

# ANÁLISE DE ITENS DE UMA PROVA DE RACIOCÍNIO ESTATÍSTICO<sup>1</sup>

Claudette Maria Medeiros Vendramini<sup>\*</sup>  
Marjorie Cristina da Silva<sup>#</sup>  
Michelle Canale<sup>¶</sup>

**RESUMO.** Este estudo objetivou analisar as 18 questões (do tipo múltipla escolha) de uma prova sobre conceitos básicos de Estatística pelas teorias clássica e moderna. Participaram 325 universitários, selecionados aleatoriamente das áreas de humanas, exatas e saúde. A análise indicou que a prova é predominantemente unidimensional e que os itens podem ser mais bem ajustados ao modelo de três parâmetros. Os índices de dificuldade, discriminação e correlação biserial apresentam valores aceitáveis. Sugere-se a inclusão de novos itens na prova, que busquem confiabilidade e validade para o contexto educacional e revelem o raciocínio estatístico de universitários ao ler representações de dados estatísticos.

**Palavras-chave:** teoria de resposta ao item, psicometria clássica, educação estatística.

## ANALYSIS OF ITEMS OF A STATISTICAL REASONING TEST

**ABSTRACT.** This study aimed at to analyze the 18 questions (of multiple choice type) of a test on basic concepts of Statistics for the classic and modern theories. The test was taken by 325 undergraduate students, randomly selected from the areas of Human, Exact and Health Sciences. The analysis indicated that the test has predominantly one dimension and that the items can be better fitting to the model of three parameters. The indexes of difficulty, discrimination and biserial correlation present acceptable values. It is suggested to include new items to the test in order to obtain reliability and validity to use it in the education context and to reveal the statistical reasoning of undergraduate students when dealing with statistical data representation.

**Key words:** item response theory, classical psychometry, statistical education.

O ensino-aprendizagem de Estatística vem ocupando um lugar importante nas instituições de ensino superior, devido à necessidade de formar profissionais capacitados e com domínio de técnicas de análise de dados que fundamentem a tomada de decisões baseada na inferência de dados amostrais. Muitas vezes, o profissional necessita lidar com grande quantidade de informações e processá-las em tempo mínimo.

A tomada de decisões baseadas em dados estatísticos exige que os profissionais raciocinem estatisticamente. O raciocínio ou pensamento estatístico é definido por Garfield e Gal (1999) como a maneira como as pessoas raciocinam com as idéias estatísticas

para dar sentido às informações recebidas. Esse raciocínio envolve interpretações baseadas em conjuntos de dados, representações gráficas e resumos estatísticos, muitas vezes combinando dados e probabilidade para fazer inferências e interpretar os resultados estatísticos. Subjacente a esse raciocínio está a compreensão conceitual de distribuição, centro, dispersão, associação, incerteza, aleatoriedade e amostragem (Garfield, 2002).

Um modelo proposto por Lalonde e Gardner (1993) para prever o desempenho em Estatística de alunos de Psicologia considera que a aprendizagem de Estatística é como a aprendizagem de uma segunda linguagem, onde a ansiedade matemática, a atitude e a

<sup>1</sup> Este artigo foi elaborado com base nos resultados de pesquisa realizada e apoiada financeiramente pela Universidade São Francisco - USF e pelo Programa Institucional de Bolsas de Iniciação Científica do CNPq - PIBIC/CNPq.

<sup>\*</sup> Estatística, Doutora em Educação pela Universidade Estadual de Campinas e Docente da graduação e do Programa de Estudos Pós-Graduados em Psicologia da USF.

<sup>#</sup> Aluna do curso de Psicologia da USF e bolsista da iniciação científica PIBIC/CNPq.

<sup>¶</sup> Aluna do curso de Matemática da USF e bolsista da iniciação científica PIBIC/CNPq.

motivação em relação à Estatística e à habilidade matemática são os componentes essenciais para entender e prever o desempenho em Estatística. Em um estudo desenvolvido por Vendramini (2000) com universitários, constatou-se que apenas 24,5% dos participantes conseguiam definir Estatística. Todavia, 80,3% declaravam ter um motivo para estudá-la e 90,0% a consideravam uma ferramenta útil. Os resultados indicaram associações positivas e significativamente diferentes de zero entre o desempenho acadêmico dos participantes na disciplina Estatística e o desempenho em uma prova específica da disciplina.

Dentro deste contexto, tanto as *atitudes* quanto o *raciocínio* estatístico são importantes, sendo necessários testes e provas de qualidade para se fazerem inferências sobre esses construtos. Na avaliação psicológica ou educacional, os testes e provas são instrumentos de medida úteis e auxiliam nas tomadas de decisão dos profissionais dessas áreas. Várias formas e procedimentos podem ser utilizados nesse processo e, quando consideradas as medidas objetivas (testes), pode-se fundamentar esta avaliação em dois modelos: a Teoria Clássica dos Testes (TCT) e a Teoria de Resposta ao Item (TRI).

Há muitas décadas que a TCT tem sido útil para o desenvolvimento dos testes psicológicos, embora padeça de várias limitações, como, por exemplo, ser dependente do conjunto de itens que compõem o instrumento de medida, limitando-se assim, a sua aplicabilidade (Andrade, Tavares & Valle, 2000). Na TRI, o procedimento de medida utilizado parte da suposição de que existe no sujeito um traço (uma característica individual determinante de como responder aos itens de um teste) que possui uma relação probabilística com cada um dos itens utilizados (Fletcher, 1994). Considerando-se que os parâmetros de cada item não dependem dos outros itens do teste, mas que a pontuação do teste se faz em função das respostas do sujeito a cada item, é possível verificar se os respondentes são mais ou menos hábeis, e da mesma forma, se os itens podem ser considerados mais fáceis ou mais difíceis, já que itens e pessoas são colocados na mesma escala de desempenho.

A TRI não entra em contradição com os princípios da psicometria clássica e traz uma nova proposta estatística, a de análise centrada nos itens, que supera as limitações da teoria clássica, além de apresentar novos recursos tecnológicos para a avaliação psicológica e educacional (Primi, 1998). Existem vários modelos possíveis de resposta ao item, que diferem em sua forma em função da característica

do item e do número de parâmetros especificados no modelo. Todos os modelos dessa teoria possuem um ou mais parâmetros que descrevem o item e um ou mais parâmetros que descrevem o respondente. O primeiro passo para uma aplicação da TRI é a estimação desses parâmetros (Hambleton, Swaminathan & Rogers, 1991).

Alguns pontos têm sido levantados na literatura sobre a adequação desta teoria na área de avaliação educacional. Dois deles considerados importantes são: a dimensionalidade do espaço de traços latentes envolvidos na avaliação e a equalização de diferentes avaliações. É necessário ressaltar que, apesar de não haver dúvidas de que a aplicação desta teoria muito contribui para a melhoria das avaliações educacionais em geral, sua disseminação efetiva depende da integração de especialistas das áreas de estatística e educação (Andrade, 2001).

Uma das vantagens da utilização da TRI na avaliação educacional é que esta possibilita análises qualitativas a partir dos resultados brutos de uma prova, fornecendo assim informações mais precisas do desempenho dos respondentes e da qualidade das questões utilizadas (itens), questões que devem ter índices de dificuldade e de discriminação aceitáveis e correlacionadas com a prova total. No Brasil, a TRI vem sendo aplicada em diversas avaliações educacionais desde 1995, como na análise de dados do Sistema Nacional do Ensino Básico - SAEB (INEP, 2002) e do Sistema de Avaliação do Rendimento Escolar do Estado de São Paulo - SARESP (Vendramini, 2002). Nessas avaliações, o objetivo é comparar o desempenho dos alunos em diferentes séries e disciplinas, utilizando questões comuns entre ciclos de aplicação e entre séries.

Convém, no entanto, ressaltar que os modelos de resposta ao item só podem ser considerados vantajosos quando o ajuste do modelo aos dados de interesse for satisfatório. Um modelo mal-ajustado não fornecerá parâmetros invariantes para os itens e para as habilidades (Hambleton, Swaminathan & Rogers, 1991).

Nesse sentido, o presente estudo objetivou analisar as questões de uma prova de matemática e estatística para universitários pela TCT e pela TRI, centrando a análise nas questões da prova e não na prova como um todo, e assim contribuir para o melhor entendimento do desempenho dos alunos quando submetidos a esse tipo de prova. Para a apresentação dos modelos matemáticos e discussão dos resultados serão utilizados os termos *itens* e *teste*, correspondentes aos termos *questões* e *prova*.

## MÉTODOS

Foram analisados os dados de pesquisa de Vendramini (2000), relativos à aplicação de uma prova de estatística em alunos regularmente matriculados no primeiro ano de uma universidade particular, que cursavam a disciplina Estatística nos anos de 1996 a 1998.

### Participantes

De um total de 29 cursos de uma universidade particular do interior do Estado de São Paulo, foram selecionados, por conveniência, dois da área de Ciências Humanas (Psicologia, Administração), dois de Ciências Exatas (Ciência da Computação e Engenharia Mecânica-Automação e Sistemas - Mecatrônica) e dois de Ciências da Saúde (Farmácia e Medicina). Participaram da pesquisa 325 alunos com idades variando de 18 a 35 anos (*Média*=21,7 anos; *DP*=3,7 anos), 56,7% do gênero feminino.

### Material e procedimento

A prova é composta por 18 questões de múltipla escolha sobre conceitos básicos de Estatística, sendo seis referentes a dados apresentados em gráficos, seis envolvendo conceitos e operações matemáticas necessárias para a resolução de problemas estatísticos (fração, porcentagem, número decimal, potenciação, arredondamento) e seis referentes a dados apresentados em uma tabela estatística. Essa prova foi aplicada em sala de aula, por auxiliares de pesquisa, após consentimento dos coordenadores de curso e a concordância dos alunos em participar da pesquisa.

### Análise de dados

Tanto para a análise da TCT quanto para a TRI foram consideradas as respostas dos participantes nas cinco alternativas de cada item, indicadas as respostas corretas, para possibilitar que os itens assumam escores do tipo certo/errado (itens dicotômicos). Embora as alternativas dos itens deste estudo não tenham sido construídas com o objetivo de evidenciar raciocínios falhos dos estudantes, a análise dos índices de discriminação e das correlações ponto biserial (item total) revelou tendências de escolha da opção errada em alguns dos itens, fornecendo indicadores de raciocínio errado dos estudantes, o que pode auxiliar o professor no processo de ensino e aprendizagem de estatística. Os dados foram analisados por programas computacionais específicos: o programa ITEMAN para a análise convencional de itens (*Assessment System Corporation*, 1998); o programa TESTFACT

para a análise da dimensionalidade da prova (Wilson, Wood & Gibbons, 1998); o programa RASCAL para análise do ajuste do modelo de um parâmetro de *Rasch* (*Assessment System Corporation*, 1995a); e o XCALIBRE para o ajuste dos modelos de dois e três parâmetros pela estimação marginal de máxima verossimilhança (*Assessment System Corporation*, 1995b).

Será apresentada inicialmente a análise pela TCT para, a seguir, serem apresentados e aplicados os três modelos matemáticos da TRI (um, dois e três parâmetros), selecionado-se aquele que se ajustar a um maior número de itens e que melhor represente as respostas dos sujeitos à prova.

## RESULTADOS DA ANÁLISE DE ITENS PELA TEORIA CLÁSSICA DOS TESTES

A análise clássica dos itens de uma prova baseia-se em parâmetros descritivos dos itens, que auxiliam na interpretação da distribuição das respostas em cada alternativa do item. As propriedades psicométricas dos itens da prova (Tabela 1) correspondem aos seguintes parâmetros: (1) índice de facilidade - proporção de participantes que responderam ao item corretamente; (2) índice de discriminação - que mede a capacidade do item de diferenciar os participantes de maior habilidade (27% dos respondentes com pontuações mais altas) daqueles de menor habilidade (27% dos respondentes com pontuações mais baixas) e corresponde à diferença entre a proporção de acertos do primeiro grupo e a do segundo grupo; (3) correlação entre ponto biserial item-total e entre a resposta correta no item e na alternativa e a pontuação total na prova; (4) a média e o desvio-padrão do número total de acertos, considerando-se as respostas deixadas em branco como erradas; (5) a média do total de acertos dos participantes que acertaram um determinado item.

Os resultados revelaram que a prova apresenta itens fáceis (com valores acima de 70% de acertos - itens 1, 2, 3 e 10) e itens difíceis (com valores inferiores a 30% de acertos - itens 7, 8, 13, 15, 16, 17 e 18). A maioria dos itens referentes à interpretação de dados apresentados em tabelas foi muito difícil, indicando a necessidade de uma análise cuidadosa do significado desses resultados. Segundo depoimento de alguns respondentes, a tabela apresentada na prova não estava clara, dificultando assim a solução dos problemas referentes a seus dados. Essas dificuldades estavam relacionadas a: apresentação dos dados numéricos em duas colunas separadas pelas categorias das variáveis; total apresentado somente no título;

apresentação de várias variáveis em uma mesma tabela (partido, sexo, idade, religião e profissão). Sugere-se que, para elaboração de novas provas, também sejam apresentadas algumas tabelas com número pequeno de variáveis, com categorias de resposta apresentadas na primeira coluna e com total na última linha da tabela, e não no título, de maneira que se possibilite ao pesquisador verificar quanto a forma de apresentação da tabela interfere na resposta do sujeito.

Os índices de discriminação e as correlações ponto bisserial (item total) revelaram uma tendência de escolha da opção errada pelos participantes que obtiveram os escores mais altos no teste como um todo nos itens 7, 8, 13, 15, 16, 17 e 18. Quando consideradas, para estes mesmos itens, as correlações ponto bisserial por alternativa de resposta (não apresentadas na Tabela 1), encontram-se valores positivos para alternativas erradas, e não negativos, como era de se esperar, indicando um padrão de respostas para uma alternativa errada não inverso ao de respostas corretas. No item 16, por exemplo, 47% dos respondentes de maior pontuação total responderam a alternativa (c), quando a correta era a (d), sendo que esse item obteve a menor correlação com o total de itens da prova e uma correlação ponto bisserial positiva para as alternativas c ( $r=0,21$ ) e d ( $r=0,10$ ).

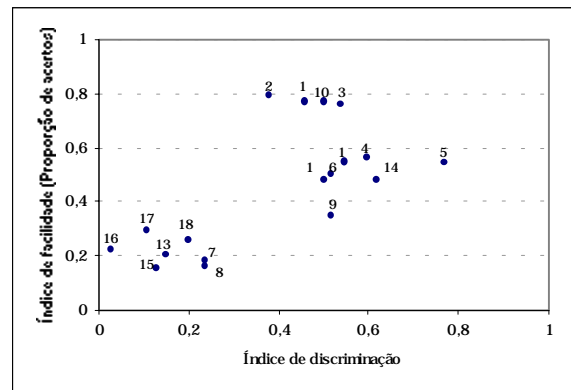
**Tabela 1.** Parâmetros descritivos dos itens da prova ( $N=325$ ).

Item	Média de acertos dos que acertaram o item	Índice de facilidade (Proporção de acertos)	Índice de discriminação total	Correlação item-Ponto-bisserial
1	8,86	0,77	0,46	0,49
2	8,73	0,79	0,38	0,43
3	9,02	0,76	0,54	0,57
4	9,47	0,56	0,60	0,52
5	9,88	0,54	0,77	0,63
6	9,49	0,50	0,52	0,46
7	9,98	0,18	0,24	0,29
8	10,15	0,16	0,24	0,29
9	10,10	0,35	0,52	0,48
10	8,90	0,77	0,50	0,52
11	9,37	0,55	0,55	0,47
12	9,50	0,48	0,50	0,45
13	9,48	0,20	0,15	0,23
14	9,78	0,48	0,62	0,53
15	9,32	0,15	0,13	0,17
16	8,64	0,22	0,03	0,10
17	8,89	0,29	0,11	0,17
18	9,16	0,26	0,20	0,21
<b>Total</b>	<b>8,04</b>	-	-	-
<b>DP</b>	<b>3,17</b>	-	-	-

Os parâmetros psicométricos encontrados atendem razoavelmente aos requisitos necessários,

pelos instrumentos de medida. A média de acertos na prova (8,0) é uma unidade inferior ao ponto médio da escala (9,0), com desvio-padrão igual a 3,2, sendo que a média de acertos na prova para os participantes que acertaram um determinado item variou de 8,64 a 10,15. Os itens foram considerados de dificuldade mediana e consistência interna razoável pela técnica de *Kuder-Richardson<sub>2</sub>* que permite verificar a fidedignidade do teste ( $K-R 20=0,683$ ).

A representação gráfica da dispersão dos itens, segundo os índices de discriminação e de facilidade está apresentada na Figura 1.



**Figura 1.** Gráfico de dispersão dos itens, índice de discriminação por índice de facilidade.

A figura indica a presença de dois grupos de itens: sete itens difíceis com baixa discriminação e 11 itens de facilidade média e alta com discriminação mais elevada. Para todos os itens foi possível discriminar os participantes com pior desempenho (pontuação de 0 a 5) e os participantes com melhor desempenho (pontuação de 11 a 18), revelando qualidade dos itens da prova.

## RESULTADOS DA ANÁLISE PELA TEORIA DE RESPOSTA AO ITEM

### Principais modelos matemáticos

Os modelos matemáticos propostos na literatura variam conforme a natureza da questão ou item (dicotômico ou não dicotômico), o número de populações envolvidas (uma ou mais de uma) e o número de traços latentes que estão sendo medidos (um ou mais de um).

Antes de empregar os modelos matemáticos da TRI, deve-se comprovar o cumprimento de dois pressupostos teóricos fundamentais para a utilização de modelos unidimensionais: (1) critério da unidimensionalidade: os itens de um teste devem medir uma única habilidade, ou

ao menos deve haver um fator dominante que influencie o desempenho dos respondentes no teste; (2) critério da independência local: as respostas dos participantes aos itens não devem ser influenciadas pelas respostas fornecidas a outros itens, para indivíduos com uma mesma habilidade.

Os parâmetros dos itens mais relevantes são: a dificuldade, a discriminação e a probabilidade de acerto por acaso (isto é, a probabilidade de um sujeito de baixa habilidade dar uma resposta correta a um item difícil). A dificuldade do item é dada na mesma escala da habilidade, referente à habilidade necessária para uma dada probabilidade de acertar o item, calculada a partir da probabilidade de acertar o item por acaso. A discriminação corresponde à inclinação da Curva Característica do Item (CCI) e indica quanto indivíduos de diferentes habilidades diferem quanto à probabilidade de acertar um item. A representação gráfica da CCI tem forma de "S", com inclinação e deslocamento na escala de habilidade definidos pelos parâmetros dos itens. No eixo das abscissas está indicado o nível observado de habilidade (traço latente) do indivíduo, designado por  $\theta$ , cujo valor pode variar de  $-\infty$  a  $+\infty$ , e no eixo das ordenadas a probabilidade de responder corretamente ao item, designado por  $P_i(\theta)$  e variando de 0 a 1. Assim, os modelos matemáticos apresentados nas três expressões a seguir, representam para cada item  $i$  a função de probabilidade de indivíduos com habilidade  $\theta$  acertarem esse item.

Na TRI, a probabilidade de resposta correta depende da habilidade do sujeito ( $\theta$ ), que permite expressar a sua resposta aos itens, e dos parâmetros dos itens ( $a_i$ ,  $b_i$  e  $c_i$ ), sendo, portanto, necessário estimar valores destes parâmetros que melhor expliquem os resultados obtidos, com base nas respostas dos participantes aos itens. Esse processo é chamado de calibração ou parametrização, e é feito com o auxílio de programas específicos (RASCAL, XCALIBRE, BILOG, entre outros).

Nas expressões dos modelos matemáticos apresentados a seguir, a probabilidade condicional  $P(X_{ij}=1|\theta_j)$  é denotada apenas por  $P_i(\theta)$ . Para cada item  $i$  a probabilidade de acerto é uma função da habilidade  $\theta$  do indivíduo, que varia de  $-\infty$  a  $+\infty$ . Conhecidos a habilidade  $\theta$  de um indivíduo e os parâmetros do item  $i$  é possível determinar a probabilidade de ele acertar esse item pelas expressões (1), (2) ou (3) de acordo com o modelo escolhido.

### Modelo logístico de três parâmetros

O modelo de três parâmetros é o modelo teórico mais completo. Considera como variáveis

que influenciam a probabilidade do indivíduo acertar o item os três parâmetros citados: a dificuldade, a discriminação e a probabilidade de acerto ao acaso, e é expresso pela função matemática (1) a seguir.

$$P(X_{ij} = 1|\theta_j) = P_i(\theta) = c_i + \frac{(1 - c_i)}{1 + e^{-Da_i(\theta - b_i)}} \quad (1)$$

com  $i = 1, 2, \dots, n$  e  $j = 1, 2, \dots, m$ , sendo:

$X_{ij}$  uma variável dicotômica que assume os valores 1 ou 0, conforme o indivíduo  $j$  responda correta ou incorretamente o item  $i$ , respectivamente;

$\theta_j$  o valor que representa a variável latente (aptidão ou habilidade) que permite explicar a resposta do  $j$ -ésimo indivíduo aos itens;

$P(X_{ij} = 1|\theta_j) = P_i(\theta)$  a probabilidade de um indivíduo  $j$  com habilidade  $\theta_j$  responder corretamente o item  $i$ ;

$c_i$  a probabilidade de acerto ao acaso;

$b_i$  o índice de dificuldade (ou parâmetro de posição) do item  $i$ , medido na mesma escala da habilidade  $\theta$ . Corresponde à habilidade necessária para uma probabilidade de acerto igual a  $(1 + c_i) / 2$ ;

$a_i$  o índice de discriminação (ou parâmetro de inclinação) do item  $i$ , com valor proporcional à inclinação da CCI no ponto  $b_i$ ;

$D$  um fator de escala constante, igual a 1 ou a 1,7 (quando se deseja que a função logística se aproxime da ogiva normal);

$e$  um número transcendental, base dos logaritmos neperianos, cujo valor é aproximadamente 2,718;

$n$  o número de itens;

$m$  o número de indivíduos.

A representação gráfica da função anterior, denominada Curva Característica do Item (CCI), é apresentada posteriormente no tópico "Análise de itens pela Teoria de Resposta ao item".

### Modelo logístico de dois parâmetros

O modelo de dois parâmetros possui em sua expressão o índice de dificuldade  $b_i$  e o de discriminação  $a_i$ , como variáveis que influenciam a probabilidade de o indivíduo acertar o item. Esse modelo pode ser entendido como um modelo de três parâmetros com o valor  $c_i = 0$  (Expressão 2).

$$P(X_{ij} = 1|\theta_j) = P_i(\theta) = \frac{1}{1 + e^{-Da_i(\theta - b_i)}} \quad (2)$$

### Modelo logístico de um parâmetro

O modelo de um parâmetro, também conhecido como modelo de *Rasch*, possui em sua expressão o índice de dificuldade  $b_i$ , que se relaciona com a probabilidade de acertar o item  $i$  por acaso. Esse modelo pode ser entendido como um modelo de três parâmetros com o valor  $c_i=0$  e mesmo valor  $a$  para todos os  $a_i$ 's (Expressão 3).

$$P(X_{ij} = 1 | \theta_j) = P_i(\theta) = \frac{1}{1 + e^{-Da(\theta - b_i)}} \quad (3)$$

### Análise da dimensionalidade da prova

Um das suposições da TRI é que a prova seja unidimensional, ou pelo menos que se possa assumir um fator predominante para se utilizarem modelos unidimensionais. O programa TESTFACT efetua a análise considerando questões do tipo certo/errado (dados dicotômicos) a partir das respostas dos participantes (em vez da matriz de correlação). Esta análise é denominada Análise Fatorial com Informação Completa (*Full Information Factor Analysis*) e inclui progressivamente fatores que indicam a contribuição do fator incluído para a explicação das correlações entre os itens, possibilitando fazer previsões das respostas dos participantes aos itens a partir das curvas dos itens. Diferentes padrões de resposta podem ser esperados quando os itens são completamente independentes, ou quando medem um único fator ou mais de um fator. A adequação de um modelo unidimensional ou multidimensional aos padrões de respostas dos participantes é verificada pelo teste Qui-quadrado. É imprescindível verificar, por este mesmo teste, se a inclusão sucessiva de fatores nos modelos tem um efeito significativo (Bock, Gibbons & Muraki, 1988).

Os resultados da análise fatorial com informação completa, efetuada para os dados em questão, indicaram uma baixa correlação média tetracórica ( $r_{tet}=0,17$ ) entre os 153 pares de combinações, dois a dois, dos itens. Ao considerar o modelo unidimensional, podem ser explicados 23,7% da variância entre os itens. Incluindo-se um segundo fator, modelo bidimensional, a variância explicada aumenta para 29,5%. A magnitude da diferença entre os padrões de resposta observados e os reproduzidos pelo modelo é estatisticamente significativa quando se acrescenta esse segundo fator ( $Qui^2 [17]=52,96$ ;  $p<0,001$ ), indicando que ele representa uma parte significativa das correlações não explicadas pelo primeiro fator.

Ao se considerar um modelo tridimensional, a variância explicada entre os itens aumenta muito pouco, passando para 31,5%. A magnitude da diferença entre os padrões de resposta observados e os reproduzidos pelos modelos, quando se acrescenta um terceiro fator, não é significativa ( $Qui^2 [16]=10,44$ ;

$p=0,843$ ), podendo-se supor que ele não seja necessário. Na Tabela 2 estão apresentadas as cargas fatoriais dos itens nos três fatores considerados para a análise fatorial dos dados.

Dado que alguns itens apresentaram carga fatorial baixa (menor que 0,30) nos três fatores (Tabela 2), esses foram excluídos passo a passo até a obtenção de um conjunto de itens com boa consistência interna pela técnica de *Kuder-Richardson* que permite verificar a fidedignidade do teste ( $K-R 20 = 0,753$ ). Desta forma, os resultados da análise fatorial com informação completa indicaram que os itens 1, 2, 3, 4, 5, 6, 9, 10, 11, 12 e 14 foram agrupados no primeiro fator, sendo que somente o item 14 referia-se aos dados apresentados em tabela estatística, e o restante, a itens referentes aos gráficos de pontos e de colunas e a itens que envolviam apenas cálculos matemáticos (vide prova no Anexo 1).

**Tabela 2.** Cargas não rotacionadas dos fatores principais obtidos na primeira extração dos fatores.

Item	Fator 1	Fator 2	Fator 3	Item	Fator 1	Fator 2	Fator 3
1	<b>0,63</b>	0,00	0,04	10	<b>0,71</b>	-0,03	0,20
2	<b>0,56</b>	-0,20	0,11	11	<b>0,49</b>	0,05	0,05
3	<b>0,79</b>	-0,16	0,00	12	<b>0,48</b>	0,09	0,14
4	<b>0,56</b>	-0,02	-0,24	13	0,05	<b>0,47</b>	-0,40
5	<b>0,78</b>	-0,01	-0,14	14	<b>0,58</b>	0,29	0,19
6	<b>0,57</b>	-0,30	-0,26	15	0,11	-0,18	<b>-0,64</b>
7	<b>0,36</b>	-0,03	0,33	16	-0,07	0,10	-0,16
8	<b>0,32</b>	0,14	-0,03	17	-0,03	<b>0,47</b>	-0,04
9	<b>0,49</b>	0,32	-0,13	18	0,06	<b>0,57</b>	0,07

Considerando-se a não-disponibilidade, para o presente estudo, de programas computacionais específicos para análise multidimensional, optou-se por estudar as dimensões da prova separadamente, o que pode ser realizado segundo Embretson (1999). Assim, uma nova análise fatorial foi realizada com os itens do primeiro fator, citados no parágrafo anterior, sendo possível confirmar que a prova passa a ser predominantemente unidimensional, com variância explicada de 36,7% e consistência interna significativa entre os itens ( $r=0,753$ ). Com este último agrupamento dos itens, se considerado o modelo bidimensional, constata-se que a variância explicada aumenta para 40,8%, mas a magnitude da diferença entre os padrões de resposta observados e os reproduzidos pelo modelo não é estatisticamente significativa quando se acrescenta um segundo fator, podendo-se supor que tal fator não é necessário.

Depois de verificada a unidimensionalidade da prova procedeu-se às análises estatísticas TCT e TRI unidimensional, com o auxílio dos softwares

ITEMAN, RASCAL e XCALIBRE, cujos resultados estão apresentados e comentados a seguir.

### Análise de itens pela Teoria de Resposta ao Item

Foram analisados os ajustes dos dados aos modelos de um, dois e três parâmetros considerando-se os 18 itens iniciais e posteriormente os itens que compõem a prova unidimensional. Os resultados das estatísticas gerais dos itens indicam a adequabilidade ou não dos itens aos modelos propostos. Essas estatísticas são calculadas a partir: do agrupamento dos escores de habilidade semelhante; da probabilidade de acerto teórica para cada subgrupo; da curva característica do item; e dos resíduos relativos à probabilidade real observada.

Conforme resultados apresentados por Vendramini (2002), os dados que não se ajustam ao modelo de Rasch são referentes aos itens 5, 13, 15, 16, 17 e 18 ( $Qui^2 [13, N=325] > 22,36; p < 0,05$ ) de maior dificuldade na prova, com exceção do item 5, que é fácil.

Após a constatação de que o modelo logístico de um parâmetro implicava a eliminação de 6 dos 18 itens da prova, por não atenderem às condições exigidas, verificou-se que essa eliminação poderia comprometer a representatividade do domínio avaliado. Por isso decidiu-se testar também os modelos de dois e três parâmetros, cujos resultados apresentam, entre outras, as seguintes informações: indicação de itens âncora, advertências sobre problemas de ajuste, índice de discriminação, índice de dificuldade, probabilidade de acerto casual, resíduos padronizados indicando o ajuste do item (cujos valores maiores que 2,0 indicam discrepância significativa), correlação ponto-bisserial entre a resposta ao item e o total da prova, correlação entre a resposta ao item e o traço latente e número de participantes considerado para a estimação dos parâmetros. Os problemas que podem ocorrer com os parâmetros dos itens correspondem a: (1) valor do índice de discriminação  $a$  abaixo do valor crítico 0,30; (2) índice de dificuldade  $b$  acima do valor crítico 2,95 ou abaixo de -2,95; (3) probabilidade de acerto casual acima do valor crítico 0,40; (4) erro de chaveamento, isto é, correlação entre uma das opções de resposta incorreta e o escore total mais alta que a correlação entre a resposta correta e o escore total; e (5) resíduos padronizados do ajuste do modelo que excedem o valor crítico 2,0.

Na TCT, o índice de facilidade de um item corresponde a sua proporção de acertos, variando de 0 a 1. Quanto maior o índice de facilidade, menor será então a dificuldade do item. Na TRI, é considerado o índice de dificuldade, estimado a partir das respostas

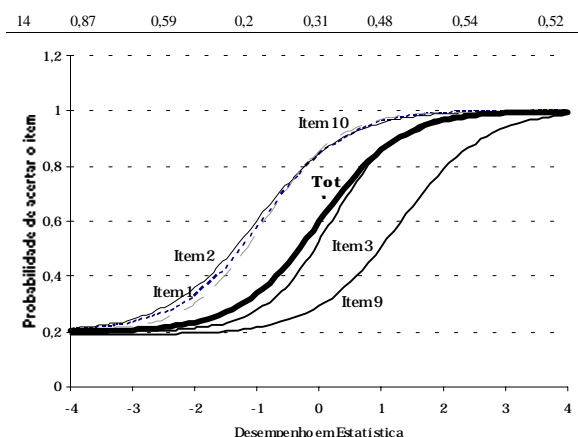
dos participantes, considerando-se a probabilidade de um indivíduo de certa habilidade responder corretamente a um item. A fase inicial de estimação é baseada na proporção de acertos e correlação bisserial, e refinada segundo procedimentos estatísticos. Os índices de dificuldade são transformados em uma escala, que pode variar de  $-\infty$  a  $\infty$ , sendo inicialmente limitados de -3 a 3. Um item com índice de dificuldade -3 é considerado extremamente fácil, zero de dificuldade média, e 3 extremamente difícil.

A partir da análise descrita nos parágrafos anteriores constatou-se que apenas o item 5, dentre os citados como não ajustados ao modelo de um parâmetro, se ajusta aos modelos de dois e três parâmetros, sendo que os itens 13 e 18 se ajustam apenas ao modelo de dois parâmetros, e os itens 15, 16 e 17 a nenhum dos modelos analisados. Para o modelo de dois parâmetros, os itens 15, 16 e 17 apresentam erro de chaveamento, sendo que o item 16 também apresenta discrepância significativa do ajuste do modelo (Resíduo > 2,0). No modelo de três parâmetros, além dos erros de chaveamento indicados para os itens 13, 15, 16, 17 e 18, os itens 15 e 16 apresentam índices de dificuldade acima do valor crítico 2,95.

Esses resultados orientaram uma outra análise, composta pelos itens da dimensão predominante, cujos resultados indicaram que nenhum dos itens apresenta erro de chaveamento (correlação entre uma das opções de resposta incorreta e o escore total mais alta que a correlação entre a resposta correta e o escore total), nem discrepância significativa do modelo ajustado (Resíduo > 2,0). A Tabela 3 e Figura 2 revelam que os itens referentes à obtenção direta das informações nos gráficos e de cálculos sem envolver números relativos (1, 2, 3 e 10) foram os mais fáceis da prova, enquanto o item 9, que envolve comparações de dados sobre distância e tempo, o mais difícil. Os itens 3, 5 e 10 foram os mais discriminativos e apresentaram correlações mais elevadas com a prova total do que com os outros itens.

**Tabela 3.** Estimativas finais dos parâmetros do fator dominante.

Item	Índice de discriminação (a)	Índice de dificuldade (b)	Probabilidade de acerto ao acaso (c)	Resíduo	Proporção de acertos	Correlação Ponto-bisserial (PBs)	Correlação bisserial (PBI)
1	0,89	-0,91	0,2	0,21	0,77	0,52	0,55
2	0,83	-1,01	0,2	0,22	0,79	0,49	0,51
3	1,09	-0,81	0,2	0,47	0,76	0,63	0,66
4	0,85	0,2	0,2	0,27	0,56	0,54	0,52
5	1,06	0,2	0,19	0,73	0,54	0,67	0,68
6	0,86	0,49	0,2	0,35	0,5	0,53	0,52
9	0,87	1,3	0,19	0,39	0,35	0,47	0,42
10	1,00	-0,88	0,2	0,28	0,77	0,58	0,6
11	0,81	0,31	0,2	0,32	0,55	0,48	0,47
12	0,82	0,67	0,2	0,25	0,48	0,49	0,45



**Figura 2.** Curvas características de alguns itens da prova de estatística.

## DISCUSSÃO

A aplicação da TRI no estudo, construção e validação de testes psicológicos e educacionais tem sido muito utilizada por vários pesquisadores e instituições (INEP, 2002; Pasquali, 2002; Ziviani & Primi, 2002; entre outros). Esse método de análise muito tem contribuído para enriquecer as análises feitas referentes aos itens de instrumentos psicológicos ou educacionais, razão pela qual se optou por essa teoria para a análise dos itens da prova de estatística.

As análises realizadas com todos os itens confirmaram o não-ajuste do reduzido número dos itens que compunham a segunda e terceira dimensão, optando-se apenas pela análise da dimensão predominante. Segundo Embretson (1999), as diferentes dimensões podem ser analisadas separadamente, na falta de programas específicos para uma análise multidimensional, embora já existam programas computacionais específicos para se proceder a essa análise.

Os instrumentos utilizados apresentaram um índice de confiabilidade satisfatório ( $\alpha=0,735$ ) (Constatou-se que, segundo a TCT, cinco das 18 questões propostas na prova apresentaram índices de discriminação baixos (menores ou iguais a 0,20) e proporções baixas de acertos (menores ou iguais a 0,29). O índice de maior discriminação (0,77) foi o da questão 5. De acordo com a TRI, verificou-se que os mesmos cinco itens apontados pela TCT como de baixa discriminação e baixa porcentagem de acertos tiveram seus ajustes rejeitados segundo a interpretação do modelo de Rasch  $Qui^2[s[13, N=325]>26,93; p's<0,0127]$ . O modelo de Rasch também não representou um bom ajuste para a questão 5  $Qui^2[s[13, N=325]=24,098; p=0,0303]$ .

Ao se ajustar o modelo de um parâmetro aos itens da prova, a maioria dos itens referentes à interpretação de dados apresentados em tabela indicou índices de dificuldade elevados, enquanto nos itens referentes à interpretação de dados apresentados em gráficos os índices de dificuldade foram os mais baixos.

Destarte, o modelo de um parâmetro mostrou-se não apropriado para 5 dos 6 itens referentes à interpretação de dados apresentados em tabela. Mesmo considerando-se os modelos de dois e três parâmetros, alguns dos itens continuaram a não se ajustar aos modelos propostos. A análise realizada indicou que a prova final de estatística, composta pelos itens 1, 2, 3, 4, 5, 6, 9, 10, 11, 12 e 14, é predominantemente unidimensional e com consistência interna significativa de 0,753. Os itens apresentaram propriedades (índice de dificuldade, de discriminação e correlação bisserial) dentro de padrões aceitáveis, garantindo parâmetros invariantes para os itens e para o desempenho em estatística, pois, conforme afirmam Hambleton, Swaminathan e Rogers (1991), só com um modelo bem ajustado isso é possível.

Além disso, os itens que permaneceram na prova final referem-se às interpretações baseadas em representações gráficas e de tabelas, atendendo ao propósito de obter informações sobre o raciocínio estatístico dos alunos, as quais, segundo Garfield e Gal (1999), são elementos essenciais para o pensamento estatístico.

Como descrito na sessão anterior, os parâmetros dos modelos da TRI, são estimados inicialmente com base na proporção de acertos e correlação bisserial, calculadas pela teoria clássica, e refinados segundo procedimentos estatísticos que buscam estimadores de máxima verossimilhança e de resíduos mínimos. Os índices de dificuldade  $b$  são transformados inicialmente em uma escala, que varia de -3 (itens extremamente fáceis) a 3 (itens extremamente difíceis); os de discriminação  $a$ , que variam de 0,5 (baixa discriminação) a 2,0 (alta discriminação); a probabilidade de acerto por acaso  $c$ , a partir do número de alternativas do item, neste estudo  $c = 0,2$  (uma de cinco alternativas).

Quando comparados os índices de facilidade da TCT com os de dificuldade da TRI, da prova final, observa-se que, quanto maior a proporção de acertos no item, menor tende a ser o seu índice de dificuldade. A vantagem da TRI é que, conhecida a habilidade de um indivíduo, não necessariamente participante da amostra, é possível determinar a probabilidade de ele acertar um item. Isto não ocorre com a TCT, cujos resultados são dependentes da amostra. Os itens da



prova aplicada são de nível mediano, não são muito difíceis ( $b_i < 1,30 < 3$ ) e não muito fáceis ( $b_i > -0,88 > -3$ ).

Os índices de discriminação da TCT são calculados a partir da diferença entre a proporção de acertos dos 27% dos participantes com maior pontuação total e a proporção de acertos dos 27% de menor pontuação total. Na TRI, esses índices são estimados a partir da correlação bisserial item total, e diferem dos encontrados pela TCT. Os valores apresentados na Tabela 4 indicam que os itens não são de grande poder discriminativo, tanto pela TRI ( $0,5 < 0,81 < a_i < 1,09 < 2$ ), quanto pela TCT (maioria dos índices próximos de 0,50).

**Tabela 4.** Índices de discriminação e dificuldade dos itens segundo a TRI e TCT.

Item	Índice de discriminação		Índice	
	TRI (a)	TCT	Dificuldade TRI (b)	Facilidade TCT (%acertos)
1	0,89	0,46	-0,91	0,77
2	0,83	0,38	-1,01	0,79
3	1,09	0,54	-0,81	0,76
4	0,85	0,60	0,20	0,56
5	1,06	0,77	0,20	0,54
6	0,86	0,52	0,49	0,50
9	0,87	0,52	1,30	0,35
10	1,00	0,50	-0,88	0,77
11	0,81	0,55	0,31	0,55
12	0,82	0,50	0,67	0,48
14	0,87	0,62	0,59	0,48

Os resultados aqui relatados indicam a necessidade de outros estudos, com a inclusão de novos itens na prova, para se chegar a uma prova que apresente índices mais elevados de confiabilidade e validade para ser utilizada no contexto educacional e que revelem o raciocínio estatístico de universitários no que diz respeito à leitura de dados apresentados em gráficos e tabelas estatísticas.

## REFERÊNCIAS

- Andrade, D. F. (2001). Comparando desempenhos de grupos de alunos por intermédio da teoria da resposta ao item. *Estudos em Avaliação Educacional*, 23, 31-69.
- Andrade, D. F., & Valle, R. C. (1998). Introdução à teoria da resposta ao item. *Estudos em Avaliação Educacional*, 18, 13-32.
- Andrade, D. F., Tavares, H. R., & Valle, R. C. (2000). *Teoria da resposta ao item: conceitos e aplicações*. São Paulo: ABE.
- Andriola, W. B. (1998). Avaliação da aprendizagem: uma análise descritiva segundo a teoria de resposta ao item (TRI). *Educação em Debate*, 20(36), 93-102.
- Assessment System Corporation (1995a). *User's manual for the RASCAL-Rasch Analysis Program*. 2nd ed., Windows version 3.50e, St. Paul, MN: Author.
- Assessment System Corporation (1995b). *User's manual for the XCALIBRE-Marginal Maximum-Likelihood IRT Parameter Estimation Program*. (2nd ed.), Windows 3.x/95/NT version, St. Paul, MN: Author.
- Assessment System Corporation (1998). *User's manual for the ITEMAN-Conventional Item Analysis Program*. 2nd ed., Windows 3.x/95/NT version, St. Paul, MN: Author.
- Bock, R. D., Gibbons, R., & Muraki, E. (1988). Full information item factor analysis. *Applied Psychological Measurement*, 12(3), 261-280.
- Embretson, S. E. (1999). Generating items during testing: Psychometric issues and models. *Psychometrika*, 64(4), 407-433.
- Fletcher, P. (1994). A teoria da resposta ao item: medidas invariantes do desempenho escolar. *Ensaio: avaliação e políticas públicas em educação*, 1(2), 21-28.
- Garfield, J. (2002). The Challenge of Developing Statistical Reasoning. *Journal of Statistics Education* 10(3). Disponível em: <www.amstat.org/publications/jse/v10n3/garfield.html> (Acessado em 21/04/2003)
- Garfield, J., & Gal, I. (1999). Teaching and assessing statistical reasoning. In L. Stiff (Ed.). *Developing mathematical reasoning in grades K-12*. (pp. 207-219). Reston, VA: National Council Teachers of Mathematics.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park: Sage Publishers.
- Instituto Nacional de Estudos e Pesquisas Educacionais - INEP (2002). *Sistema nacional de avaliação da educação básica: relatório nacional 2001*. Brasília: INEP.
- Lalonde, R. N., & Gardner, R. C. (1993). Statistics as a second language? A model for predicting performance in psychology students. *Canadian Journal Behavioral Science*, 25(1), 108-125.
- Pasquali, L. (1999). *Instrumentos psicológicos: manual prático de elaboração*. Brasília: Laboratório de Pesquisa em Avaliação e Medida - LabPAM.
- Pasquali, L. (2002). Provão (ENC) de Psicologia 2000 e 2001: Análise dos parâmetros psicométricos. Em R. Primi (Org.). *Temas em avaliação psicológica* (pp.152-178). Campinas: Instituto Brasileiro de Avaliação Psicológica.
- Primi, R. (1998). *Desenvolvimento de um instrumento informatizado para a avaliação do raciocínio analítico*. Tese Doutorado Não-Publicada. Instituto de Psicologia. Universidade de São Paulo.
- Primi, R., Almeida, L. S. (2001). Teoria de Resposta ao Item. Em: E. M. Fernandes & L. S. Almeida (Orgs.). *Métodos e técnicas de avaliação: contributos para a prática e investigação psicológicas*. (pp. 205-232). Braga: Centro de Estudos em Educação e Psicologia, Universidade do Minho.
- Vendramini, C. M. M. (2000). *Implicações das atitudes e das habilidades matemáticas na aprendizagem dos conceitos de Estatística*. Tese de Doutorado Não-Publicada. Faculdade de Educação. Universidade Estadual de Campinas.
- Vendramini, C. M. M. (2002). Aplicação da teoria de resposta ao item na avaliação educacional. Em Primi R. (Org.). *Temas em avaliação psicológica* (pp.116-130). Campinas: Instituto Brasileiro de Avaliação Psicológica.
- Wilson, D. T., Wood, R., & Gibbons, R. (1998). *TESTFACT 2 - test scoring, item statistics, and item factor*. Chicago: Scientific Software International.

Ziviani, C., & Primi, R. (2002). Teoria de resposta ao item e o modelo Rasch de mensuração: uma análise do provão de Psicologia. Em R. Primi (Org.). *Temas em avaliação psicológica* (pp.131-151). Campinas: Instituto Brasileiro de Avaliação Psicológica.

Recebido em 17/11/2003

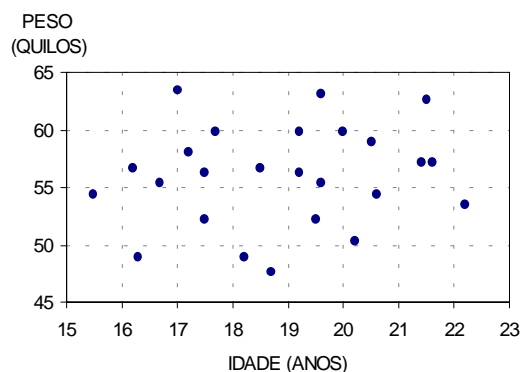
Aceito em 09/07/2004

## ANEXO – PROVA DE ESTATÍSTICA

Início: \_\_\_\_\_ horas \_\_\_\_\_ minutos

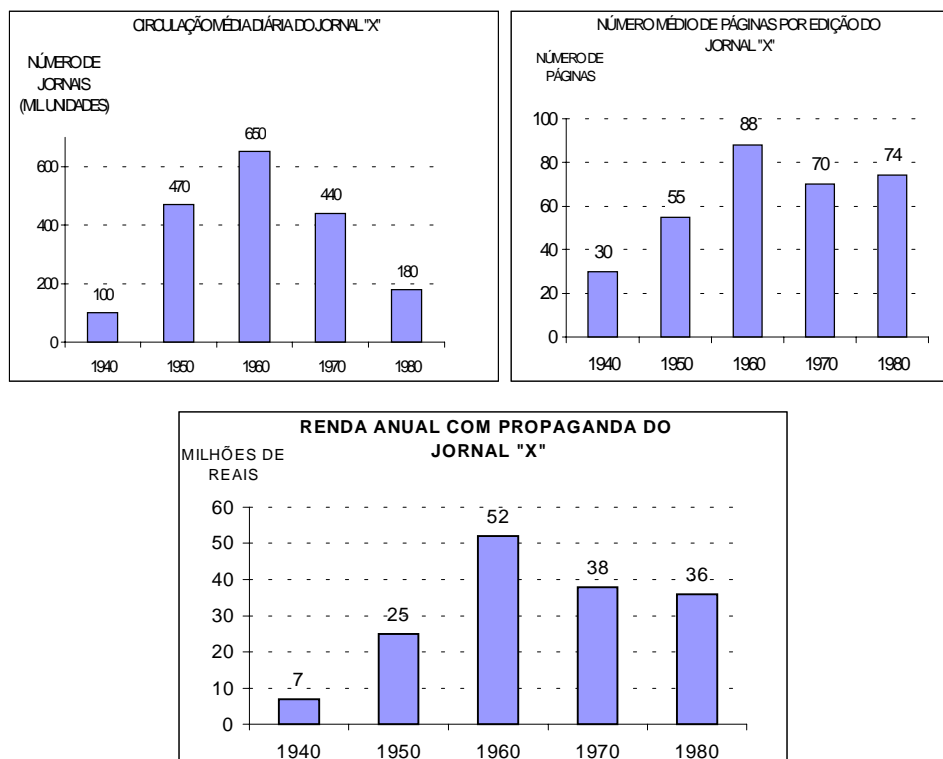
Assinale a alternativa mais adequada das cinco alternativas propostas

### IDADE E PESO DE 25 ESTUDANTES



1. Os pontos do gráfico a seguir indicam peso e idade de 25 estudantes. Qual a porcentagem desses estudantes com idade inferior a 19 anos e peso superior a 50 quilos?
- a) 48%                      b) 88%                      c) 36%                      d) 52%                      e) n. d. a.

As questões de número 2 a 6 referem-se aos seguintes gráficos



2. Em quantos dos anos mostrados o número médio de páginas por edição foi no mínimo o dobro da média obtida em 1940?  
a) quatro b) três c) dois d) um e) nenhum
3. Em 1950, se o custo de impressão por jornal fosse igual a R\$ 0,05, qual seria o custo médio total de impressão de uma circulação diária?  
a) R\$ 23.500,00 b) R\$ 32.500,00 c) R\$ 26.000,00 d) R\$ 22.000,00 e) R\$ 2.350,00
4. Em 1980, quantas vezes a renda anual com propaganda foi maior que a circulação média diária?  
a) 500 b) 200 c) 100 d) 50 e) 20
5. A porcentagem de decréscimo ocorrido na circulação média diária de 1960 para 1970 foi aproximadamente de:  
a) 10% b) 12% c) 20% d) 26% e) 32%
6. Qual das afirmações a seguir podem ser inferidas dos dados?  
I. O maior acréscimo ocorrido na renda anual com propaganda sobre algum período de 10 anos foi igual a R\$ 27 milhões.  
II. Em cada um dos períodos de 10 anos em que ocorreu um decréscimo na renda anual com propaganda a circulação média diária também decresceu.  
III. De 1970 para 1980 o número médio de páginas por jornal cresceu em 10 unidades.  
a) I somente b) II somente c) III somente d) I e II e) II e III
7. Qual das alternativas a seguir é igual a  $1/4$  de 0,01 por cento?  
a) 0,000025 b) 0,00025 c) 0,0025 d) 0,025 e) 0,25
8. Se  $d = 5,03$  e  $g$  a expressão de  $d$  arredondada para décimos.  
Comparando os valores  $d$  e  $g$  podemos afirmar que:  
a)  $d$  0,6% superior a  $g$  b)  $d$  0,03% superior a  $g$  c)  $d$  é igual a  $g$   
d)  $d$  1,006% superior a  $g$  e)  $d$  0,006% superior a  $g$
9. Se  $d_1$  = tempo necessário para viajar  $d$  quilômetros a  $s$  quilômetros por hora;  $d_2$  = tempo necessário para viajar  $d/2$  quilômetros a  $2s$  quilômetros por hora. Considerando ( $d_2 \neq 0$ ) então: a)  $d_1$  maior que  $d_2$  b)  $d_1$  menor que  $d_2$  c)  $d_1 = d_2$  d)  $d_1 = 0$  e) não é possível responder
10. Numa certa loja, os cadernos são normalmente vendidos por R\$ 0,59 cada um. Numa segunda loja são vendidos dois desses cadernos por R\$ 0,99. Quanto pode ser economizado na compra de 10 cadernos na segunda loja?  
a) R\$ 0,85 b) R\$ 0,95 c) R\$ 1,10 d) R\$ 1,15 e) R\$ 2,00
11. Se a média (média aritmética) de 5 números inteiros consecutivos é 12, qual é a soma do menor com o maior desses 5 números?  
a) 24 b) 14 c) 12 d) 11 e) 10
12. O número 0,01 é quantas vezes maior que o número  $(0,0001)^2$ ?  
a)  $10^2$  b)  $10^4$  c)  $10^6$  d)  $10^8$  e)  $10^{10}$

## As questões de número 13 a 18 referem-se aos dados a seguir

PERFIL DO CONGRESSO AMERICANO NO ANO "X" - (total de membros: 535)

Total de representantes		Senado	Total de representantes		Senado
Partido			Profissão		
292	Democrático.....	62	215	Advogado.....	63
143	Republicano.....	38	81	Executivo/Banqueiro	15
435	Total.....	100	45	Educador.....	6
	<b>Sexo</b>		14	Agricultor/Fazendeiro	6
418	Masculino.....	100	22	Oficial do Governo	0
17	Feminino.....	0	24	Jornalista, Executivo de Comunicação.....	4
	<b>Idade</b>		2	Médico.....	0
27	Mais Novo.....	34	1	Veterinário.....	1
77	Mais Velho.....	80	0	Geólogo.....	2
48	Média(média aritmética)	54	6	Operário, Comerciante Especializado....	0
	<b>Religião</b>		25	Outra.....	3
255	Protestante.....	69			
107	Católico.....	12			
18	Judeu.....	5			
4	Mórmon.....	3			
51	Outra.....	11			

13. Qual das afirmações a seguir são, com certeza, verdadeiras?
- I. A variabilidade total de idades do Senado é maior que do total de Representantes
  - II. Mais de 30% dos Congressistas são democráticos.
  - III. 50% do Senado tem idades superiores ou iguais a 54 anos.
- a) I somente    b) II somente    c) III somente    d) I e II    e) II e III
14. No Senado, se 25 homens forem substituídos por 25 mulheres, a razão de homens para mulheres será:
- a) 4 para 1    b) 3 para 1    c) 3 para 2    d) 2 para 1    e) 1 para 1
15. Qual porcentagem de membros do Congresso que são advogados?
- a) 63%    b) 58%    c) 56%    d) 52%    e) 49%
16. Se 5 senadores são Democratas Católicos, quantos senadores não são Católicos e nem democratas?
- a) 79    b) 74    c) 69    d) 31    e) 21
17. Se todos os advogados e todas as mulheres do total de Representantes votam pela passagem de um projeto de lei, quantos votos mais serão necessários para a maioria?
- a) 435    b) 220    c) 3    d) 0    e) Não pode ser determinado pela informação dada.
18. O que se pode inferir a partir das informações apresentadas?
- I. Mais de 80% dos homens do Congresso são membros Representantes.
  - II. A porcentagem de membros que são Agricultores ou Fazendeiros é maior para o total de Representantes do que para o Senado.
  - III. A idade mediana do Senado é 57.
- a) I somente    b) II somente    c) III somente    d) I e II    e) II e III

Término: \_\_\_\_ horas \_\_\_\_ minutos

---

**Endereço para correspondência:** Claudette Maria Medeiros Vendramini: Rua Herculano Pupo Nogueira, 309, Vila Belém, CEP 13256-300, Itatiba-SP. E-mail: cvendramini@uol.com.br