

# **MAE-212: Introdução à Probabilidade e à Estatística II**

## **Aula 03**

### **Inferência Estatística**

# Distribuição amostral de uma Proporção

$X = 1$ , se o indivíduo for portador da característica;  
 $= 0$ , caso contrário;

É fácil ver que  $\mu_x = E(X) = p$  ,  
 $\sigma_x^2 = \text{Var}(X) = p(1-p)$

Retirada de uma AAS(amostragem aleatória simples) dessa população, e indicando por  $Y_n$  o total de indivíduos portadores da característica na amostra, então

$$Y_n = X_1 + X_2 + \dots + X_n,$$

$$Y_n \sim b(n, p)$$

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right).$$

$$\hat{p} = \frac{Y_n}{n}.$$

**Exemplo 10.12** Suponha que  $p = 30\%$  dos estudantes de uma escola sejam mulheres. Colhemos uma AAS de  $n = 10$  estudantes e calculamos  $\hat{p}$  = proporção de mulheres na amostra. Qual a probabilidade de que  $\hat{p}$  difira de  $p$  em menos de 0,01? Temos que essa probabilidade é dada por

$$P(|\hat{p} - p| < 0,01) = P(-0,01 < \hat{p} - p < 0,01).$$

Mas,  $\hat{p} - p \sim N\left(0, \frac{p(1-p)}{n}\right)$  e como  $p = 0,3$ , temos que

$$\text{Var}(\hat{p}) = (0,3)(0,7)/10 = 0,021,$$

e, portanto, a probabilidade pedida é igual a

$$P\left(\frac{-0,01}{\sqrt{0,021}} < Z < \frac{0,01}{\sqrt{0,021}}\right) = P(-0,07 < Z < 0,07) = 0,056.$$

# Outras Distribuições Amostrais

**Exemplo 10.13.** Na Tabela 10.6 apresentamos a distribuição de três outras estatísticas; a variância da amostra,

$$S^2 = \frac{1}{(n - 1)} \sum_{i=1}^n (X_i - \bar{X})^2,$$

a mediana amostral,  $md$ , e o estimador

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2,$$

**Tabela 10.6:** Distribuição amostral de algumas estatísticas obtidas de amostra de tamanho  $n = 3$ , retiradas da população  $\{1, 3, 5, 5, 7\}$  ( $\mu = 4,2$ ,  $\sigma^2 = 4,16$  e  $Md = 5$ ).

Tipo de amostra	Frequência (prob. $\times$ 125)	Soma	Soma dos quadrados	Média $\bar{x}$	Mediana $md$	Variância	
						$s^2$	$\hat{\sigma}^2$
111	1	3	3	1,00	1	0	0
113	3	5	11	1,67	1	4/3	8/9
115	6	7	27	2,33	1	16/3	32/9
117	3	9	51	3,00	1	12	8
133	3	7	19	2,33	3	4/3	8/9
135	12	9	35	3,00	3	4	8/3
137	6	11	59	3,67	3	28/3	56/9
155	12	11	51	3,67	5	16/3	32/9
157	12	13	75	4,33	5	28/3	56/9
177	3	15	99	5,00	7	12	8
333	1	9	27	3,00	3	0	0
335	6	11	43	3,67	3	4/3	8/9
337	3	13	67	4,33	3	16/3	32/9
355	12	13	59	4,33	5	4/3	8/9
357	12	15	83	5,00	5	4	8/3
377	3	17	107	5,67	7	16/3	32/9
555	8	15	75	5,00	5	0	0
557	12	17	99	5,67	5	4/3	8/9
577	6	19	123	6,33	7	4/3	8/9
777	1	21	147	7,00	7	0	0
Total	125						

**Tabela 10.7:** Distribuição amostral da variância  $S^2$ , para amostras de tamanho 3, retiradas da população  $\{1, 3, 5, 5, 7\}$ .

$s^2$	0,00	1,33	4,00	5,33	9,33	12,00
$P(S^2 = s^2)$	11/125	42/125	24/125	24/125	18/125	6/125

$$E(S^2) = 4,16, \quad \text{Var}(S^2) = 11,28.$$

**Tabela 10.8:** Distribuição amostral da mediana da amostra  $md$  para amostras de tamanho 3, retiradas da população  $\{1, 3, 5, 5, 7\}$ .

$md$	1	3	5	7
Prob.	13/125	31/125	68/125	13/125

$$E(md) = 4,30, \quad \text{Var}(md) = 2,54.$$

**Tabela 10.9:** Distribuição amostral da variância  $\hat{\sigma}^2$ , para amostras de tamanho 3, retiradas da população  $\{1, 3, 5, 5, 7\}$ .

$\hat{\sigma}^2$	0,00	0,89	2,67	3,56	6,22	8,00
Prob.	11/125	42/125	24/125	24/125	18/125	6/125

$$E(\hat{\sigma}^2) = 2,77, \quad \text{Var}(\hat{\sigma}^2) = 5,04.$$

Figura 10.6: Distribuição amostral de  $S^2$  para amostras de tamanho  $n = 3$  extraídas de  $\{1, 3, 5, 5, 7\}$ .

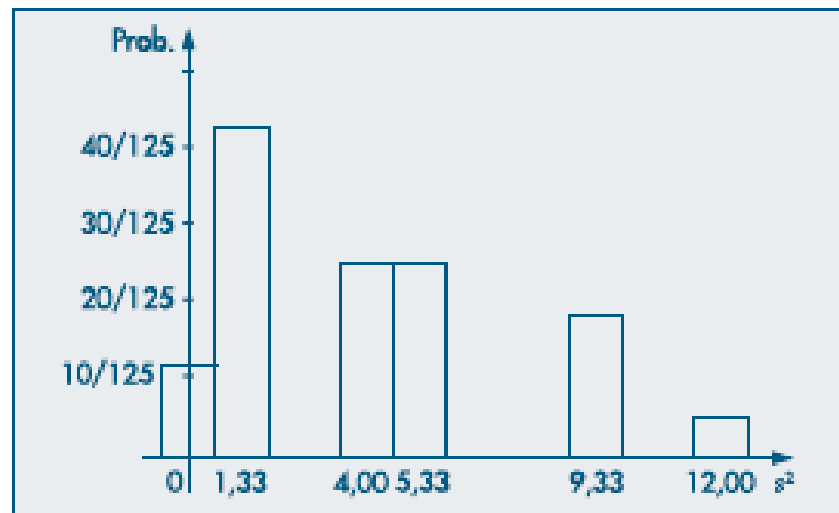


Figura 10.7: Distribuição amostral de md para amostras de tamanho  $n = 3$  de  $\{1, 3, 5, 5, 7\}$ .

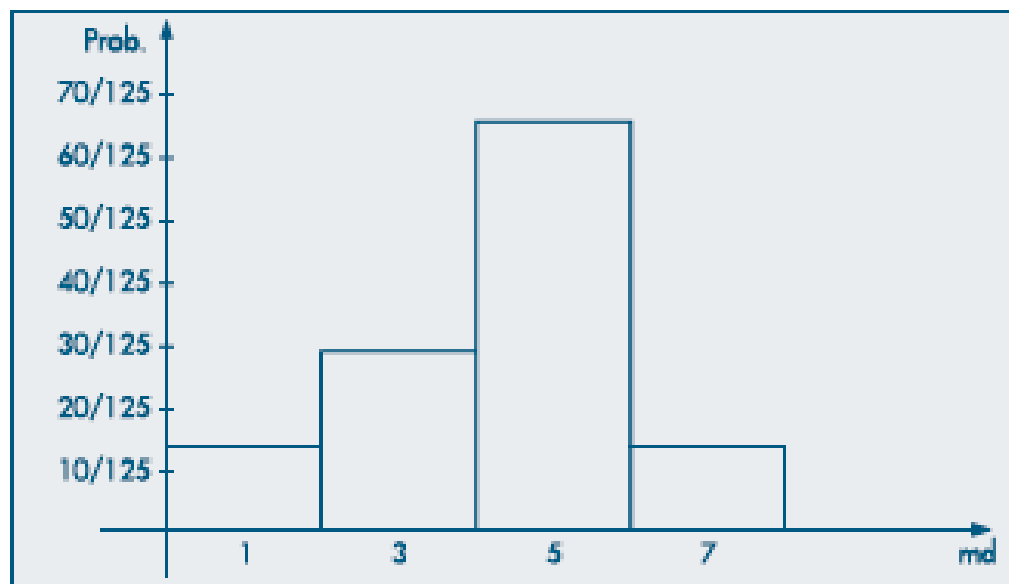
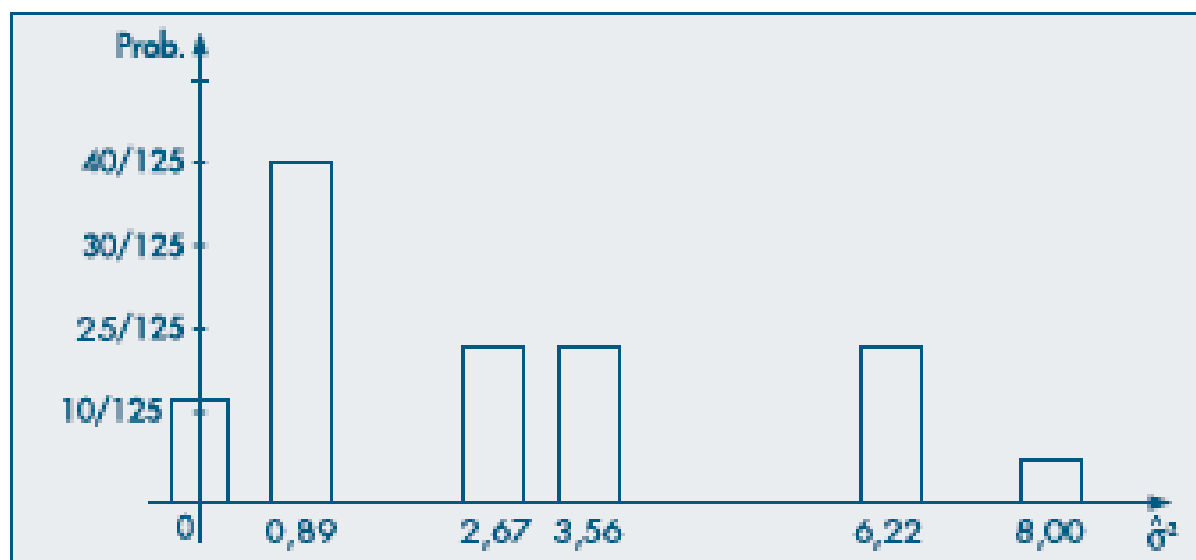


Figura 10.8: Distribuição amostral de  $\hat{\sigma}^2$  para amostras de tamanho  $n = 3$  extraídas de  $\{1, 3, 5, 5, 7\}$ .



Por exemplo, note que  $E(S^2) = 4,16 = \sigma^2$ , logo  $S^2$  satisfaz uma propriedade análoga a  $E(\bar{X}) = \mu$ ; dizemos que  $\bar{X}$  e  $S^2$  são estimadores *não-viesados* dos respectivos parâmetros  $\mu$  e  $\sigma^2$ . Esta propriedade já não vale para  $md$  e  $\hat{\sigma}^2$ , pois  $E(md) = 4,3$ , enquanto  $Md = 5,0$  e  $E(\hat{\sigma}^2) = 2,77$  e não 4,16. Vemos que  $\hat{\sigma}^2$  sub-estima a verdadeira variância.



# **Dimensionamento da amostra**

Seja  $P(\varepsilon) = \gamma$ , a probabilidade da média amostral  $\bar{X}$  estar a uma distância de, no máximo  $\varepsilon$ , da média populacional  $\mu$  (desconhecida), ou seja,

$$\begin{aligned}\gamma &= P\left(\left|\bar{X} - \mu\right| \leq \varepsilon\right) = P\left(\mu - \varepsilon \leq \bar{X} \leq \mu + \varepsilon\right) \\ &= P\left(\frac{-\varepsilon}{\frac{\sigma}{\sqrt{n}}} \leq \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq \frac{\varepsilon}{\frac{\sigma}{\sqrt{n}}}\right) \approx P\left(\frac{-\varepsilon \sqrt{n}}{\sigma} \leq Z \leq \frac{\varepsilon \sqrt{n}}{\sigma}\right),\end{aligned}$$

sendo  $Z \sim N(0,1)$  .

# Dimensionamento da amostra

A partir da relação  $\varepsilon = z \frac{\sigma}{\sqrt{n}}$ ,

o tamanho da amostra  $n$  é determinado por

$$n = \left( \frac{z}{\varepsilon} \right)^2 \sigma^2,$$

conhecendo-se o desvio padrão  $\sigma$  de  $X$ , o erro  $\varepsilon$  da estimativa e o coeficiente de confiança  $\gamma$  do intervalo, sendo  $z$  tal que

$$\gamma = P(-z \leq Z \leq z) \text{ e } Z \sim N(0,1).$$

## Exemplo 4:

A renda per-capita domiciliar numa certa região tem distribuição normal com desvio padrão  $\sigma = 250$  reais e média  $\mu$  desconhecida. Se desejamos estimar a renda média  $\mu$  com erro  $\varepsilon = 50$  reais e com uma confiança  $\gamma = 95\%$ , quantos domicílios devemos consultar?

$X$  : renda per-capita domiciliar na região

$$X \sim N(\mu; 250^2)$$

$n = ??$  tal que  $\varepsilon = 50$  reais,

$$\gamma = 0,95 \Rightarrow z = 1,96$$

**Então,**

$$\begin{aligned} \mathbf{n} &= \left( \frac{\mathbf{z}}{\varepsilon} \right)^2 \sigma^2 \\ &= \left( \frac{\mathbf{1,96}}{\mathbf{50}} \right)^2 (\mathbf{250})^2 \\ &= \mathbf{96,04} \end{aligned}$$

**Aproximadamente 97 domicílios devem ser consultados.**

## Exemplo 5:

A quantidade de colesterol  $X$  no sangue das alunas de uma universidade segue uma distribuição de probabilidades com desvio padrão  $\sigma = 50$  mg/dl e média  $\mu$  desconhecida. Se desejamos estimar a quantidade média  $\mu$  de colesterol com erro  $\varepsilon = 20$  mg/dl e confiança de 90%, quantas alunas devem realizar o exame de sangue?

$X$ : quantidade de colesterol no sangue das alunas da universidade

$\sigma = 50$  mg/dl

$n = ??$  tal que  $\varepsilon = 20$  mg/dl

$\gamma = 0,90 \Rightarrow z = 1,65$

**Supondo que o tamanho da amostra a ser selecionada é suficientemente grande, pelo Teorema do Limite Central temos:**

$$\begin{aligned} \mathbf{n} &= \left( \frac{\mathbf{z}}{\varepsilon} \right)^2 \sigma^2 \\ &= \left( \frac{1,65}{20} \right)^2 (50)^2 \\ &= 17,02 \end{aligned}$$

**Assim, aproximadamente 18 alunas devem realizar o exame de sangue.**

**No caso de proporção, usando a aproximação normal para proporção amostral, temos**

$$n = \left( \frac{z}{\varepsilon} \right)^2 p(1-p)$$

**Como não conhecemos  $p$ , podemos usar o fato de que  $p(1-p) \leq 1/4$ , para todos  $p$ .**



# Dimensionamento da amostra

Da relação  $\varepsilon = z \sqrt{\frac{p(1-p)}{n}}$ ,

segue que o **tamanho amostral**  $n$ , dados  $\gamma$  e a margem de erro  $\varepsilon$ , tem a forma

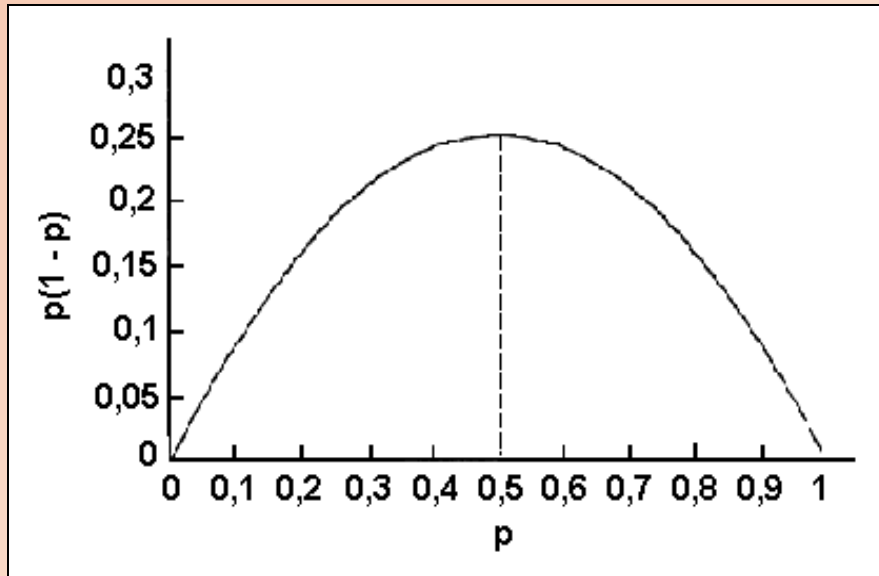
$$n = \left( \frac{z}{\varepsilon} \right)^2 p(1-p),$$

onde  $z$  é tal que  $\gamma = P(-z \leq Z \leq z)$  e  $Z \sim N(0,1)$ .

Entretanto, nesta expressão,  $n$  depende de  $p(1-p)$ , que é desconhecido.

- **Como calcular o valor de  $n$ ?**

## Gráfico da função $p(1-p)$ , para $0 \leq p \leq 1$ .



Pela figura observamos que:

- a função  $p(1-p)$  é uma parábola simétrica em torno de  $p = 0,5$ ;
- o máximo de  $p(1-p)$  é 0,25, alcançado quando  $p = 0,5$ .

Assim, na prática, substituímos  $p(1-p)$  por seu valor máximo, obtendo

$$n = \left( \frac{z}{\varepsilon} \right)^2 0,25 ,$$

que pode fornecer um valor de  $n$  maior do que o necessário.

## Exemplo 6:

No exemplo da USP (Exemplo 1) suponha que nenhuma amostra foi coletada. Quantos estudantes precisamos consultar de modo que a estimativa pontual esteja, no máximo, a 0,02 da proporção verdadeira  $p$ , com uma probabilidade de 0,95?

**Dados do problema:**

$\varepsilon = 0,02$  (erro da estimativa);

$P(\varepsilon) = \gamma = 0,95 \Rightarrow z = 1,96.$

$$n = \left( \frac{1,96}{0,02} \right)^2 p(1-p) \leq \left( \frac{1,96}{0,02} \right)^2 0,25 = 2401 \text{ estudantes.}$$

**Pergunta:** *É possível reduzir o tamanho da amostra quando temos alguma informação a respeito de  $p$ ?*

**Por exemplo, sabemos que:**

- $p$  não é superior a 0,30, ou
- $p$  é pelo menos 0,80, ou
- $p$  está entre 0,30 e 0,60.

**Resposta:** *Depende do tipo de informação sobre  $p$ .*

**Em alguns casos, podemos substituir a informação  $p(1-p)$ , que aparece na expressão de  $n$ , por um valor menor que 0,25.**

# Redução do tamanho da amostra

Vimos que, se nada sabemos sobre o valor de  $p$ , no cálculo de  $n$ , substituímos  $p(1-p)$  por seu valor máximo, e calculamos

$$n = \left( \frac{z}{\varepsilon} \right)^2 0,25 .$$

Se temos a informação de que  $p$  é *no máximo*  $0,30$  ( $p \leq 0,30$ ), então o valor máximo de  $p(1-p)$  será dado por  $0,3 \times 0,7 = 0,21$ .

Logo, reduzimos o valor de  $n$  para

$$n = \left( \frac{z}{\varepsilon} \right)^2 0,21 .$$

**Agora, se  $p$  é pelo menos 0,80 ( $p \geq 0,80$ ), então o máximo de  $p(1-p)$  é  $0,8 \times 0,2 = 0,16$  e temos**

$$n = \left( \frac{z}{\varepsilon} \right)^2 0,16 .$$

**Mas, se  $0,30 \leq p \leq 0,60$ , o máximo de  $p(1-p)$  é  $0,5 \times 0,5 = 0,25$  e, neste caso, não há redução, ou seja,**

$$n = \left( \frac{z}{\varepsilon} \right)^2 0,25 .$$

## Exemplo 7:

No Exemplo 6, suponha que temos a informação de que no máximo 30% dos alunos da USP foram ao teatro no último mês. Portanto, temos que  $p \leq 0,30$  e, como vimos, o máximo de  $p(1-p)$  neste caso é 0,21.

Assim, precisamos amostrar

$$n = \left( \frac{z}{\varepsilon} \right)^2 0,21 = \left( \frac{1,96}{0,02} \right)^2 0,21 = 2017 \text{ estudantes,}$$

conseguindo uma redução de  $2401 - 2017 = 384$  estudantes.





# Distribuição Normal : Valores de $P( Z \leq z ) = A(z)$

## Segunda decimal de z

	0	1	2	3	4	5	6	7	8	9
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990
3.1	0.9990	0.9991	0.9991	0.9991	0.9992	0.9992	0.9992	0.9992	0.9993	0.9993
3.2	0.9993	0.9993	0.9994	0.9994	0.9994	0.9994	0.9994	0.9995	0.9995	0.9995
3.3	0.9995	0.9995	0.9995	0.9996	0.9996	0.9996	0.9996	0.9996	0.9996	0.9997
3.4	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9998
3.5	0.9998	0.9998	0.9998	0.9998	0.9998	0.9998	0.9998	0.9998	0.9998	0.9998
3.6	0.9998	0.9998	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999
3.7	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999
3.8	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999	0.9999
3.9	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000

Parte inteira e primeira decimal de z