

Prof. Dr. Ivo Mossig
Institut für Geographie
Tel: 0421 / 218 67410
ivo.mossig@uni-bremen.de

Dipl.-Geogr. Lars Schieber
Institut für Geographie
Tel: 0421 / 218 67113
lars.schieber@uni-bremen.de

Deskriptive Statistik

WS2009/10

III Induktive Statistik: Schätzen und Testen

7 Anwendung bestimmter theoretischer Verteilungen

1. Beispiel für Fragestellungen der Wahrscheinlichkeitsrechnung

Von einem Fluss sind langjährig die Pegelstände gemessen worden. Aus diesen Daten lässt sich nun der mittlere Pegel (arithmetisches Mittel) und die mittlere Schwankung der Pegelstände (Standardabweichung) berechnen.

Zum Schutz vor Überschwemmungen wurde die Eindeichung des Flusses beschlossen, wobei die Höhe des Deiches diskutiert wird:

Einerseits soll der Deich so hoch sein, dass das Überschwemmungsrisiko gering und die Schutzfunktion des Deiches gewährleistet ist.

Andererseits sollen die Kosten für den Deichbau so niedrig wie möglich sein.

Gesucht wird ein Kompromiss, der ein gewisses Überschwemmungsrisiko toleriert

Angenommen man einigt sich auf eine 99%ige Sicherheit.

Das heißt, dass 99% der bisherigen Jahreshöchststände unterhalb der Deichkrone bleiben und nur ein so genanntes Jahrhundertereignis, das einmal in 100 Jahren auftritt, würde zur Überschwemmung führen.

Mit Hilfe der **Wahrscheinlichkeitsrechnung** lässt sich aus dem Daten der Messreihe berechnen, welche Höhe des Deiches diesen 99%igen Schutz gewährleistet.

2. Beispiel für Fragestellungen der Wahrscheinlichkeitsrechnung

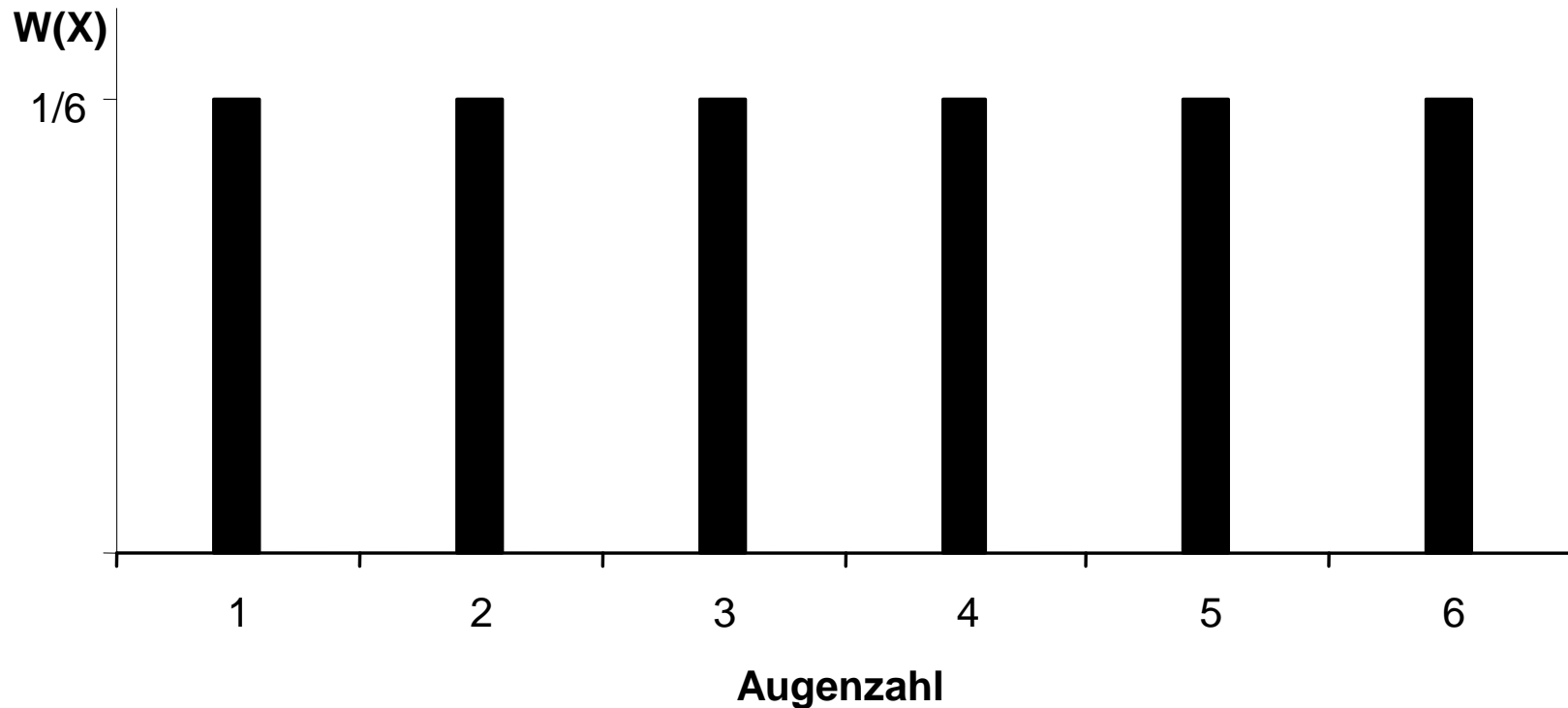
Die Oberböden in der Nordeifel werden anhand von 20 Proben auf ihren Zinkgehalt geprüft und mit dem Standardwert für Zink in Tongesteinen verglichen.

Der Mittelwert der Proben liegt unterhalb des Standardwerts für Zink in Tongesteinen.

Frage:

Ist dieser Unterschied zufällig, weil die Zahl der Proben zu gering ist und die Proben somit bezüglich ihres Zinkgehaltes nicht repräsentativ sind oder ist der Unterschied trotz der geringen Probenzahl nicht zufällig (signifikant) und der Boden in der Nordeifel kann als relativ arm an Zink eingestuft werden?

Wahrscheinlichkeiten der Ereignisse beim einmaligen Würfelwurf



$$W(X = a) = \frac{\text{Anzahl der günstigen Fälle}}{\text{Anzahl der möglichen Fälle}}$$

Es gilt:

- Die Wahrscheinlichkeit ist immer eine Zahl zwischen 0 und 1:
 $0 \leq W \leq 1$.

- Die Summe der Einzelwahrscheinlichkeiten ist immer 1

Beispiel einmaliges Würfeln:

$$W(X=1) + W(X=2) + \dots + W(X=6) = 1/6 + 1/6 + \dots + 1/6 = 1$$

- Ein sicheres Ereignis hat die Wahrscheinlichkeit 1, also

$$W(X=a) = 1$$

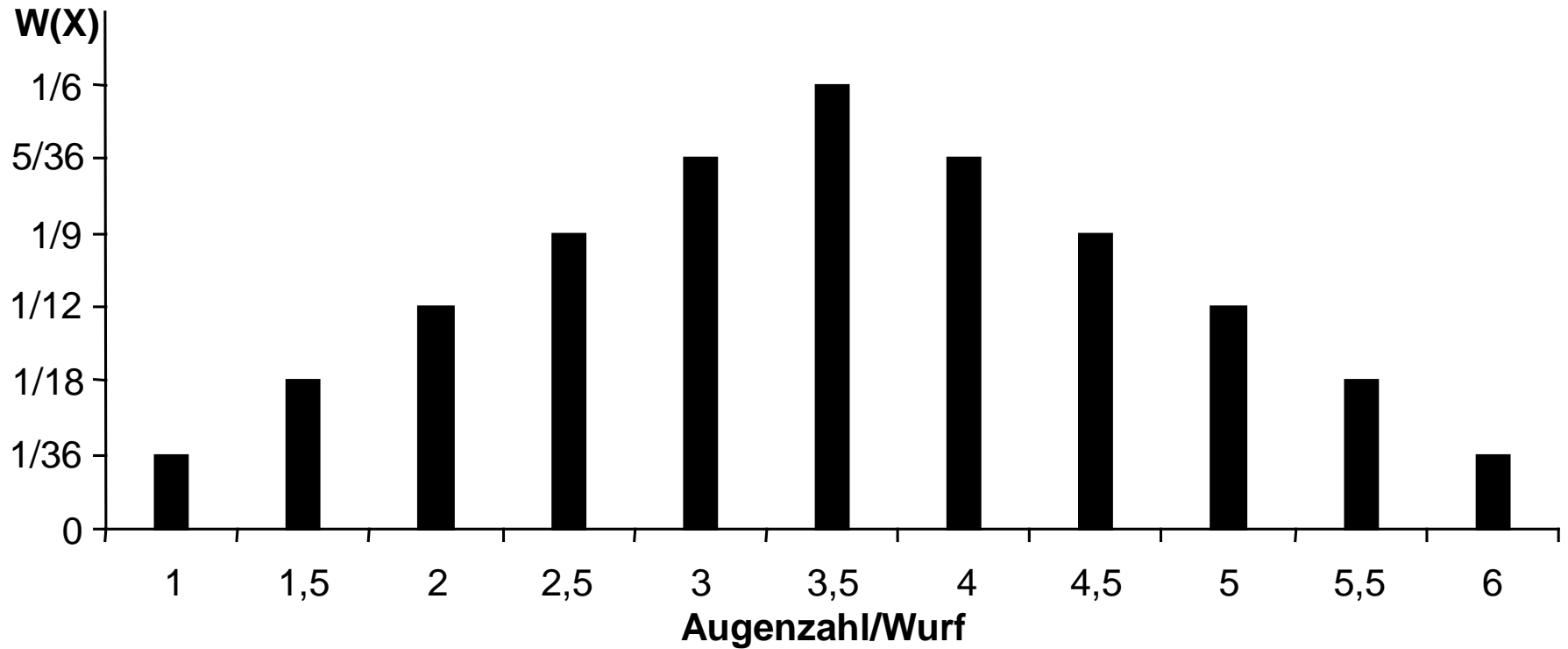
Beispiel einmaliges Würfeln: $W(X=1,2,3,4,5,6) = 1$

- Ein unmögliches Ereignis hat die Wahrscheinlichkeit 0, also

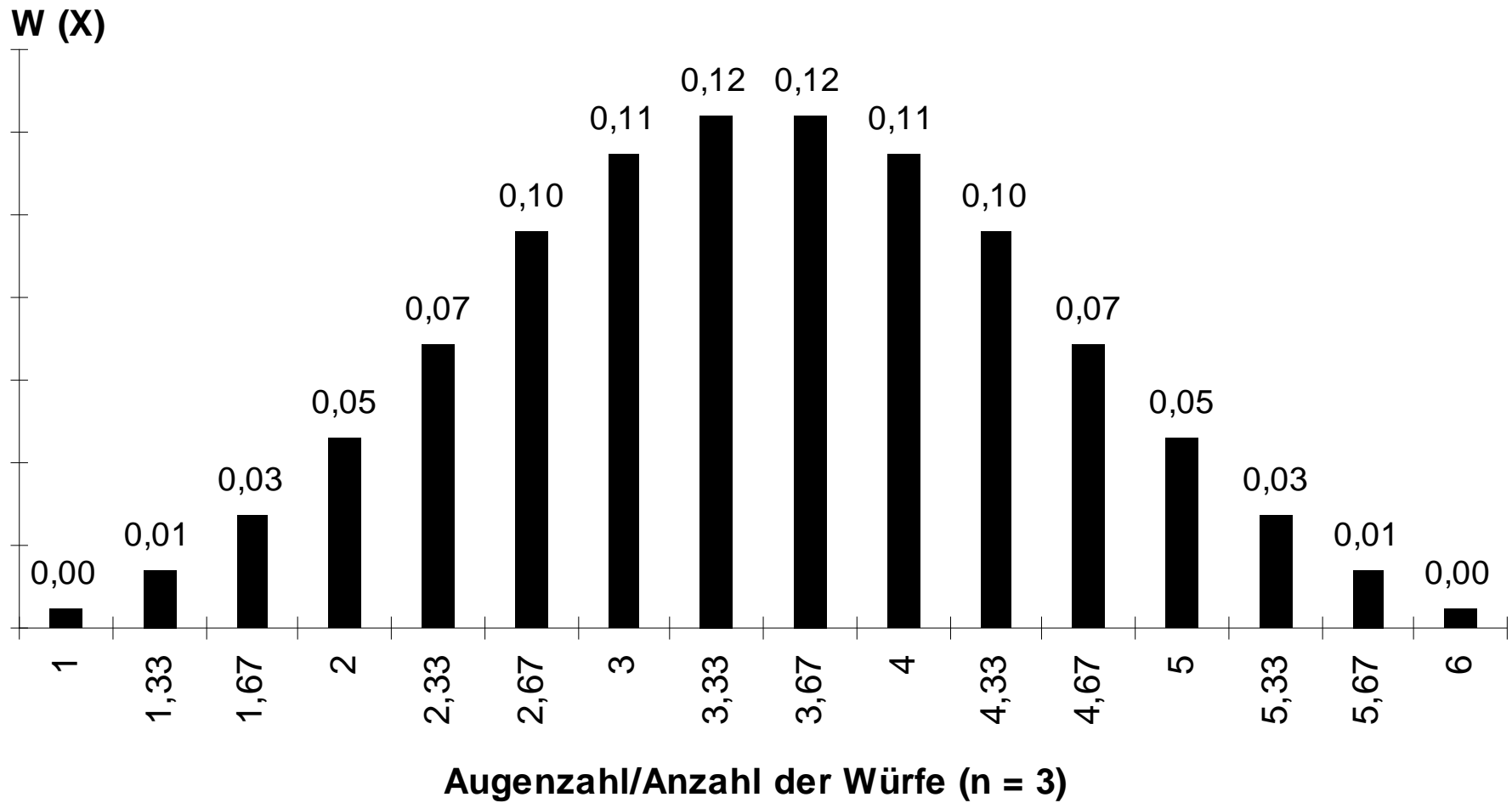
$$W(X=a) = 0$$

Beispiel einmaliges Würfeln: $W(X=7) = 0$.

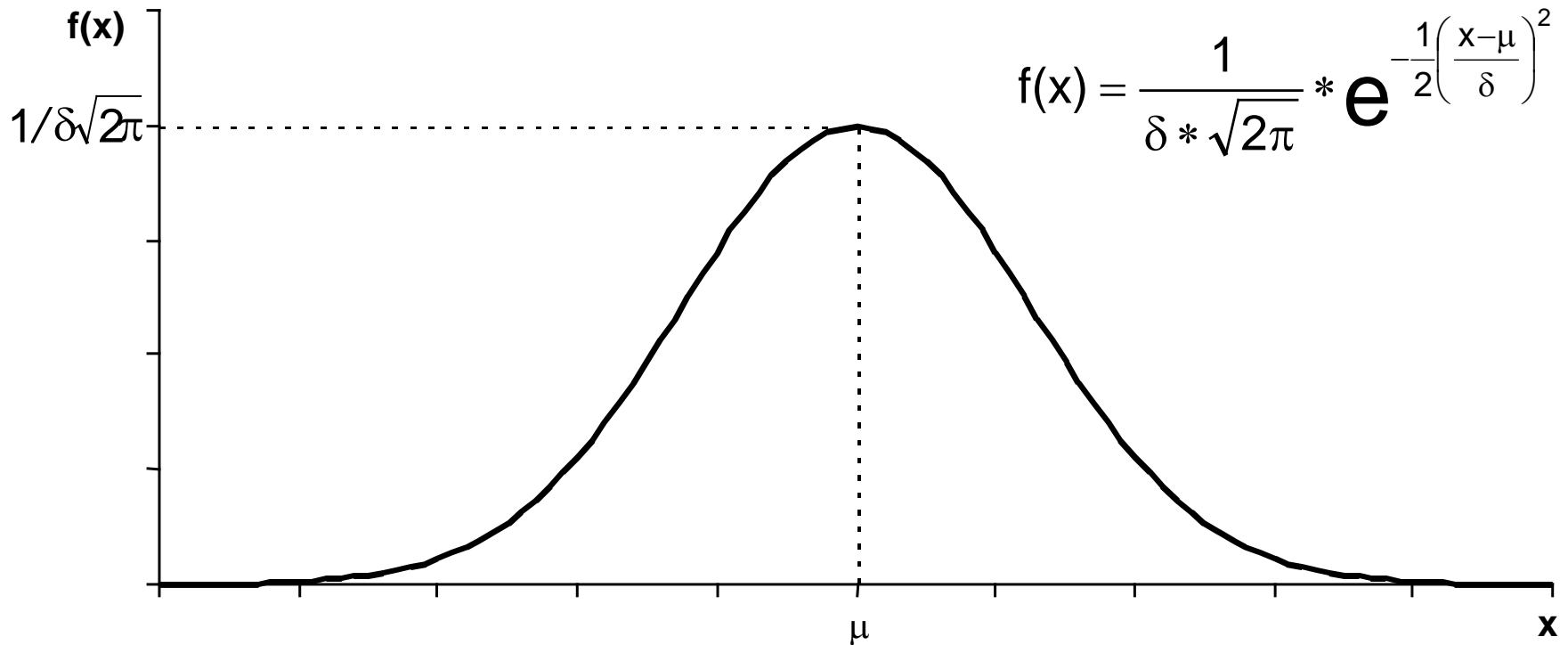
Wahrscheinlichkeiten der Ereignisse beim zweimaligen Würfelwurf



Wahrscheinlichkeiten der Ereignisse beim dreimaligen Würfelwurf



Dichtefunktion der Normalverteilung (Gauß'sche Glockenkurve)

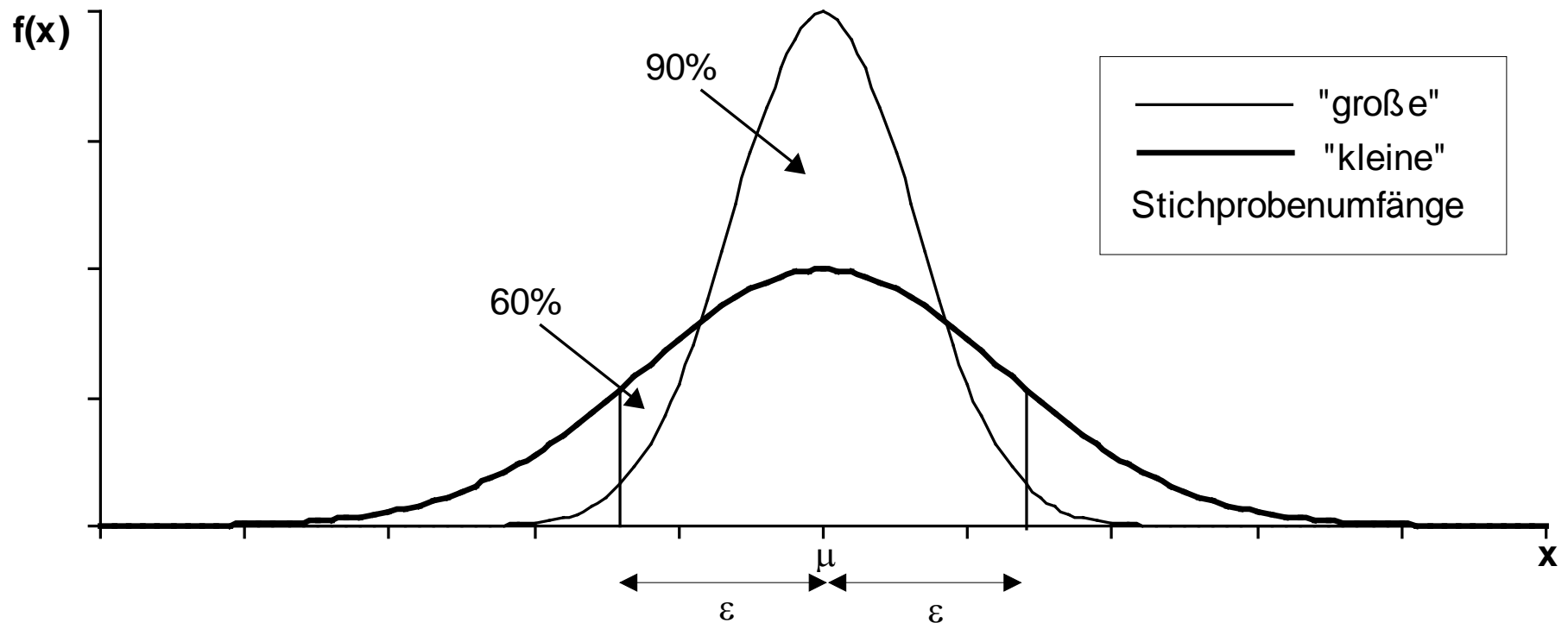


mit

δ = Standardabweichung und
 μ = tatsächlicher Mittelwert.

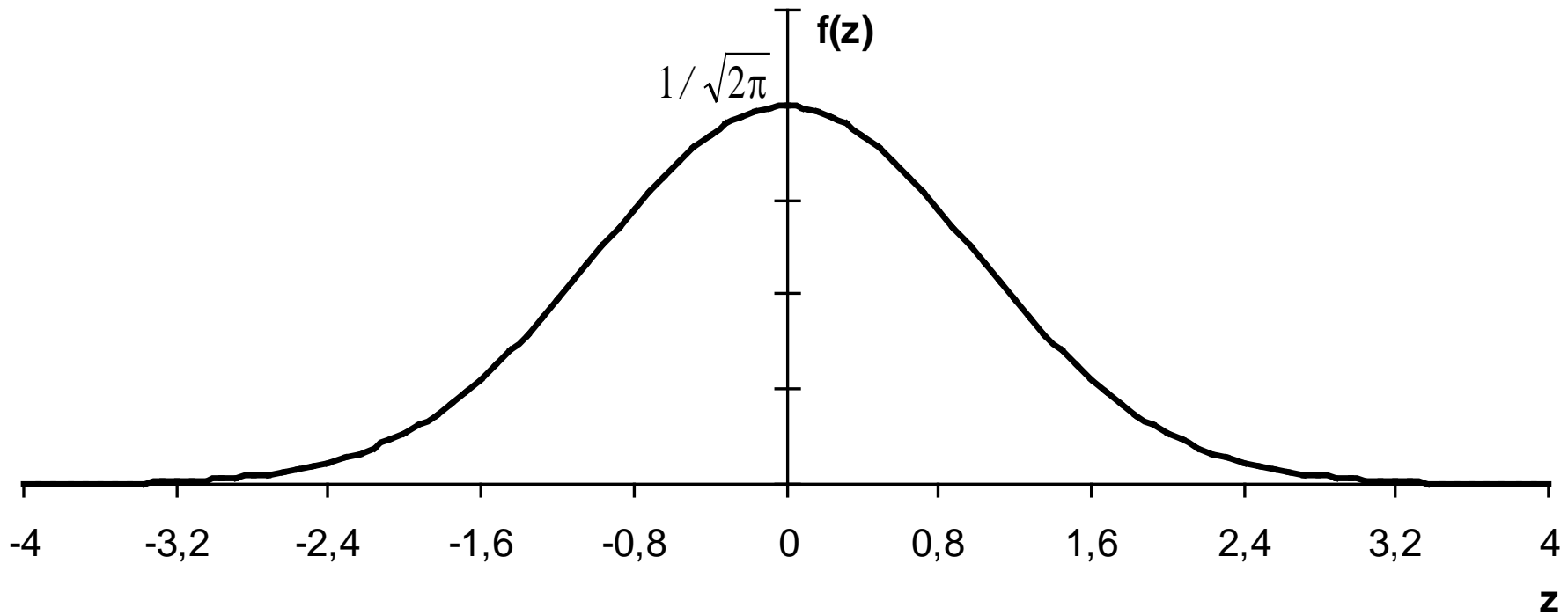
Verlauf der Dichtefunktion für „kleine“ und für „große“ Stichprobenumfänge

(Prozentangaben zum jeweiligen Flächeninhalt sind geschätzt)



Die Standardnormalverteilungsfunktion

führe z-Transformation (Standardisierung) durch: $z = \frac{x - \mu}{\sigma}$



$$f(z) = \frac{1}{\sqrt{2\pi}} * e^{-\frac{1}{2}(z)^2}$$

Tafel 2 Die Verteilungsfunktion $\Phi(z)$ und die Funktion $D(z)$ der Standardnormalverteilung

Quelle: KREYSZIG 1968, S. 393-394

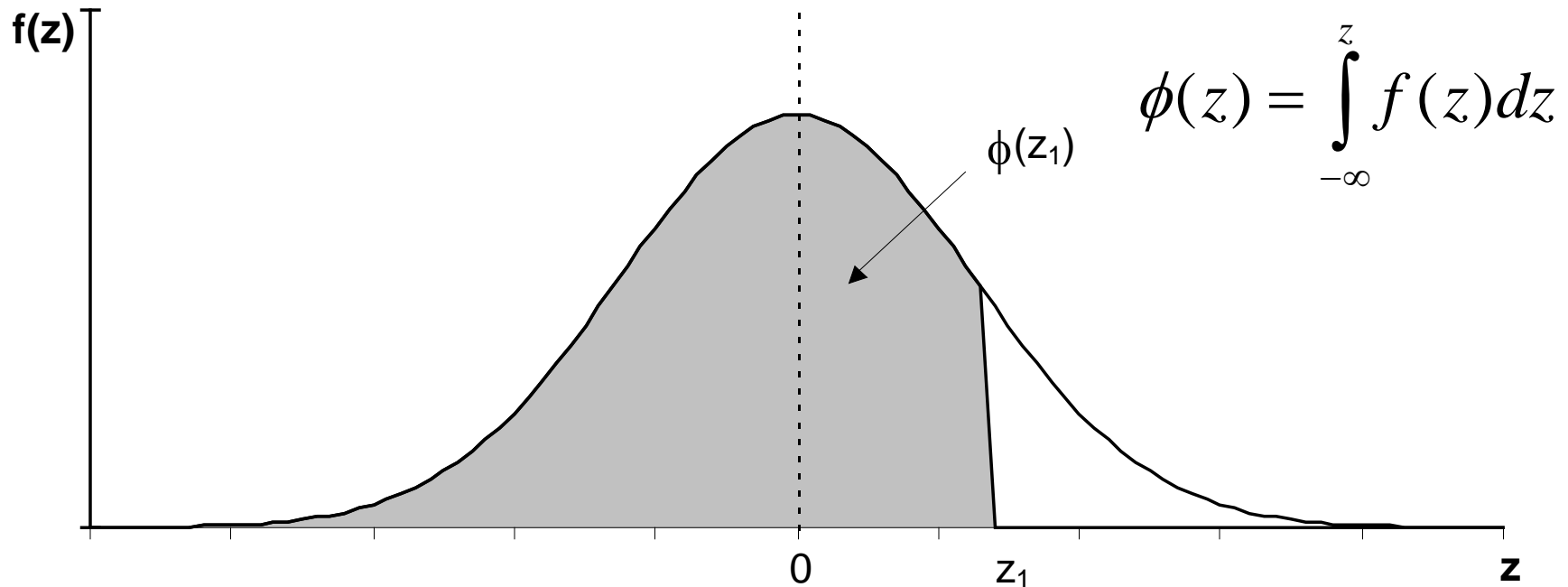
$$\Phi(-z) = 1 - \Phi(z) \quad D(z) = \Phi(z) - \Phi(-z)$$

z	$\Phi(z)$	$D(z)$	z	$\Phi(z)$	$D(z)$	z	$\Phi(z)$	$D(z)$	z	$\Phi(z)$	$D(z)$
0,01	5040	0080	0,43	6664	3328	0,85	8023	6047	1,27	8980	7959
0,02	5080	0160	0,44	6700	3401	0,86	8051	6102	1,28	8997	7995
0,03	5120	0239	0,45	6736	3473	0,87	8078	6157	1,29	9015	8029
0,04	5160	0319	0,46	6772	3545	0,88	8106	6211	1,30	9032	8064
0,05	5199	0399	0,47	6608	3616	0,89	8133	6265	1,31	9049	8098
0,06	5239	0478	0,48	6844	3688	0,90	8159	6319	1,32	9066	8132
0,07	5279	0558	0,49	6879	3759	0,91	8186	6372	1,33	9082	8165
0,08	5319	0638	0,50	6915	3829	0,92	8212	6424	1,34	9099	8198
0,09	5359	0717	0,51	6950	3899	0,93	8238	6476	1,35	9115	8230
0,10	5398	0797	0,52	6958	3969	0,94	8264	6528	1,36	9131	8262
0,11	5438	0876	0,53	7019	4039	0,95	8289	6579	1,37	9147	8293
0,12	5478	0955	0,54	7054	4108	0,96	8315	6629	1,38	9162	8324
0,13	5517	1034	0,55	7088	4177	0,97	8340	6680	1,39	9177	8355
0,14	5557	1113	0,56	7123	4245	0,98	8365	6729	1,40	9192	8385
0,15	5596	1192	0,57	7157	4313	0,99	8389	6778	1,41	9207	8415
0,16	5636	1271	0,58	7190	4381	1,00	8413	6827	1,42	9222	8444
0,17	5675	1350	0,59	7224	4448	1,01	8438	6875	1,43	9236	8473
0,18	5714	1428	0,60	7257	4515	1,02	8461	6923	1,44	9251	8501
0,19	5753	1507	0,61	7291	4581	1,03	8485	6970	1,45	9265	8529
0,20	5793	1585	0,62	7324	4647	1,04	8508	7017	1,46	9279	8557
0,21	5832	1663	0,63	7357	4713	1,05	8531	7063	1,47	9292	8584
0,22	5871	1741	0,64	7389	4778	1,06	8554	7109	1,48	9306	8611
0,23	5910	1819	0,65	7422	4843	1,07	8577	7154	1,49	9319	8638
0,24	5948	1897	0,66	7454	4907	1,08	8599	7199	1,50	9332	8664
0,25	5987	1974	0,67	7486	4971	1,09	8621	7243	1,51	9345	8690
0,26	6026	2051	0,68	7517	5035	1,10	8643	7287	1,52	9357	8715
0,27	6064	2128	0,69	7549	5098	1,11	8665	7330	1,53	9370	8740
0,28	6103	2205	0,70	7580	5161	1,12	8686	7373	1,54	9382	8764
0,29	6141	2282	0,71	7611	5223	1,13	8708	7415	1,55	9394	8789
0,30	6179	2358	0,72	7642	5285	1,14	8729	7457	1,56	9406	8812
0,31	6217	2434	0,73	7673	5346	1,15	8749	7499	1,57	9418	8836
0,32	6255	2510	0,74	7704	5407	1,16	8770	7540	1,58	9429	8859
0,33	6293	2586	0,75	7734	5467	1,17	8790	7580	1,59	9441	8882
0,34	6331	2661	0,76	7764	5527	1,18	8810	7620	1,60	9452	8904
0,35	6368	2737	0,77	7794	5587	1,19	8830	7660	1,61	9463	8926
0,36	6406	2812	0,78	7823	5646	1,20	8849	7699	1,62	9474	8948
0,37	6443	2886	0,79	7852	5705	1,21	8869	7737	1,63	9484	8969
0,38	6480	2961	0,80	7881	5763	1,22	8888	7775	1,64	9495	8990
0,39	6517	3035	0,81	7910	5821	1,23	8907	7813	1,65	9505	9011
0,40	6554	3108	0,82	7939	5878	1,24	8925	7850	1,66	9515	9031
0,41	6591	3182	0,83	7967	5935	1,25	8944	7887	1,67	9525	9051
0,42	6628	3255	0,84	7995	5991	1,26	8962	7923	1,68	9535	9070

Tafel 2 (Fortsetzung)

z	$\Phi(z)$	$D(z)$	z	$\Phi(z)$	$D(z)$	z	$\Phi(z)$	$D(z)$
1,69	9545	9090	2,13	9834	9668	2,57	9949	9898
1,70	9554	9109	2,14	9838	9675	2,58	9951	9901
1,71	9564	9127	2,15	9842	9684	2,59	9952	9904
1,72	9573	9146	2,16	9846	9692	2,60	9953	9907
1,73	9582	9164	2,17	9850	9700	2,61	9955	9909
1,74	9591	9181	2,18	9854	9707	2,62	9956	9912
1,75	9599	9199	2,19	9857	9715	2,63	9957	9915
1,76	9608	9216	2,20	9861	9722	2,64	9959	9917
1,77	9616	9233	2,21	9864	9729	2,65	9960	9920
1,78	9625	9249	2,22	9868	9736	2,66	9961	9922
1,79	9633	9265	2,23	9871	9743	2,67	9962	9924
1,80	9641	9281	2,24	9875	9749	2,68	9963	9926
1,81	9649	9297	2,25	9878	9756	2,69	9964	9929
1,82	9656	9312	2,26	9881	9762	2,70	9965	9931
1,83	9664	9328	2,27	9884	9768	2,71	9966	9933
1,84	9671	9342	2,28	9887	9774	2,72	9967	9935
1,85	9678	9357	2,29	9890	9780	2,73	9968	9937
1,86	9686	9371	2,30	9893	9786	2,74	9969	9939
1,87	9693	9385	2,31	9896	9791	2,75	9970	9940
1,88	9699	9399	2,32	9898	9797	2,76	9971	9942
1,89	9706	9412	2,33	9901	9802	2,77	9972	9944
1,90	9713	9426	2,34	9904	9807	2,78	9973	9946
1,91	9719	9439	2,35	9906	9812	2,79	9974	9947
1,92	9726	9451	2,36	9909	9817	2,80	9974	9949
1,93	9732	9464	2,37	9911	9822	2,81	9975	9950
1,94	9738	9476	2,38	9913	9827	2,82	9976	9952
1,95	9744	9488	2,39	9916	9832	2,83	9977	9953
1,96	9750	9500	2,40	9918	9836	2,84	9977	9955
1,97	9756	9512	2,41	9920	9840	2,85	9978	9956
1,98	9761	9523	2,42	9922	9845	2,86	9979	9958
1,99	9767	9534	2,43	9925	9849	2,87	9979	9959
2,00	9772	9545	2,44	9927	9853	2,88	9980	9960
2,01	9778	9556	2,45	9929	9857	2,89	9981	9961
2,02	9783	9566	2,46	9931	9861	2,90	9981	9963
2,03	9788	9576	2,47	9932	9865	2,91	9982	9964
2,04	9793	9586	2,48	9934	9869	2,92	9982	9965
2,05	9798	9596	2,49	9936	9872	2,93	9983	9966
2,06	9803	9606	2,50	9938	9876	2,94	9984	9967
2,07	9808	9615	2,51	9940	9879	2,95	9984	9968
2,08	9812	9625	2,52	9941	9883	2,96	9985	9969
2,09	9817	9634	2,53	9943	9886	2,97	9985	9970
2,10	9821	9643	2,54	9945	9889	2,98	9986	9971
2,11	9826	9651	2,55	9946	9892	2,99	9986	9972
2,12	9830	9660	2,56	9948	9895	3,00	9987	9973

Linksseitige Wahrscheinlichkeit der Standardnormalverteilung



Es gilt: $\phi(-z) = 1 - \phi(z)$

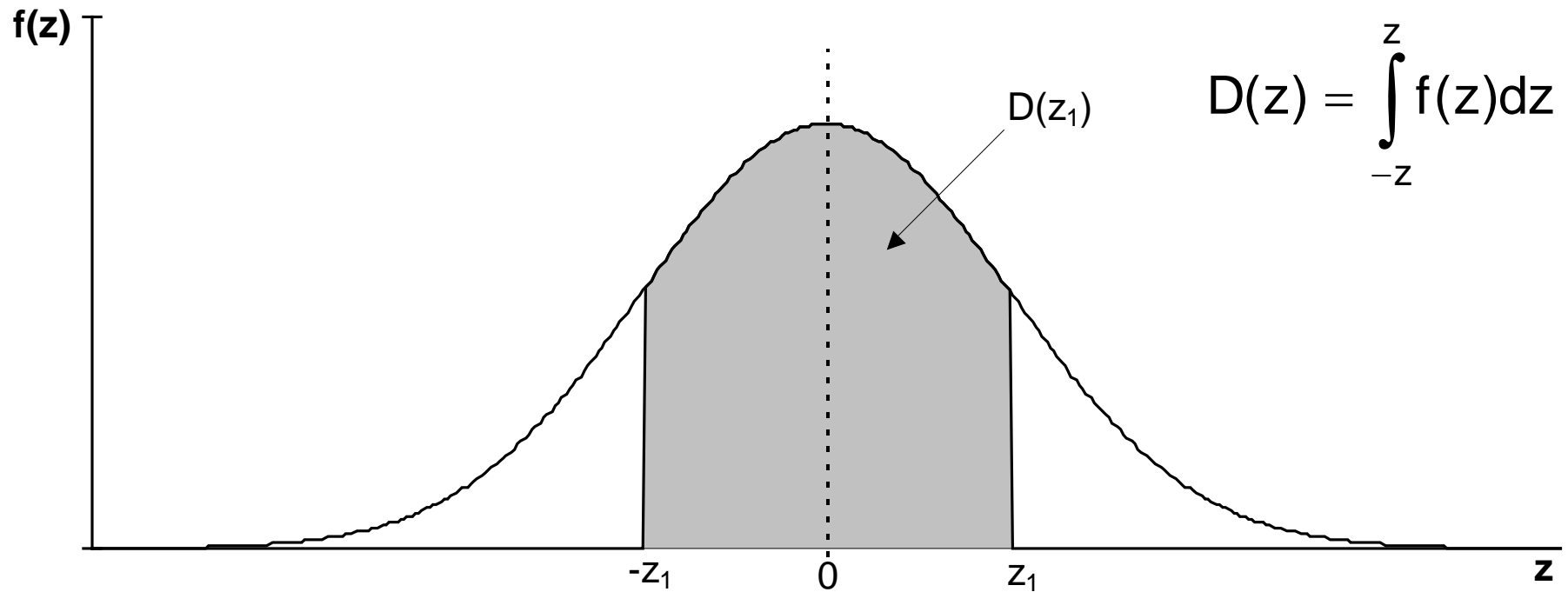
Beispiele für linksseitige Wahrscheinlichkeiten

a) $\phi(0,5) = W(Z \leq 0,5) = 0,6915 = 69,15\%$

b) $\phi(2) = W(Z \leq 2) = 0,9772 = 97,72\%$

c) $\phi(-1) = W(Z \leq -1) = 1 - \phi(1) = 1 - 0,8413 = 0,1587 = 15,87\%$

Zentrale Wahrscheinlichkeit der Standardnormalverteilung



Beispiele für zentrale Wahrscheinlichkeiten:

a) $D(0,5) = W(-0,5 \leq Z \leq 0,5) = 0,3829 = 38,29\%$

b) $D(1,0) = W(-1,0 \leq Z \leq 1,0) = 0,6827 = 68,27\%$

Rechenbeispiel für eine Über- und Unterschreitungswahrscheinlichkeit:

(übernommen von Dr. Merja Hoppe, Universität Marburg)

Ein Fluss soll eingedeicht werden, um ein Gebiet vor Überschwemmungen zu schützen. Aus langjährigen Messreihen der jährlichen Höchstwasserstände wurden folgende Werte ermittelt:

Mittelwert: $\mu = 4 \text{ m}$

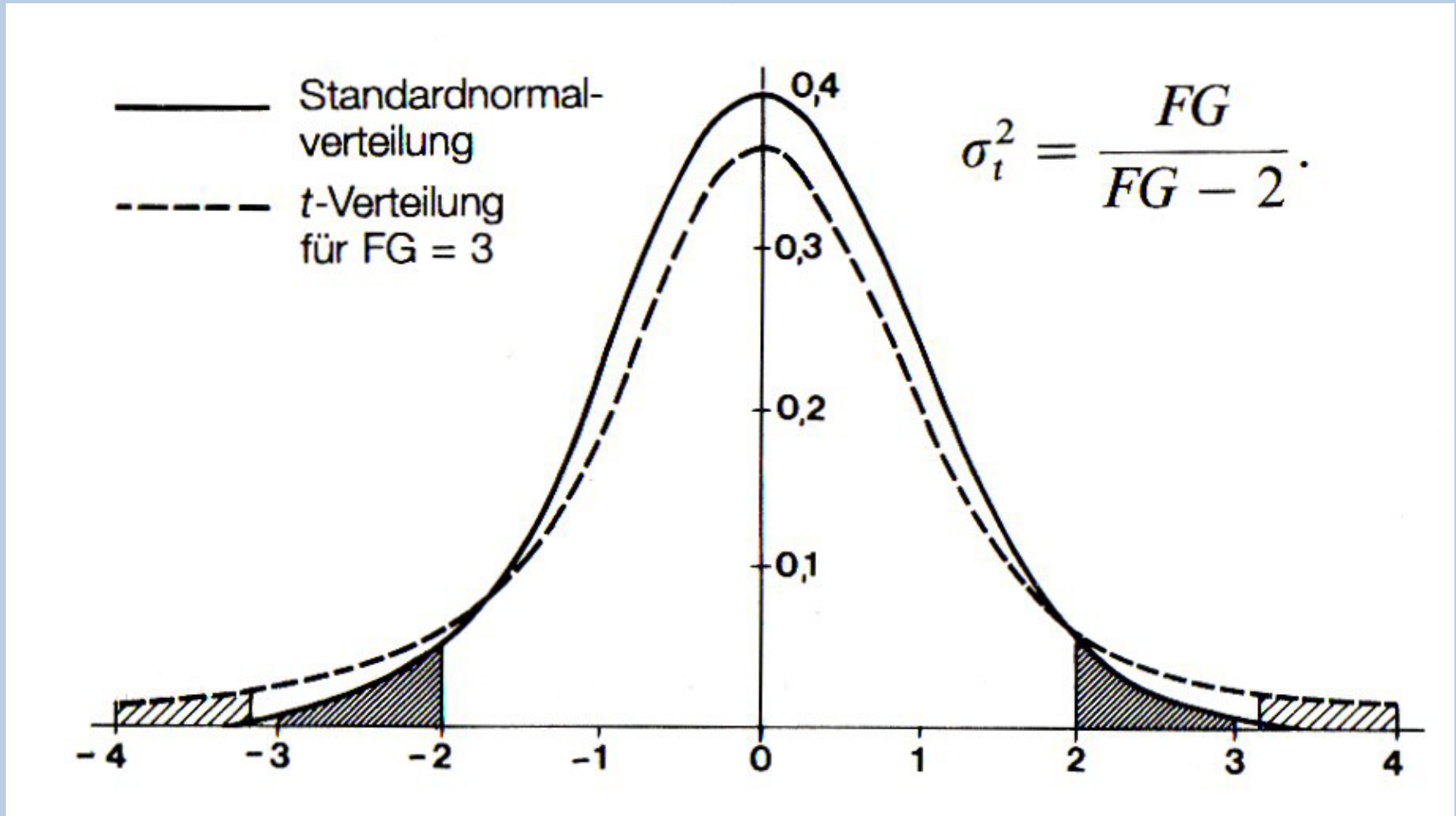
Standardabweichung: $\sigma = 2 \text{ m}$

- a) Mit welcher Wahrscheinlichkeit überschreitet der jährliche Höchstwasserstand die Marke von 7m?
- b) Welchen Mindestwert erreicht das sogenannte 100jährige Hochwasser?
- c) Wie groß ist die Wahrscheinlichkeit, dass die jährlichen Höchstwasserstände zwischen 4m und 7m liegen?

Gehen Sie davon aus, dass die jährlichen Höchststände normalverteilt sind.

7.2 t-Verteilung und t-Test

Die t-Verteilung (für FG=3) im Vergleich zur Standardnormalverteilung



t-Verteilung für verschiedene Signifikanzniveaus und Freiheitsgrade

Es gilt $F(t) = 1 - \alpha$, z.B. $F(6,31) = 0,95$ bei $FG = 1$. Bei zweiseitiger Fragestellung müssen die Signifikanzniveaus verdoppelt werden

FG	Signifikanzniveaus in %					
	10%	5%	2,5%	1%	0,5%	0,1%
1	3,08	6,31	12,7	31,8	63,7	318,3
2	1,89	2,92	4,30	6,97	9,93	22,3
3	1,64	2,35	3,18	4,54	5,84	10,2
4	1,53	2,13	2,78	3,75	4,60	7,17
5	1,48	2,02	2,57	3,37	4,03	5,89
6	1,44	1,94	2,45	3,14	3,71	5,21
7	1,42	1,90	2,37	3,00	3,50	4,79
8	1,40	1,86	2,31	2,90	3,36	4,50
9	1,38	1,83	2,26	2,82	3,25	4,30
10	1,37	1,81	2,23	2,76	3,17	4,14
11	1,36	1,80	2,20	2,72	3,11	4,03
12	1,36	1,78	2,18	2,68	3,06	3,93
13	1,35	1,77	2,16	2,65	3,01	3,85
14	1,35	1,76	2,15	2,62	2,98	3,79
15	1,34	1,75	2,13	2,60	2,95	3,73
16	1,34	1,75	2,12	2,58	2,92	3,69
17	1,33	1,74	2,11	2,57	2,90	3,65
18	1,33	1,73	2,10	2,55	2,88	3,61
19	1,33	1,73	2,09	2,54	2,86	3,58
20	1,33	1,73	2,09	2,53	2,85	3,55
22	1,32	1,72	2,07	2,51	2,82	3,51
24	1,32	1,71	2,06	2,49	2,80	3,47
26	1,32	1,71	2,06	2,48	2,78	3,44
28	1,31	1,70	2,05	2,47	2,76	3,41
30	1,31	1,70	2,04	2,46	2,75	3,39
40	1,30	1,68	2,02	2,42	2,70	3,31
50	1,30	1,68	2,01	2,40	2,68	3,26
100	1,29	1,66	1,98	2,37	2,63	3,17
∞	1,28	1,65	1,96	2,33	2,60	3,09

Freiheitsgrade (FG)

Der Freiheitsgrad einer Prüfverteilung ist die Anzahl der Stichprobenelemente, die zur Berechnung eines statistischen Parameters aus der Stichprobe notwendig sind.

Beispiel: Wird als Prüfgröße der Mittelwertes einer Stichprobe errechnet, so sind dazu alle Stichprobenelemente n erforderlich. Entsprechend ist der Freiheitsgrad $FG = n$.

Für die Standardabweichung werden noch $n-1$ Elemente der Stichprobe gebraucht, denn das letzte Element lässt sich theoretisch aus dem zuvor bestimmten Mittelwert berechnen. Also $FG = n-1$.

Allgemein gilt : Der Freiheitsgrad ist gleich der Anzahl der Stichprobenelemente minus der Anzahl der bereits bestimmten statistischen Parameter, die zur Berechnung weiterer Parameter notwendig sind.

Allgemeines Testverfahren: Anwendungsbeispiel t-Test

Die Oberböden in der Nordeifel werden anhand von $n = 20$ Proben auf ihren Zinkgehalt geprüft und mit dem Standardwert für Zink in Tongesteinen verglichen. Dieser Standardwert liegt bei $\mu = 95$ ppm (Tongesteinstandard)

Der Mittelwert der 20 Proben unterscheidet sich vom Standardwert.
Er liegt bei $\bar{x} = 91$ ppm,
die Standardabweichung der 20 Proben beträgt $\sigma = 10$ ppm.

Frage: Liegt dieser Unterschied im Zufallsbereich (z.B. wegen einer zu geringen Zahl der Proben) und die Proben sind bezüglich ihres Zinkgehaltes nicht repräsentativ oder ist der Unterschied nicht zufällig (signifikant) und der Boden in der Nordeifel ist relativ arm an Zink?

Gegeben:

Standardwert für Zink in Tongesteinen: $a = 95 \text{ ppm}$

20 Proben mit $\mu = 91 \text{ ppm}$,

und $\sigma = 10 \text{ ppm}$.

Frage: Liegt dieser Unterschied im Zufallsbereich (z.B. wegen einer zu geringen Zahl der Proben) und die Proben sind bezüglich ihres Zinkgehaltes nicht repräsentativ oder ist der Unterschied nicht zufällig (signifikant) und der Boden in der Nordeifel ist relativ arm an Zink?

Lösungsweg:

Führe einseitigen t-Test durch, weil die Richtung der Abweichung ($<$ oder $>$)

In diesem Fall klar ist, denn der Mittelwert der Messwerte ist kleiner als der Referenzwert. Abweichungen nach oben müssen also nicht betrachtet werden.

a) Formulierung der Hypothesen H_0 (Nullhypothese) und H_A (Alternativhypothese)

H_0 : Zwischen dem Tongesteinstandardwert und dem Mittelwert der Stichprobe besteht kein signifikanter Unterschied. (Die Abweichung ist zufällig durch die Stichprobenauswahl bestimmt).

H_A : Der Stichprobenwert ist signifikant kleiner als der Standardwert.

b) Festlegung des Signifikanzniveaus (α)

$\alpha = 5\%$:	signifikant *
$\alpha = 1\%$:	sehr signifikant **
$\alpha = 0,1\%$:	hochsignifikant ***

Beispiel:

Der Mittelwert der Stichprobe soll zum $\alpha = 5\%$ Signifikanz-Niveau getestet werden.

c) Bestimmung der Prüfgröße hinsichtlich der Abweichung und des entsprechenden Vergleichswertes aus der Tabelle der t-Verteilung

Die Prüfgröße für einen solchen Mittelwerttest lautet:

$$t = \frac{\mu - a}{\sigma / \sqrt{n}}$$

mit

μ = Mittelwert der Stichprobe

a = Vergleichswert, an dem die Abweichung bemessen werden soll

σ = Standardabweichung der Stichprobe

n = Stichprobenumfang.

Die Zahl der Freiheitsgrade beträgt $FG = n-1 = 19$,
weil in die Prüfgröße die Standardabweichung der Stichprobe enthält.

Lese aus der Tabelle für die t-Verteilung den entsprechendem Wert für $FG=19$ und ein Signifikanzniveau $\alpha = 5\%$ ab:

$$t_{5\%;19} = 1,73$$

Wenn nun die Abweichung der Prüfgröße den Wert $t_{5\%;19} = 1,73$ übersteigt, dann liegt keine zufallsbedingte, sondern eine systematische (signifikante) Abweichung vor.

d) Berechne t für die Werte der Stichprobe

$$t = \frac{\mu - a}{\sigma / \sqrt{n}} = \frac{91 - 95}{10 / \sqrt{20}} = \frac{4}{10 / 4,472} = \frac{4}{2,236} = 1,7889$$

e) Ergebnis und Antwort:

$$t = 1,7889 > 1,73, = t_{5\%,19}$$

Der errechnete Wert liegt außerhalb der vorgegebenen Prüfgröße. Die Zinkarmut in den 20 Bodenproben der Nordeifel ist also gegenüber dem Standardwert keine zufällige, sondern eine signifikante Abweichung (Signifikanzniveau 5%).

Mit anderen Worten:

Die Irrtumswahrscheinlichkeit, dass die Böden der Nordeifel zinkarm sind, beträgt 5%.

Frage: Ist die Abweichung auch sehr signifikant ** oder gar hochsignifikant *?**

Vergleiche die Prüfgröße $t = 1,7889$ mit den Vergleichswerten für ein Signifikanzniveau $\alpha = 1\%$ und $\alpha = 0,1\%$

Es gilt: $t = 1,7889 < t_{1\%,19} = 2,54$
und $t = 1,7889 < t_{0,1\%,19} = 3,58$

Daher ist die Abweichung nur signifikant zum Signifikanzniveau $\alpha = 5\%$, aber nicht sehr Signifikant (Signifikanzniveau $\alpha = 1\%$) oder gar hochsignifikant ($\alpha = 0,1\%$).

2. Beispiel

Der durchschnittliche Intelligenzquotient (IQ) liegt bei 100.

In den beiden Statistik-Übungen wurden für die Teilnehmer die folgende Werte ermittelt:

Kurs A: $n = 33$

$$\mu = 105$$

$$\sigma = 11$$

Kurs B: $n = 71$

$$\mu = 103$$

$$\sigma = 9$$

Testen Sie für jeden Kurs, ob die Abweichungen signifikant sind und wenn ja, zu welchem Signifikanzniveau.