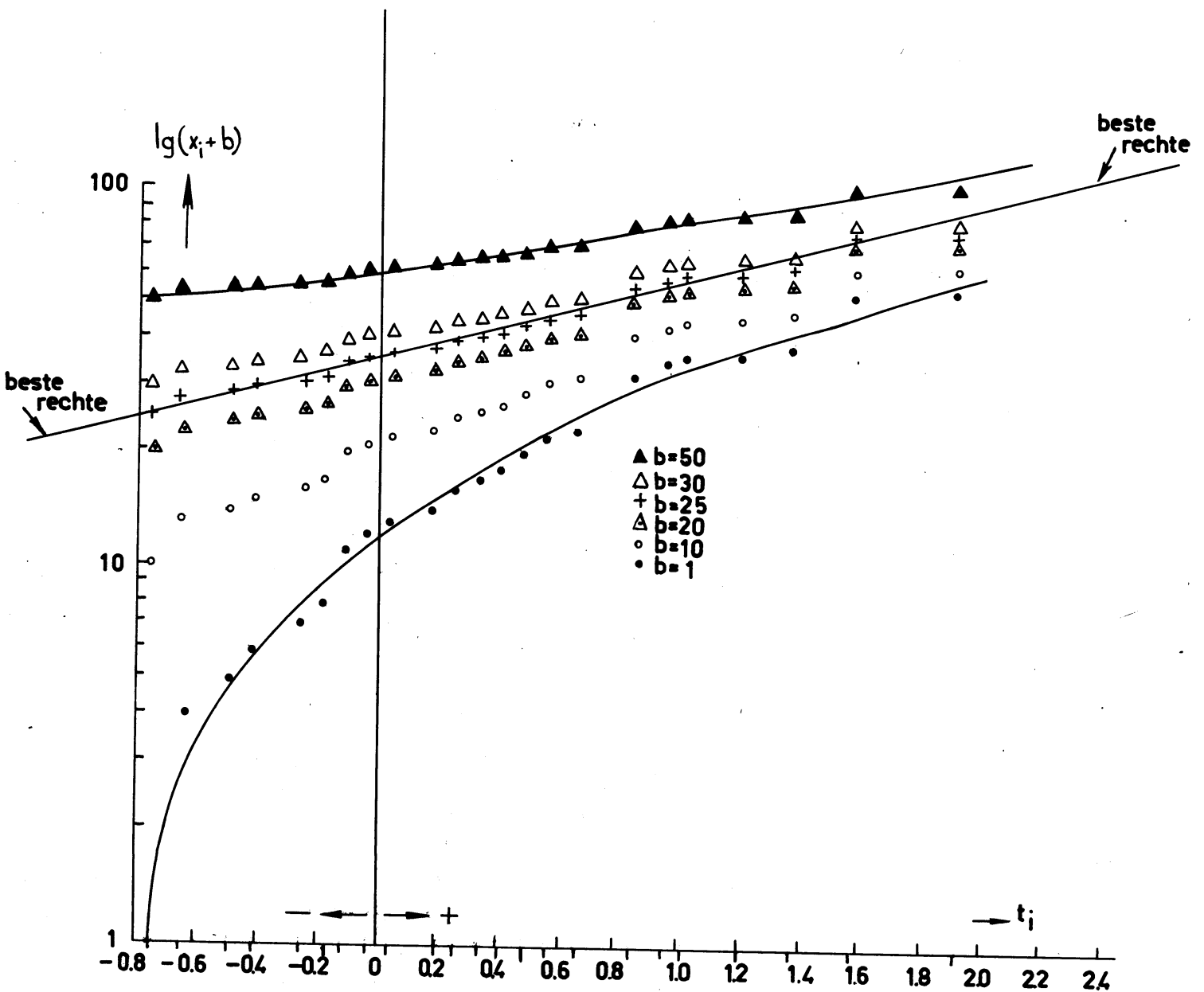
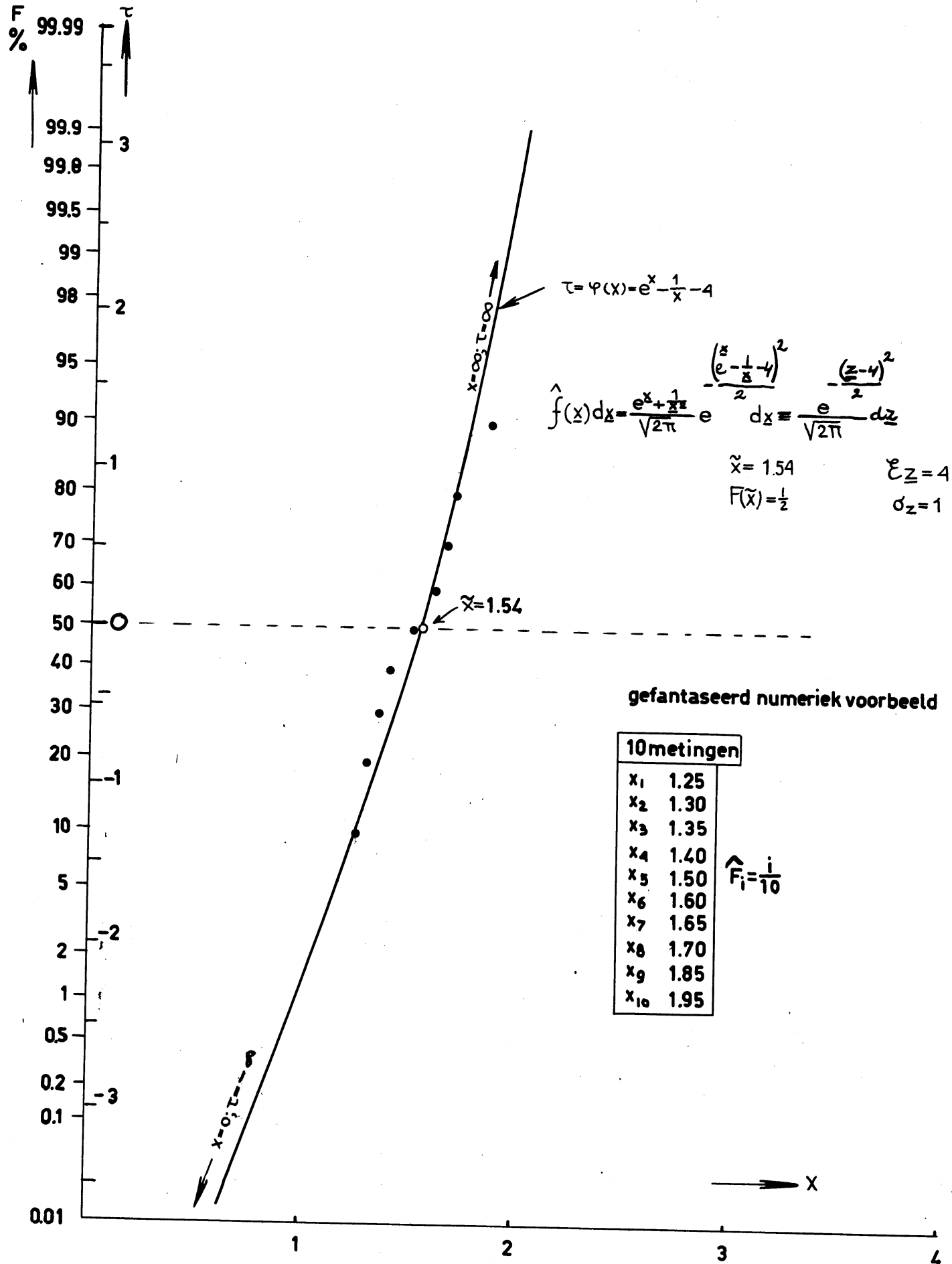


Fig 1



gefantaseerd numeriek voorbeeld

10 metingen	
x_1	1.25
x_2	1.30
x_3	1.35
x_4	1.40
x_5	1.50
x_6	1.60
x_7	1.65
x_8	1.70
x_9	1.85
x_{10}	1.95

$\hat{F}_i = \frac{i}{10}$

Fig 2