

PRINCIPALI DISTRIBUZIONI DI PROBABILITA'

Psicometria 1 - Lezione 9
Lucidi presentati a lezione

AA 2000/2001 dott. Corrado Caudek

DISTRIBUZIONE BINOMIALE

Possiamo definire un *processo bernoulliano* come una sequenza di n prove di un esperimento aleatorio tali per cui

1. ciascuna prova ha solo due esiti, che chiameremo successo e insuccesso.
2. La probabilità p di un successo in ciascuna prova è la stessa per tutte le prove e non è influenzata dagli esiti precedenti (le prove sono indipendenti). La probabilità di un insuccesso è $q = 1 - p$.

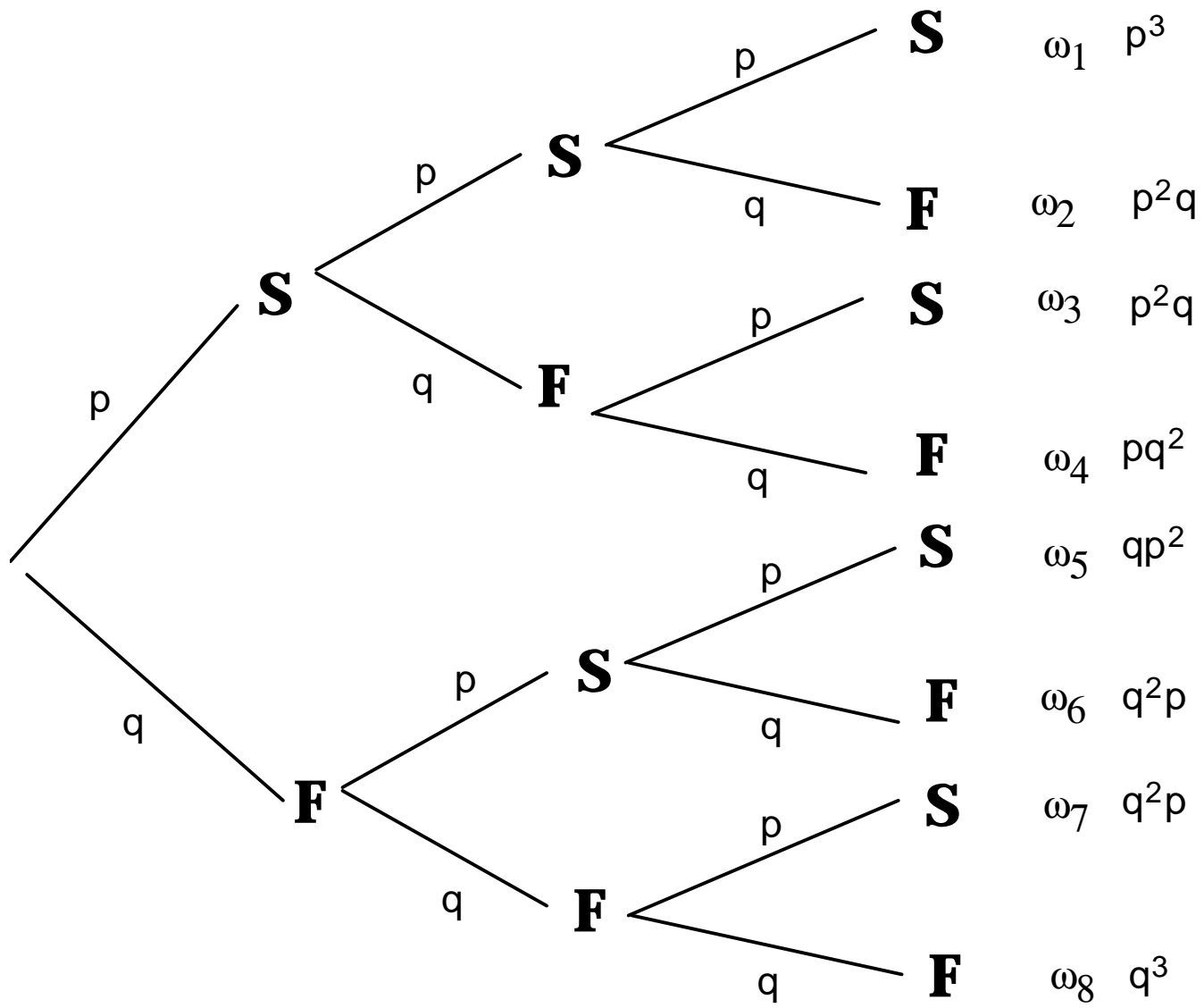
Esempi.

- Una moneta è lanciata per 10 volte. I due esiti possibili sono T e C. La probabilità di T è $1/2$.
- Un sondaggio è eseguito chiedendo a 1000 persone, scelte in maniera casuale, se sono favorevoli a una certa proposta di legge. I due esiti possibili sono si e no. La probabilità p di una risposta affermativa (cioè un successo) indica la proporzione di persone nella popolazione che sono favorevoli al disegno di legge.
- Uno giocatore d'azzardo fa una serie di puntate di 100 000 lire sul rosso o sul nero nel gioco della roulette. Qui un successo è vincere 100 000 lire e un fallimento è perdere 100 000 lire. Dal momento che il giocatore vince se la pallina si ferma su 18 delle possibili 37 posizioni della ruota, la probabilità di vincere è $p = 18/37 = .486$.

Il problema che ci poniamo è di stabilire la probabilità che può essere associata all'evento

"si osservano r successi in n prove di un processo bernoulliano".

Consideriamo ad esempio una sequenza di 3 lanci di una moneta con probabilità di successo (esito "testa") uguale a p .



Come si calcola la probabilità di ottenere una specifica sequenza di S e F ?

Le prove che costituiscono una sequenza sono, per definizione, *eventi indipendenti*.

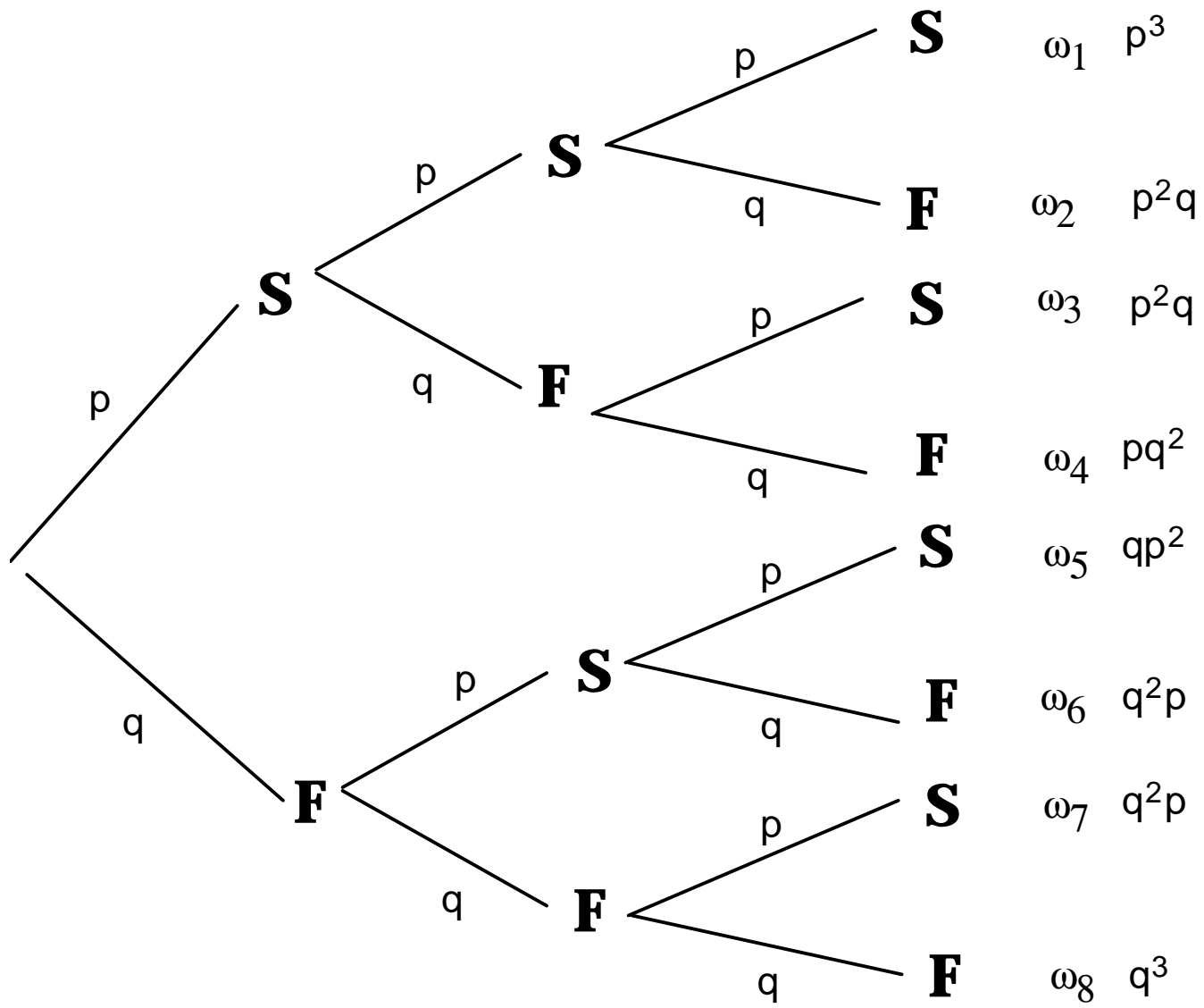
La probabilità che si verifichi una certa sequenza di eventi indipendenti è data dal prodotto delle probabilità degli eventi che la costituiscono.

Per esempio, alla sequenza SFS possiamo associare la probabilità $pqp = p^2q$.

Di solito, però, non siamo interessati alla probabilità di osservare *una specifica sequenza* di S e F . Invece, vogliamo conoscere la probabilità di osservare un dato numero di successi all'interno di una sequenza di n prove.

Vogliamo sapere, ad esempio, quale è la probabilità di osservare 2 successi in 3 lanci di una moneta, indipendentemente dall'ordine in cui S e F compaiono nella sequenza.

Per trovare la probabilità in questione dobbiamo stabilire *quanti sono i rami dell'albero* che rappresenta lo spazio campione che contengono 2 S e 1 F .



Esaminando la figura possiamo vedere che ci sono 3 rami dell'albero che rappresentano sequenze con 2 successi e 1 insuccesso.

A ciascuna di queste 3 sequenze è associata la stessa probabilità: p^2q .

Quindi, la probabilità di osservare 2 successi in 3 prove è dunque uguale a $3p^2q$.

[... la somma delle probabilità degli eventi w_i deve essere uguale a 1]

In questo modo possiamo trovare la probabilità di osservare 0, 1, 2, o 3 successi in 3 prove di un processo bernoulliano:

<i>Numero di successi</i>	<i>p</i>
0	$p^0 q^3$
1	$3 p q^2$
2	$3 p^2 q$
3	$p^3 q^0$

Poniamoci ora l'obiettivo di trovare una formula che ci consenta di produrre questa distribuzione di probabilità senza dovere costruire un diagramma ad albero come abbiamo fatto in precedenza.

E' chiaro che la probabilità di ottenere r successi in n prove, con probabilità di successo uguale a p e probabilità di un insuccesso uguale a $q = 1 - p$ in una specifica sequenza di S e F , è uguale a:

$$p^r q^{n-r}$$

Resta da stabilire quante sequenze si possono costruire in modo tale da produrre esattamente r successi e $(n - r)$ insuccessi in n prove.

Con 2 successi e 1 insuccesso, ad esempio,

$$U = \{S, S, F\}$$

possiamo creare 3 permutazioni:

$$\{S, S, F\}, \{S, F, S\}, \{F, S, S\}$$

Il numero di permutazioni di n elementi non tutti diversi tra loro si trova utilizzando il *coefficiente multinomiale*.

Nel caso presente, abbiamo un sottogruppo di r elementi che corrispondono all'esito "testa" (successo) e un sottogruppo di $(n - r)$ elementi che corrispondono all'esito "croce" (insuccesso).

Applicando la formula del coefficiente multinomiale otteniamo:

$$\frac{n!}{r!(n-r)!} = \frac{3!}{2! \times 1!} = 3$$

E' evidente che, nel caso di 2 soli sottoinsiemi, la formula è identica a quella del coefficiente binomiale:

$$\binom{n}{r} = \frac{n!}{r!(n-r)!}$$

In conclusione, la probabilità di osservare r successi in n prove di un processo bernoulliano si può calcolare con la formula:

$$b(n, p, r) = \binom{n}{r} p^r q^{n-r}$$

con $q = 1 - p$.

Esempio Un dado è lanciato 4 volte. Quale è la probabilità di ottenere un solo 6?

Trattiamo questo esperimento come un processo bernoulliano in cui la probabilità di un successo (ottenere un 6) è $1/6$ e la probabilità di un insuccesso (ottenere un numero che non è un 6) è $5/6$.

$$b(n = 4, p = 1/6, r = 1) = \binom{4}{1} \left(\frac{1}{6}\right)^1 \left(\frac{5}{6}\right)^3 = .3125$$

Esempio. Supponiamo che, tra quelli che seguono questo corso, soltanto uno studente su dieci giochi a tennis.

Quale è la probabilità di trovare 3 studenti che giocano a tennis in un campione casuale costituito da 5 studenti?

$$b(n = 5, p = .1, r = 3) = \binom{5}{3} \left(\frac{1}{10}\right)^3 \left(\frac{9}{10}\right)^2 = .00810$$

Esempio. Si trovi la probabilità di osservare 2 esiti "testa" in 5 lanci di una moneta onesta.

$$b(n = 5, p = .5, r = 2) = \binom{5}{2} (.5)^2 (.5)^3 = .3125$$

L'insieme costituito da tutti i possibili numeri di successi che si possono ottenere in n prove di un processo bernoulliano, insieme alle relative probabilità, si dice *distribuzione binomiale*.

Numero delle volte in cui si osserva l'esito "testa"	Probabilità ($p = 1/2$)
0	$\binom{5}{0} \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^5 = .0312$
1	$\binom{5}{1} \left(\frac{1}{2}\right)^1 \left(\frac{1}{2}\right)^4 = .1562$
2	$\binom{5}{2} \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^3 = .3125$
3	$\binom{5}{3} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^2 = .3125$
4	$\binom{5}{4} \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^1 = .1562$
5	$\binom{5}{5} \left(\frac{1}{2}\right)^5 \left(\frac{1}{2}\right)^0 = .0312$
1.000	

$$p(X = 0; N = 5) = \binom{5}{0} \left(\frac{1}{10}\right)^0 \left(\frac{9}{10}\right)^5 = .59049$$

$$p(X = 1; N = 5) = \binom{5}{1} \left(\frac{1}{10}\right)^1 \left(\frac{9}{10}\right)^4 = .32850$$

$$p(X = 2; N = 5) = \binom{5}{2} \left(\frac{1}{10}\right)^2 \left(\frac{9}{10}\right)^3 = .07290$$

$$p(X = 3; N = 5) = \binom{5}{3} \left(\frac{1}{10}\right)^3 \left(\frac{9}{10}\right)^2 = .00810$$

$$p(X = 4; N = 5) = \binom{5}{4} \left(\frac{1}{10}\right)^4 \left(\frac{9}{10}\right)^1 = .00045$$

$$p(X = 5; N = 5) = \binom{5}{5} \left(\frac{1}{10}\right)^5 \left(\frac{9}{10}\right)^0 = .00001$$

ESERCIZI

E1. Nel corso della sua carriera il giocatore di basket Rossi ha segnato in media 3 canestri su 10 lanci.

Solitamente, in una partita, Rossi esegue 4 tiri a canestro.

Si trovi la probabilità che in una partita Rossi segni 2 canestri.

E2. Giovanni pretende di avere poteri extrasensoriali e dice di potere predire meglio del caso quale di due simboli è presente in una carta coperta. Per verificare la sua pretesa accetta di sottoporsi a un test costituito da una sequenza di 10 prove.

Decidiamo di credere a Giovanni *soltanto se le sue prestazioni (o un risultato ancora più estremo)* sono equivalenti a quelle che verrebbero ottenute in meno del 5% dei casi da una persona che si limitasse a tirare ad indovinare.

Trovate il numero di prove corrette che Giovanni deve effettuare affinché risulti credibile in base al criterio stabilito in precedenza.

E3. Un ristorante offre torte di ciliege e di mele e compra un eguale numero di torte di entrambi i tipi.

Ogni giorno 10 clienti chiedono una fetta di torta.

I clienti scelgono con eguale probabilità i due tipi di torta.

Quante fette di ciascun tipo di torta dovrebbero essere tenute nel frigorifero dal proprietario in maniera tale da garantire di avere una probabilità di circa .95 che ciascun cliente riceva la torta di suo gradimento?

Nel 19esimo secolo, Galton inventò un dispositivo chiamato *quincunx* che può essere usato per ottenere empiricamente la distribuzione binomiale.

<http://www.users.on.net/zhcchz/java/quincunx/quincunx.1.html>

VALORE ATTESO DELLA DISTRIBUZIONE BINOMIALE

Sia X_i una variabile bernoulliana, ovvero una variabile che assume il valore 1 nel caso di un successo e 0 nel caso di un insuccesso.

La probabilità di un successo è p .

La probabilità di un insuccesso è $q = 1 - p$.

Il valore atteso di **una singola prova** X_i è:

$$E(X_i) = 1 p + 0 (1 - p) = p$$

Definiamo la variabile aleatoria S_n come il numero di successi in n prove di un processo bernoulliano.

Quale è il valore atteso di S_n ?

Possiamo definire S_n come la somma di n variabili aleatorie X_i , ciascuna delle quali rappresenta una singola prova bernoulliana: $S_n = X_1 + X_2 + \dots + X_n$

$$\begin{aligned} E(S_n) &= E(X_1 + X_2 + \dots + X_n) = \\ &= E(X_1) + E(X_2) + \dots + E(X_n) = np \end{aligned}$$

In conclusione, il valore atteso di una variabile aleatoria S_n che rappresenta il *numero di successi in n prove* di un processo bernoulliano con probabilità di successo uguale a p è:

$$E(S_n) = np$$

VARIANZA DELLA DISTRIBUZIONE BINOMIALE

Sia X_i la variabile aleatoria che rappresenta l'esito di una singola prova di un processo bernoulliano ($X = 1$ nel caso di un successo, $X = 0$ nel caso di un insuccesso). Si calcoli la varianza di X .

$$E(X_i) = 1p + 0q = p$$

$$E(X_i^2) = 1^2 p + 0^2 q = p$$

$$V(X_i) = E(X_i^2) - (E(X_i))^2$$

$$= p - p^2 = p(1 - p) = pq$$

In precedenza abbiamo visto che, se X e Y sono due variabile aleatorie indipendenti, allora

$$V(X + Y) = V(X) + V(Y)$$

Definiamo la variabile aleatoria S_n come il numero di successi in n prove **indipendenti** di un processo bernoulliano. $S_n = X_1 + X_2 + \dots + X_n$

Quale è la varianza di S_n ?

$$\begin{aligned} V(S_n) &= V(X_1 + X_2 + \dots + X_n) \\ &= V(X_1) + V(X_2) + \dots + V(X_n) \end{aligned}$$

dato che tutte queste varianze sono uguali

$$= nV(X_i) = npq$$

In conclusione, la varianza di una variabile aleatoria S_n che rappresenta il numero di successi in n prove di un processo bernulliano con probabilità di successo uguale a p è:

$$V(S_n) = npq$$

ESERCIZIO

E4. Un gruppo di archeologi dispone di finanziamenti adeguati per 10 spedizioni. La probabilità che una spedizione abbia successo è 0.1. Assumete che le spedizioni siano indipendenti. I costi dell'equipaggiamento che viene sempre utilizzato in ciascuna spedizione sono di 20 000 Euro. Una spedizione che porta a dei ritrovamenti (successo) costa 30 000 Euro. Una spedizione senza ritrovamenti (insuccesso) costa 15 000 Euro.

1) Si trovi la media e la varianza del numero di spedizioni che hanno successo.

2) Si trovi il costo totale per le 10 spedizioni.

DISTRIBUZIONE NORMALE

A differenza della binomiale, la distribuzione normale riguarda **variabili continue**.

La sua importanza e il suo uso in statistica dipende da due fattori:

1. Molte variabili psicologiche sono distribuite in questo modo.
2. Il teorema del limite centrale.

La distribuzione normale è definita dalla seguente funzione di probabilità:

$$f(x) = \frac{1}{\mathbf{s} \sqrt{2\mathbf{p}}} e^{-(x-\mathbf{m})^2 / 2\mathbf{s}^2}$$

dove e è una costante con il valore di 2.771828, \mathbf{m} è la media della distribuzione e \mathbf{s}^2 è la varianza della distribuzione.

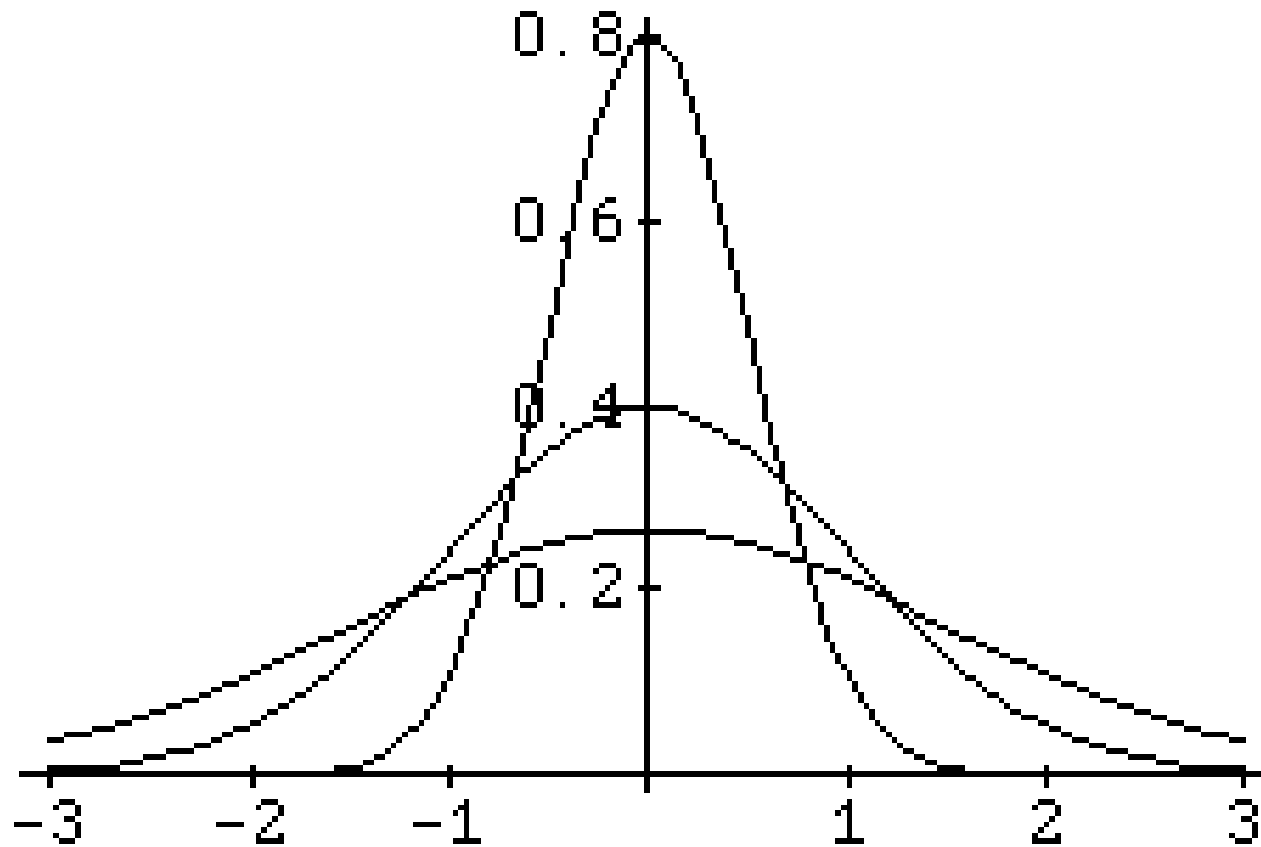
Questa funzione di probabilità fornisce un *modello teorico* che approssima la distribuzione di frequenze (relative) di molte variabili empiriche.

Gauss ha dimostrato, ad esempio, che gli *errori* di natura accidentale sono distribuiti in questo modo.

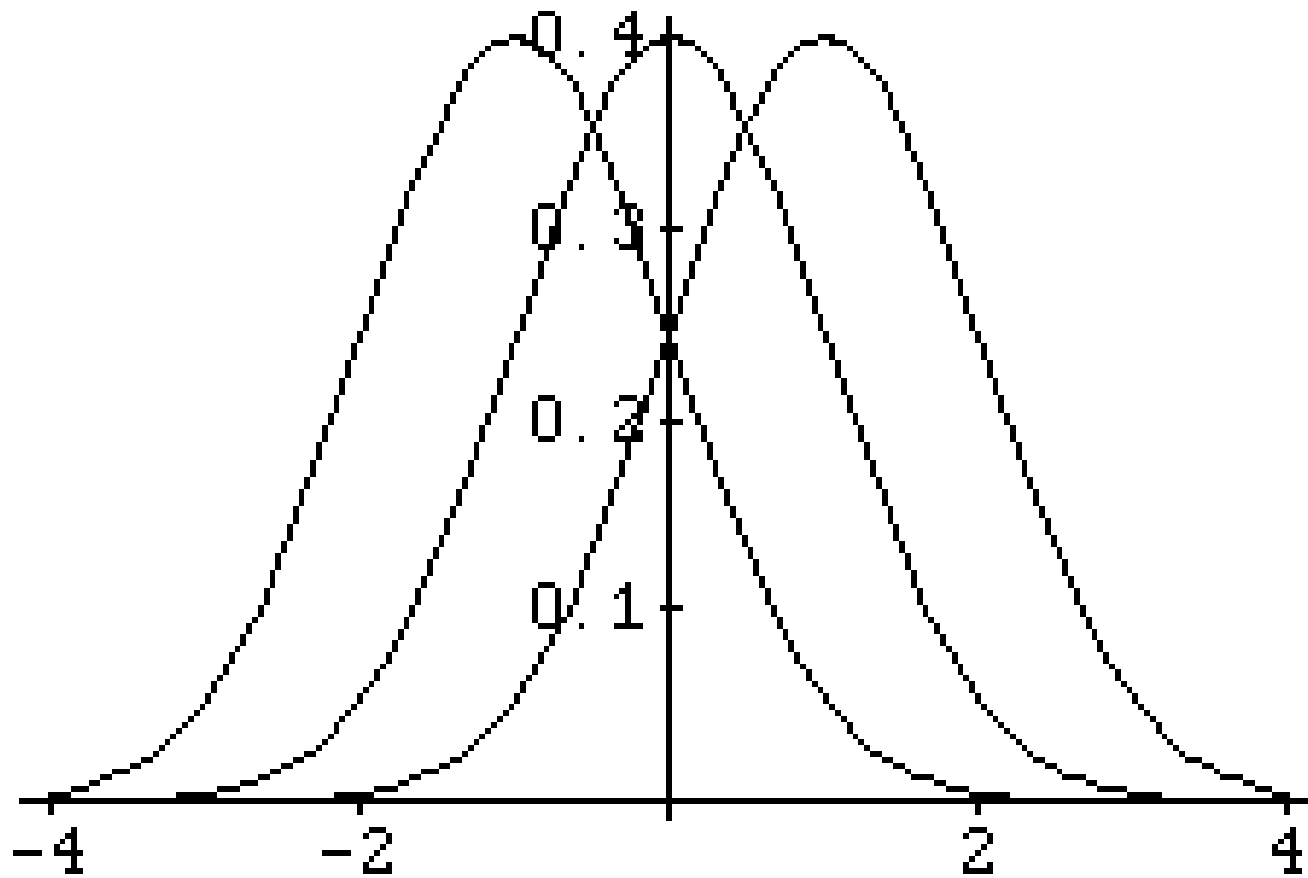
Come nel caso di tutte le distribuzioni di densità di probabilità, l'area sottesa alla curva in un qualunque intervallo della variabile aleatoria può essere interpretata come una *probabilità*.

L'area totale sottesa all'intera curva, dunque, è uguale a 1.0 per qualunque valore dei parametri μ e σ .

$\mu=0, \quad \sigma=.5, 1, 1.5$



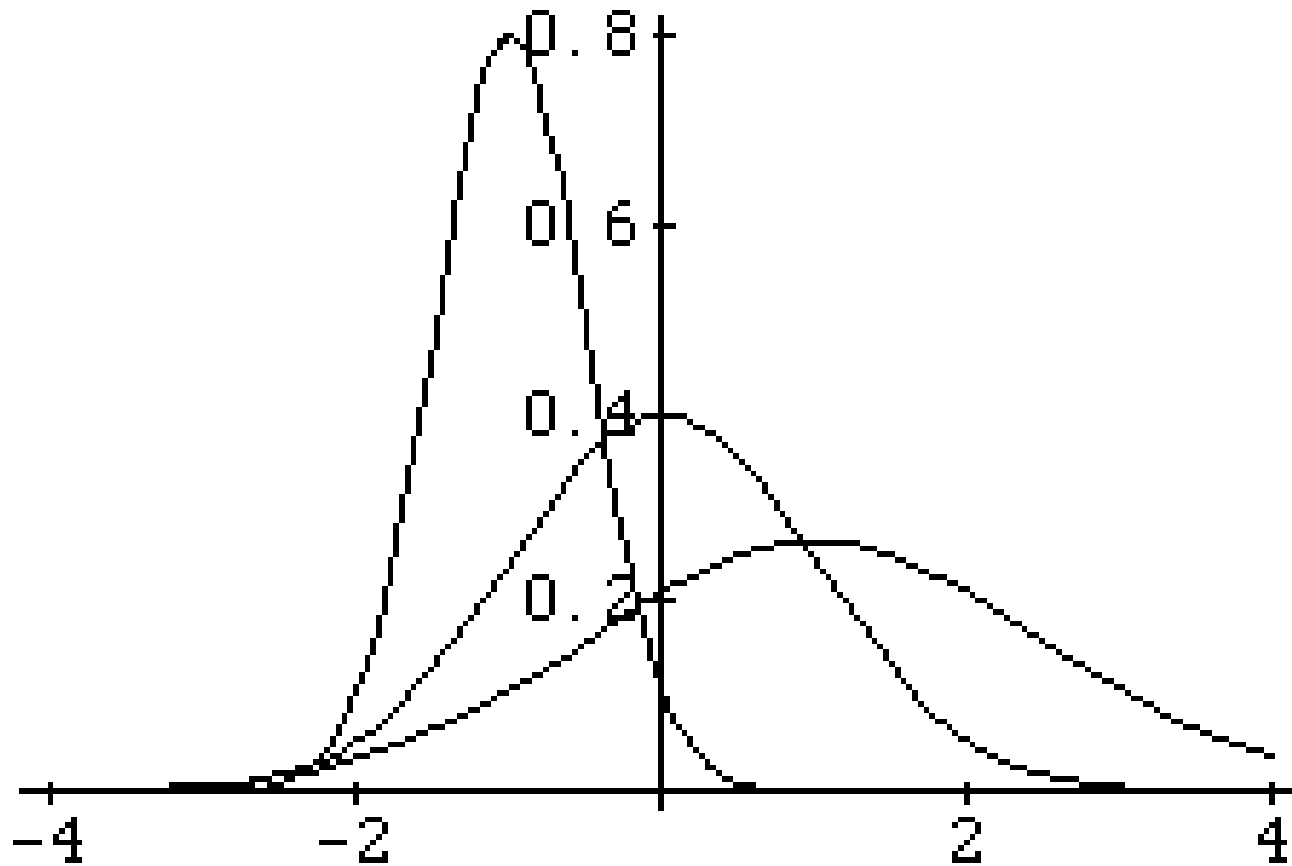
$$\mu = -1, 0, 1 \quad \sigma = 1$$



$\mu=0, \sigma=1$

$\mu=-1, \sigma=.5$

$\mu=1, \sigma=1.5$



La funzione che definisce la distribuzione normale associa a ciascun valore della variabile X un valore di densità di probabilità.

Valori di densità di probabilità diversi vengono trovati variando i due parametri μ e σ .

Questo significa che, cambiando questi due parametri, otteniamo delle curve diverse, come abbiamo visto negli esempi precedenti.

Come nel caso della distribuzione binomiale, dunque, la distribuzione normale in realtà è una famiglia di distribuzioni.

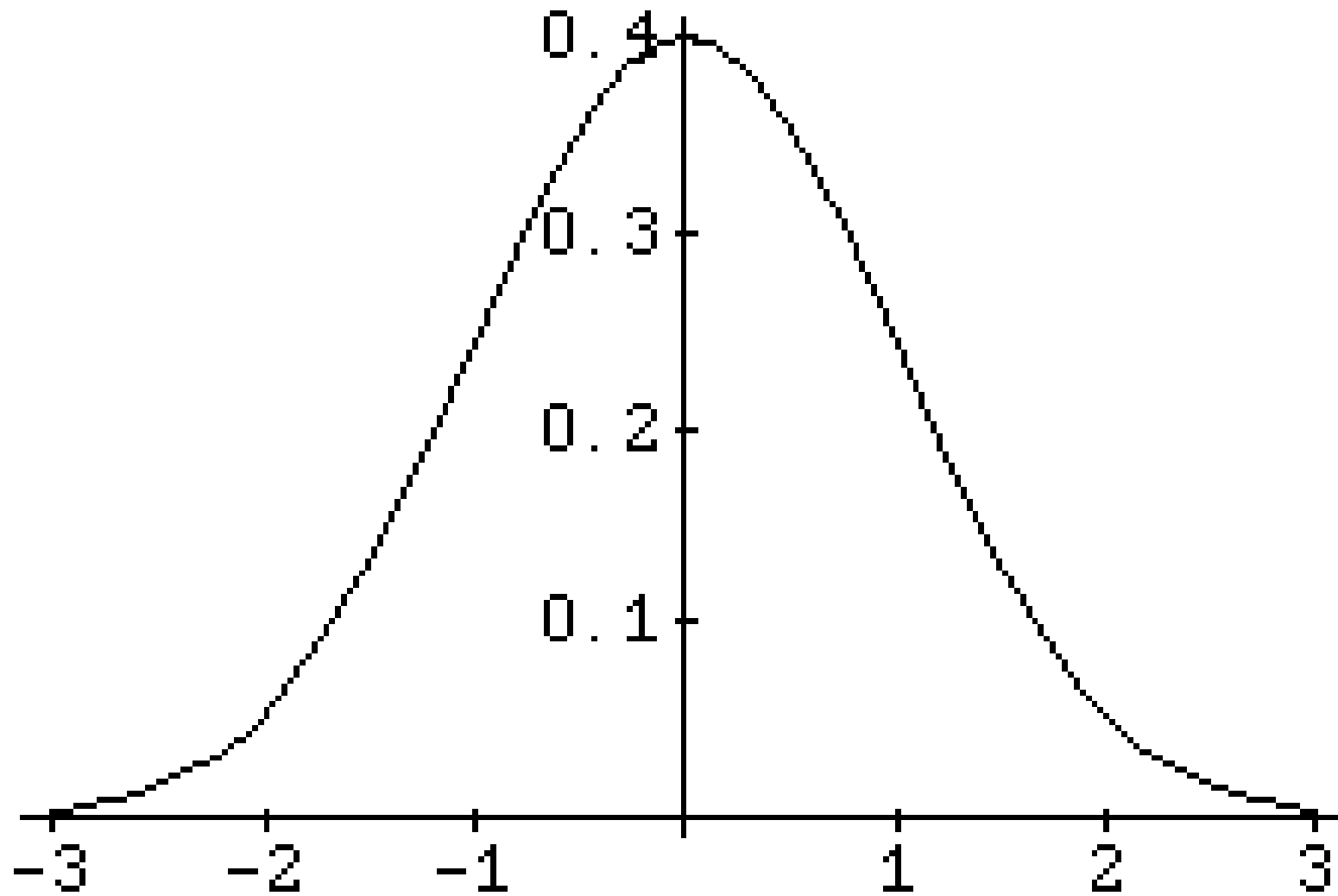
La notazione che si usa per fare riferimento alla distribuzione normale con media \mathbf{m} e varianza \mathbf{s}^2 è:

$$N(\mathbf{m}, \mathbf{s}^2)$$

Dato che la distribuzione normale specifica una famiglia di distribuzioni è vantaggioso considerare questa distribuzione nel caso di punteggi standardizzati.

Nel caso di una variabile aleatoria standardizzata (con media uguale $\mu = 0$ e scarto quadratico medio $\sigma = 1$) la funzione di probabilità della distribuzione normale si semplifica nel modo seguente:

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



Nell'uso che faremo della distribuzione normale ci occuperemo delle probabilità cumulative e delle probabilità definite in un intervallo di valori della variabile aleatoria.

La **probabilità cumulativa** (ovvero la probabilità che la variabile aleatoria assuma un valore compreso tra $-\infty$ e il valore a) $F(a) = p(X \leq a)$ corrisponde all'area sottesa alla curva normale nell'intervallo tra $-\infty$ e a .

Per mezzo delle probabilità cumulative possiamo calcolare la probabilità che la variabile aleatoria assuma un valore compreso all'interno di un qualsiasi intervallo di valori.

<http://www-stat.stanford.edu/~naras/jsm/NormalDensity/NormalDensity.html>

Esempio. Se una distribuzione normale ha una media di 50 e deviazione standard uguale a 5, quale è la probabilità cumulativa corrispondente al punteggio di 57.5?

Innanzitutto, calcoliamo il punteggio standardizzato:

$$z = \frac{57.5 - 50}{5} = 1.5$$

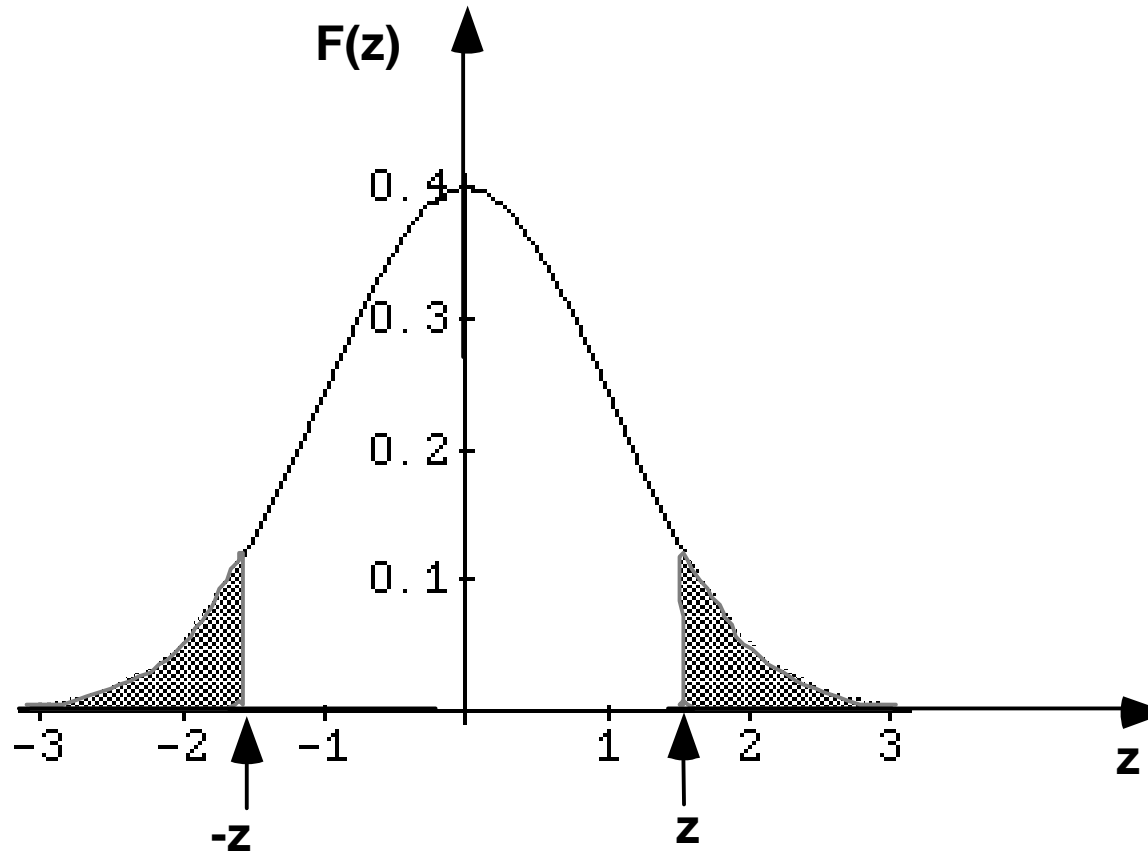
A questo punto possiamo consultare le tabelle della distribuzione normale standardizzata e trovare il valore della probabilità cumulativa corrispondente ad un punto z di 1.5. Questo valore è $\approx .9332$.

Esempio. Consideriamo una distribuzione normale con media 107 e deviazione standard uguale a 70.

Si trovi la probabilità di estrarre a caso dalla distribuzione un'osservazione con un valore minore di 100.

$$z = \frac{100 - 107}{70} = -.1$$

Come possiamo usare le tabelle nel caso di un valore z negativo, dato che le tabelle riportano soltanto valori positivi di z ?



la distribuzione normale è simmetrica e dunque la probabilità cumulativa di $-z$ è uguale alla probabilità cumulativa di $1 - z$:

$$F(-z) = 1 - F(z)$$

Le tabelle ci dicono che la probabilità cumulativa per $z = .1$ è $\approx .5398$.

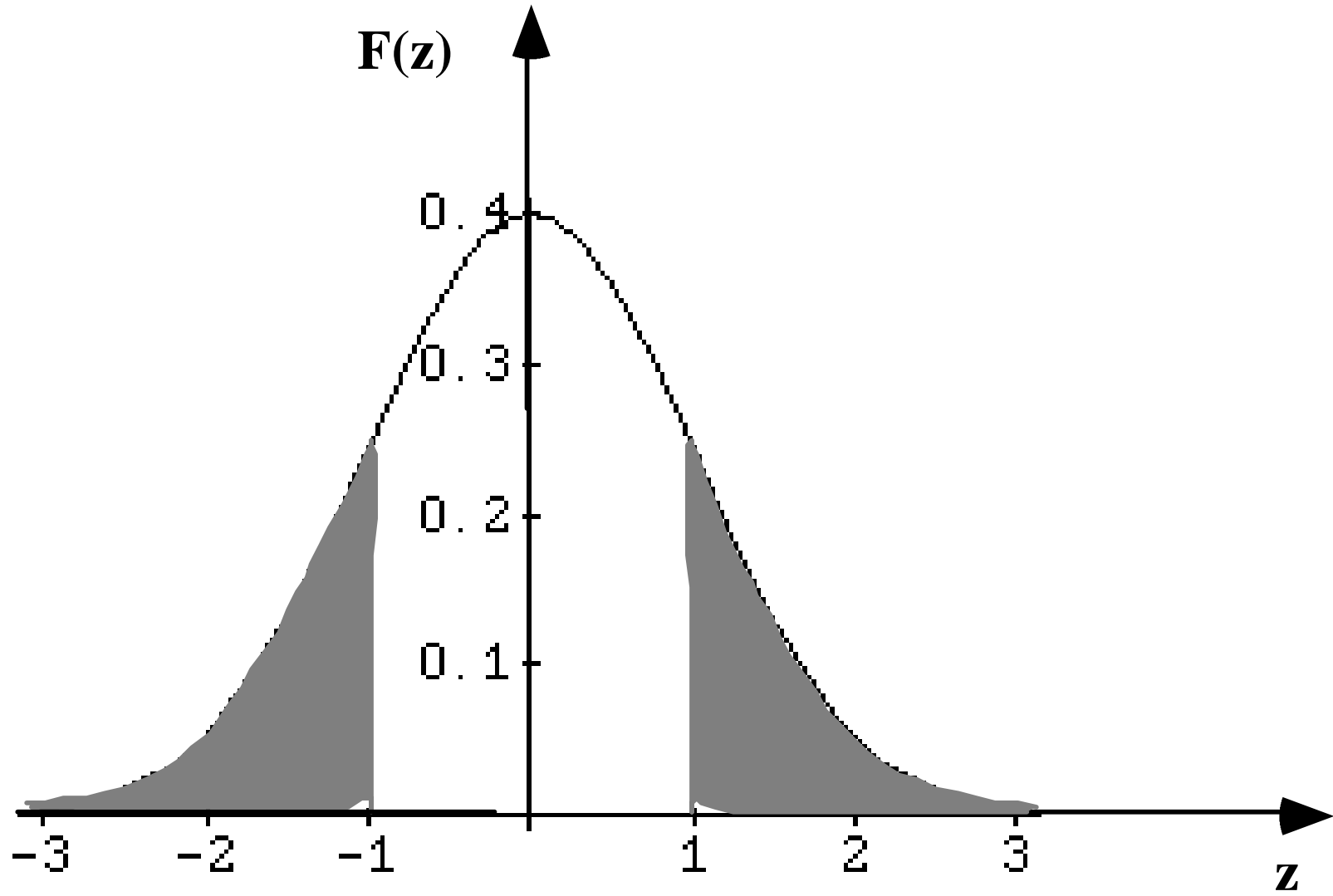
Questo valore corrisponde all'area sottesa alla curva nell'intervallo tra $-\infty$ e $.1$.

Il valore che cerchiamo, però, è quello corrispondente all'area sottesa alla curva nell'intervallo tra $-\infty$ e $-.1$.

In base alla relazione $F(-z) = 1 - F(z)$, valore cercato sarà uguale a $1 - .5398 \approx .4602$.

Esempio. Si trovi la probabilità che un valore estratto a caso da una distribuzione normale abbia un valore compreso tra più e meno una deviazione standard dalla media.

Ovvero, quale è l'area sottesa alla curva normale tra più e meno una deviazione standard dalla media?



La probabilità associata all'intervallo tra $-\infty$ e $+1$:

$$F(1) \approx .8413.$$

La probabilità associata all'intervallo $+1$ e $+\infty$ è

$$1 - F(1) = 1 - .8413 \approx .1587.$$

Possiamo ora calcolare l'area della curva nell'intervallo tra -1 e $+1$. Dall'area totale sottraiamo le due aree ombreggiate:

$$p(-1 \leq z \leq 1) = 1 - 2(.1587) \approx .6826.$$

Questo significa che circa il 68% dei casi in una distribuzione normale standardizzata si trova nell'intervallo compreso tra più e meno uno scarto quadratico medio dalla media.

In altre parole, supponiamo di estrarre a caso un osservazione da una distribuzione normale standardizzata.

Supponiamo inoltre di ripetere questa operazione tantissime volte.

Contiamo quante volte l'osservazione che abbiamo estratto ha un valore compreso tra ± 1 .

Se questo esperimento venisse veramente eseguito si troverebbe che l'osservazione estratta a caso ha un valore compreso tra ± 1 in circa il 68% dei casi.

ESERCIZI

E5 Sia Z una variabile aleatoria distribuita normalmente con media uguale a 0 e varianza uguale a 1.

Si trovi

$$P(Z > 2)$$

$$P(-2 \leq Z \leq 2)$$

$$P(0 \leq Z \leq 1.73)$$

E6 I punteggi ottenuti nel test XYZ sono distribuiti normalmente con media 75 e deviazione standard 10. Quale è la proporzione di punteggi compresi tra 80 e 90?

Distribuzione χ^2

Sia Y_1, Y_2, \dots, Y_n un campione casuale di n osservazioni indipendenti estratte da una popolazione normale con media μ e varianza σ^2 .

Standardizzando queste osservazioni si ottengono n variabili aleatorie indipendenti $Z_i = (Y_i - \bar{m})/s$ normalmente distribuite con media uguale a zero e varianza unitaria.

Può essere dimostrato che

$$X^2 = Z_1^2 + Z_2^2 + \dots + Z_n^2$$

possiede una distribuzione χ^2 con n gradi di libertà.

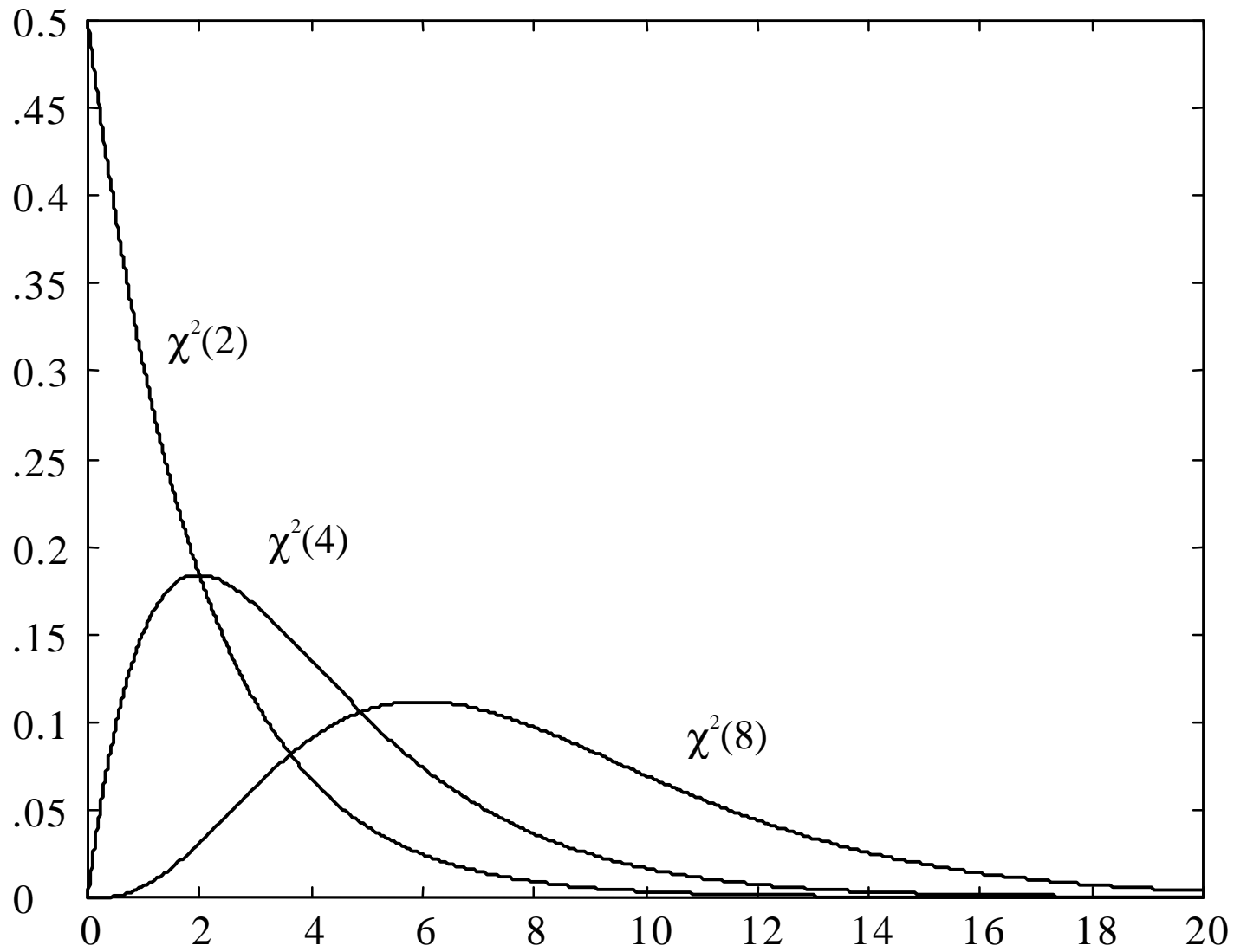
La forma della distribuzione χ^2 dipende dal numero di gradi di libertà e per ogni valore ν si avrà una diversa distribuzione di probabilità

Il parametro ν corrisponde al numero di variabili *indipendenti* sulla base delle quali la statistica $\sum z_i^2$ viene calcolata.

Il suo valore è dato dalla differenza tra il numero n di variabili z_i e il numero dei vincoli che intercorrono tra esse.

Se le variabili z_i sono tutte tra loro indipendenti, $\nu = n$.

Se $\nu \rightarrow \infty$, la distribuzione χ^2 tende alla distribuzione normale standardizzata.



Il valore atteso e la varianza di una variabile χ^2
sono $E(X^2)=n$ e $V(X^2)=2n$.

Approssimazione normale alla distribuzione C^2

La distribuzione \mathbf{c}^2 si approssima alla distribuzione normale con il crescere dei gradi di libertà.

Per campioni di grandi dimensioni, quindi,

$$z = \frac{\mathbf{c}^2 - \mathbf{n}}{\sqrt{2\mathbf{n}}}$$

ovvero

$$\mathbf{c}^2 = z\sqrt{2\mathbf{n}} + \mathbf{n}$$

Esempio 10.5 Usando l'approssimazione normale, si trovi il valore critico in corrispondenza della coda di destra della distribuzione \mathbf{C}^2 ponendo $\alpha = .05$ e $\nu = 100$.

In base alle tavole, il valore cercato è 124.342.

Usando l'approssimazione alla normale si trova:

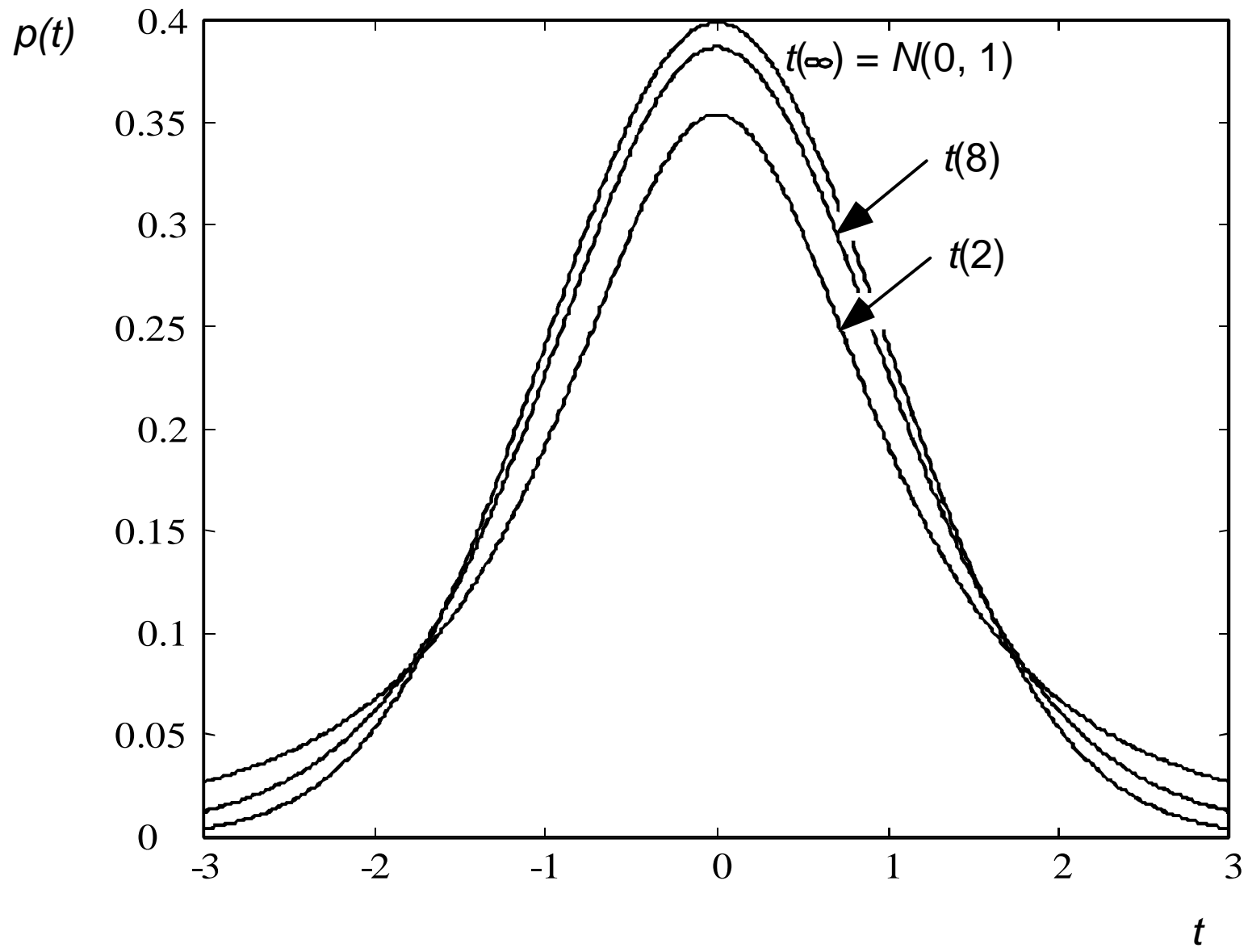
$$\mathbf{c}^2 = z\sqrt{2\mathbf{n}} + \mathbf{n} = 1,65\sqrt{2 \cdot 100} + 100 = 123.334$$

DISTRIBUZIONE t DI STUDENT

Si può dimostrare che la statistica

$$t = \frac{\bar{Y} - m}{s / \sqrt{n}}$$

segue la distribuzione t con $(n - 1)$ gradi di libertà.



$$E(t) = 0 \text{ e } V(t) = \frac{n}{n-2}.$$

Una variabile aleatoria che segue la distribuzione t di Student ha quindi lo stesso valore atteso di una variabile normale standardizzata ma una varianza sempre maggiore di 1.

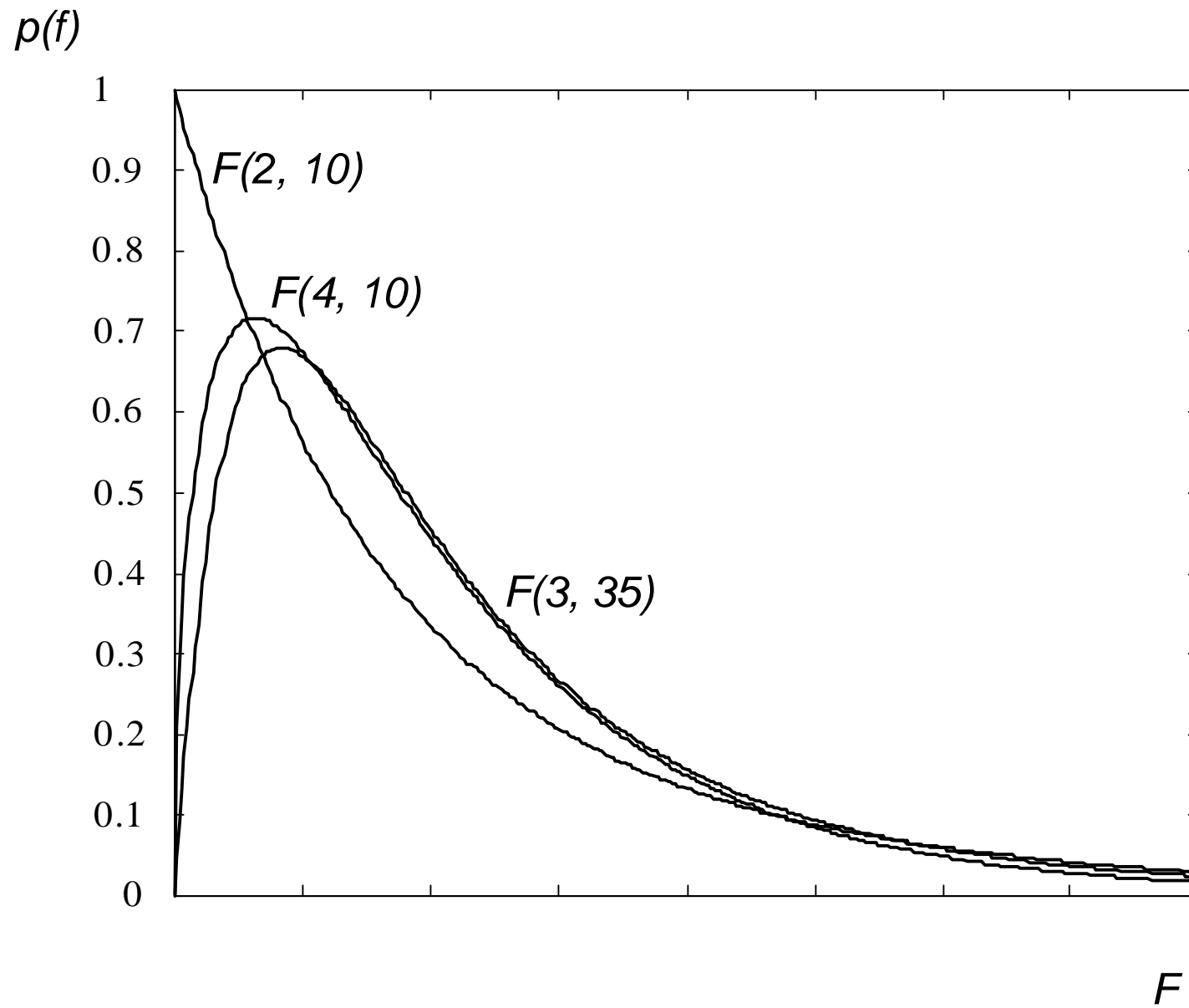
Al crescere dei gradi di libertà, la distribuzione t si approssima sempre più alla distribuzione normale:
 $t(\infty) = N(0,1)$.

DISTRIBUZIONE *F*

Se W_1 e W_2 sono due variabili aleatorie indipendenti distribuite come \mathbf{C}^2 con \mathbf{n}_1 e \mathbf{n}_2 gradi di libertà, può essere dimostrato che il rapporto

$$F = \frac{W_1/\mathbf{n}_1}{W_2/\mathbf{n}_2}$$

segue la distribuzione F con \mathbf{n}_1 gradi di libertà al numeratore e \mathbf{n}_2 gradi di libertà al denominatore.



$$E(F) = \mathbf{n}_2 / (\mathbf{n}_2 - 2)$$

$$V(F) = \frac{2\mathbf{n}_2^2 (\mathbf{n}_1 + \mathbf{n}_2 - 2)}{\mathbf{n}_1 (\mathbf{n}_2 - 2)^2 (\mathbf{n}_2 - 4)}$$

La distribuzione F comprende come sottocasi tutte le distribuzioni esaminate in precedenza.

Se $\nu_1 = 1$ e $\nu_2 \rightarrow \infty$, la distribuzione F tende alla distribuzione normale standardizzata.

Se $\nu_1 = 1$ e $\nu_2 = \nu$, la distribuzione F è uguale al quadrato della distribuzione t di Student: $F_{(1,\nu)} = t_{(\nu)}^2$.

$$t_{(\nu)}^2 = \frac{z^2}{\mathbf{c}_{(\nu)}^2/\mathbf{n}} = \frac{\mathbf{c}_{(1)}^2/1}{\mathbf{c}_{(\nu)}^2/\mathbf{n}} = F_{(1,\nu)}$$

Se $\nu_1 = \nu$ e $\nu_2 \rightarrow \infty$, la distribuzione F tende alla distribuzione χ^2 .