

# Análise Descritiva de Dados de Pesquisa

Luiz Pasquali

## Preliminar

Tipicamente uma pesquisa científica produz no final uma série (para não dizer, uma carrada) de números. O que é que esses números representam? Eles querem representar a informação que você queria obter quando quis realizar uma pesquisa. Os números são a linguagem (matemática) na qual a informação da pesquisa está sendo descrita, codificada. Agora, os números, quando usados para descrever fenômenos da natureza, podem ser um tanto traiçoeiros, porque eles não são simplesmente os números com os quais os matemáticos trabalham, os quais são claros e precisos, definidos pelos 27 axiomas que Whitehead e Russell (1937) conseguiram descobrir sobre os mesmos. Os números, quando são usados pelo cientista que estuda a natureza (este cientista não estuda os números em si), perdem muito das características que os 27 axiomas expressam. Embora os números dos matemáticos e dos cientistas da natureza sejam graficamente idênticos, o significado deles é quase sempre muito distinto. Assim, quando você tem uma série, uma escala, de números, tal como:

1 2 3 4 5 6 7 8 9 10

para o matemático ela tem um sentido muito claro e único; mas para o cientista ela pode significar coisas muito diferentes. Como assim? É que tais escalas de números para o cientista podem ter diferentes níveis de qualidade numérica, isto é, de axiomas do número; e esses diferentes níveis de qualidade afetam muito o tipo de informação que os números trazem sobre o objeto que o cientista está estudando. Assim, a primeira coisa que devemos discutir neste capítulo de descrição dos dados de uma pesquisa por meio dos números consiste em dar uma olhada nessa história das escalas de números ou escalas de medida.

# Parte I

## As Escalas de Números

### (Escalas de Medida)

A qualidade de uma escala de medida depende da quantidade de axiomas (ordem, intervalo, origem etc.) do número matemático que tal escala salva. Assim, dependendo da quantidade de axiomas do número que a medida salva, resultam vários níveis de medida, conhecidos como escalas de medida. São três os axiomas básicos do número: identidade, ordem e aditividade. O último apresenta dois aspectos úteis para o presente problema: origem e intervalo ou distância. Quanto mais axiomas do número a medida salvaguardar, maior será o seu nível. Assim, podemos considerar cinco elementos numéricos para definir o nível da medida: identidade, ordem, intervalo, origem, e unidade de medida. Desses cinco elementos, os mais discriminativos dos níveis são a origem e o intervalo, dado que a ordem é uma condição necessária para que realmente haja medida. Se a medida somente salva a identidade do número, na verdade não se trata de medida, mas sim de classificação e contagem. Nesse caso (*escala nominal*), os números não são atribuídos a atributos dos objetos, mas o próprio objeto é identificado por rótulo numérico. Esse rótulo nem precisaria ser numérico dado que não importa que símbolo ou rabisco pode ser utilizado com a mesma função de distinguir objetos um do outro ou classe de objetos de outra classe. A única condição necessária é que se salvasse a identidade do símbolo, isto é, um mesmo símbolo não pode ser duplicado para identificar objetos diferentes, como também diferentes símbolos não podem ser usados para identificar objetos idênticos. Embora não estejamos nesse caso medindo, a escala numérica que resulta dessa rotulação adquire direito ao nome escala, dado que ela corresponde em parte à definição genérica de medida que reza “medir é atribuir números às coisas empíricas” (Pasquali, 2003).

O esquema a seguir ilustra como se originam as várias escalas de medida:

		Origem	
		Não-Natural	Natural
Intervalo	Não-Igual	Ordinal	Ordinal
	Igual	Intervalar	Razão

Assim, uma medida de uma propriedade de um objeto natural que não tem uma origem natural (exemplos: aroma, QI, amizade) não pode começar em zero (0), porque não se conhece um valor zero de aroma ou de QI. Mesmo se você usa o zero na medida de tais atributos, este é um zero fictício, não natural. Dessa forma, uma escala de medida de tais atributos pode começar com qualquer número, inclusive o zero, sendo esse número a origem da escala e o próximo número tem que ser maior (se a escala for ascendente) porque a escala precisa salvar, pelo menos, a ordem natural dos números. Tal escala seria chamada de *ordinal*, onde a origem é arbitrária e a distância entre os números não seria igual. Conseqüentemente, as seguintes escalas são equivalentes, produzem exatamente a mesma informação:

3	4	5	6	7
3	5	6	10	100
0	1	2	3	4
-3	0	15	30	31

A única coisa relevante que distingue essas escalas é uma questão de estética, sendo provavelmente a mais elegante a escala 0 1 2 3 4. Mas elegância é questão de gosto e “de gustibus non est disputandum” (não se briga por gostos). Agora, seria um erro transformar essas escalas na seguinte:

0	1	2	4	3
---	---	---	---	---

porque se perderia a ordem (monotônica crescente).

Se nessa mesma medida você puder salvar a origem natural, isto é, o zero, mas não puder salvar o intervalo igual entre os números da

escala, você ainda estaria medindo apenas ao nível ordinal. Por exemplo: medir o peso de diferentes objetos sem ter uma balança. Nesse caso, você pode pedir a um ou vários sujeitos para ordenar os objetos em termos de mais pesado, surgindo daí uma ordenação dos mesmos pelo peso sem se poder dizer quanto um objeto é mais pesado que o anterior. Peso, na verdade, é um atributo da natureza que permite o valor 0, mas o processo de medida, como descrito, não permite dizer mais do que um objeto ser mais pesado que o outro, sem se poder definir quanto mais. Se você pudesse ou puder definir quanto mais pesado ele é, então você já estaria medindo ao nível de *escala de razão* que, além de ter uma origem natural 0, tem intervalos iguais entre os números da escala. Tal escala sempre começa com 0 e seus números estão a distâncias iguais entre si. Exemplo:

0	1	2	3	4
0	2	4	6	8
0	5	10	15	20

Nessa escala o que muda é apenas a unidade de medida (o tamanho do intervalo), sendo ela, no presente caso, sucessivamente de 1, 2 e 5.

Um exemplo de uma escala simplesmente *intervalar* e suas transformações legítimas seria a seguinte:

0	2	4	6	8
2	4	6	8	10
-5	0	5	10	15

onde são salvos a ordem dos números e o tamanho do intervalo entre eles, sendo arbitrária a sua origem.

Note que o fator que define o nível da medida não é o número, mas sim a característica do atributo medido da natureza (da realidade): se ele permite ou não uma ordem natural, o 0 (tais como, peso, comprimento,...), se permite definir distâncias iguais ou não (muitos pesquisadores afirmam que nenhum atributo não extensivo da natureza, como todos os atributos psicossociais, permite medida intervalar!). Os números, por natureza, têm todas estas características:

origem natural (o 0), ordem e distâncias iguais entre si. Assim, para os números, todas as escalas são de razão, mas para a medida, isto é, os números utilizados para descrever fenômenos naturais, nem sempre é possível salvarem-se essas características dos números.

A Tabela 7-1 sumaria as características de cada escala.

**Tabela 7-1. Características das escalas numéricas de medida**

Escala	Axiomas Salvos	Invariâncias	Liberdades	Transformações Permitidas	Estatísticas Apropriadas
Nominal	- identidade		- ordem - intervalo - origem - unidade	Permutação (troca 1 por 1)	Frequências: f, %, p, Mo, qui <sup>2</sup> , C
Ordinal	- identidade - ordem	- ordem	- intervalo - origem - unidade	Monotônica crescente (isotonia)	Não-paramétricas: Md, r <sub>s</sub> , U, etc.
Intervalar	- identidade - ordem - aditividade	- ordem - intervalo	- origem - unidade	Linear de tipo Y = a+bx	Paramétricas: M, DP, r, t, f, etc.
Razão	- identidade - ordem - aditividade	- ordem - intervalo - origem	- unidade	Linear de tipo y = bx (similaridade)	M geométrica, Coef. variação, Logaritmos

f = frequência; % = percentagem; p = proporção; C = coeficiente de contingência;  
Md = mediana; DP = desvio padrão; r<sub>s</sub> = correção de Spearman; U = teste de Mann-Whitney;  
r = correlação produto-momneto de Pearson; f = teste de Fisher (análise da variância)

Como já insinuado, uma escala numérica pode ser transformada numa outra equivalente se forem respeitados os elementos da invariância nessa transformação. Uma escala de maior nível pode utilizar as operações estatísticas de uma escala inferior, mas perde informação dado que as estatísticas próprias de uma escala inferior são menos eficientes, isto é, são menos robustas. Por exemplo, posso organizar o leque de idades dos sujeitos em quatro grupos etários (adolescência, jovem-adulto, adulto e terceira idade); mas, nesse caso, a partir desses grupos não posso saber a média das idades da amostra. Não é permitido (é erro) utilizar estatísticas de uma escala de nível superior numa inferior, dado que esta última não satisfaz os requisitos necessários para se utilizarem procedimentos estatísticos superiores. São chamados paramétricos os procedimentos estatísticos da escala intervalar porque os números nela possuem caráter métrico, isto é, são adicionáveis, enquanto os não-paramétricos não são métricos, dado que representam somente postos e não quantidades somáveis.

## Sumarizando:

*Escala nominal:* os números são usados simplesmente como numerais, isto é, símbolos gráficos que diferem entre si, como uma fotografia de uma maçã difere da de uma pedra ou um círculo difere de uma linha reta. Eles se distinguem apenas em termos de qualidade (gráfica) e não em termos de quantidade. Assim, se você chama de grupo 1 os sujeitos masculinos de uma classe de alunos e de grupo 2 os femininos, o 2, neste caso, não é maior do que o 1, ele é apenas diferente do 1. Dessa forma, você bem podia chamar de grupo 1 os femininos e de grupo 2 os masculinos, sem mudar nada na informação que eles trazem. Isto é, os números viram nomes ou rótulos. Mesmo assim, existem tratamentos estatísticos que podem ser utilizados com escalas nominais ou categóricas, como veremos em capítulos subseqüentes.

*Escala ordinal:* o número aqui já é utilizado como número, isto é, ele expressa uma quantidade ou magnitude. E, conseqüentemente, os números da escala já diferem em termos de tamanho e não podem ser intercambiados, porque eles têm que seguir a ordem natural dos números, onde um é maior que o outro. Os estatísticos dizem isso da seguinte maneira: se tiver dois objetos,  $o_1$  e  $o_2$ , que diferem numa magnitude  $m$  qualquer (como peso), então resulta uma escala ordinal se,

- 1)  $m(o_1) \neq m(o_2)$  i. é, o peso do  $o_1$  é diferente do peso do  $o_2$
- 2)  $m(o_1) < m(o_2)$  i. é, o peso do  $o_2$  é maior que o peso do  $o_1$ .

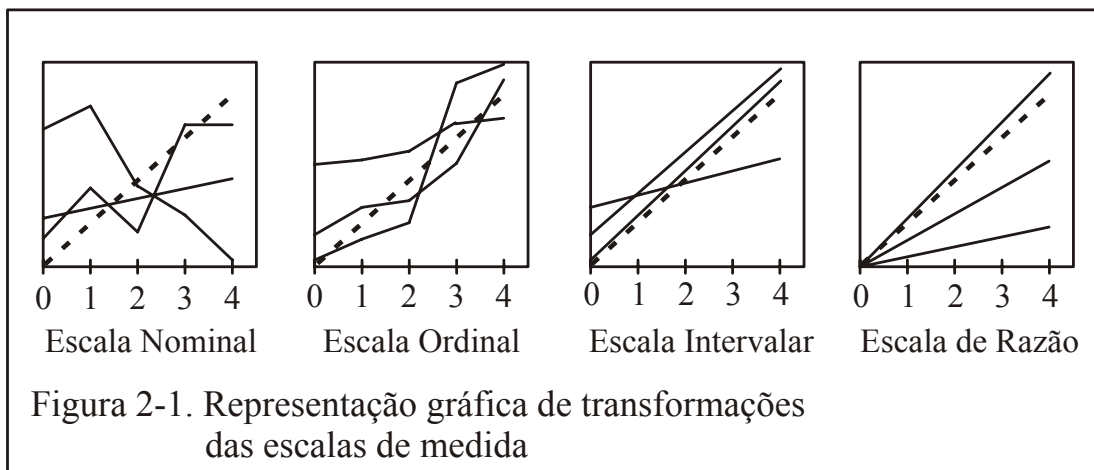
*Escala intervalar:* Além de seguir a ordem natural, os números nessa escala se distanciam igualmente um do outro, isto é, a distância ou o intervalo entre eles é sempre a mesma. Os estatísticos, novamente, dizem isso da seguinte maneira: se tiver dois objetos,  $o_1$  e  $o_2$ , que diferem em magnitude  $m$ , então resulta uma escala intervalar se,

- 1)  $m(o_1) \neq m(o_2)$
- 2)  $m(o_1) < m(o_2)$
- 3)  $m(o_1) = a[m(o_2)] + b$ ; isto é, os números da escala diferem entre si pela multiplicação por uma constante  $a$  e pela soma de outra constante  $b$ . Isso quer dizer que um segundo objeto difere do primeiro numa propriedade qualquer em  $a$  vezes o tamanho da propriedade do primeiro mais um tamanho constante qualquer adicionado (o  $b$ ). Isto é, há uma relação linear entre o peso dos dois objetos.

*Escala de razão:* Além de seguir a ordem natural e manter a distância igual entre os números da escala, os números na escala de razão têm uma origem natural, isto é, o 0. Assim, uma escala de razão sempre começa em 0 e todos os valores dela são referenciados a 0. Os estatísticos, novamente, dizem isso da seguinte maneira: se tiver dois objetos,  $o_1$  e  $o_2$ , que diferem em magnitude  $m$ , então resulta uma escala de razão se,

- 1)  $m(o_1) \neq m(o_2)$
- 2)  $m(o_1) < m(o_2)$
- 3)  $m(o_1) = a[m(o_2)]$ ; isto é, o  $o_2$  é  $a$  vezes maior que o  $o_1$  e, assim, eu posso dividi-los:  $m(o_2)/m(o_1)$ , isto é, estabelecer razões entre eles. O  $b$  foi eliminado da condição, porque qualquer transformação de uma escala de razão deve sempre começar em 0 e o  $b$  faria com que a escala pudesse começar em outro ponto da ordenada, que não fosse o 0, como é o caso com escalas de intervalo.

Veja a figura 2-1, onde a escala original representada pela reta grossa pontilhada pode ser transformada em qualquer outra escala (linhas inteiras), dependendo do tipo ou do nível da escala original ser nominal, ordinal etc.



Essa história do nível da medida é muito grave em psicologia e em ciências empíricas de um modo geral. Os matemáticos e os estatísticos trabalham os números como números; acontece, porém, que o cientista (empírico) utiliza os números para descrever a sua realidade; então, os números têm que ter relação com o objeto que estão pretendendo descrever. Assim, se a escala é de nível ordinal ou intervalar não é um problema que o matemático ou o estatístico pode



resolver; esse problema é do pesquisador que deve demonstrar que o seu objeto permite a história de magnitudes diferentes e que o instrumento que ele utilizou para aplicar os números a essas magnitudes dos objetos permite uma escala ordinal ou intervalar. Para o estatístico a escala ser deste ou daquele nível é uma suposição. Somente que essa suposição implica em operações muito distintas no tratamento dos números da medida. Em psicologia, particularmente, admitir-se medida ao nível ordinal não suscita maiores diatribes; já uma medida ao nível de intervalo é bastante controversa. De qualquer forma, os estatísticos inventaram séries de tratamentos de todas essas escalas, incluindo a escala nominal. De sorte que em tudo que for dito daqui por diante no tratamento dos dados, esta questão do nível da medida está sempre presente. De um modo particular, é preciso atender, sobretudo, ao fato da medida ser ao nível puramente ordinal ou se ela é de nível intervalar; isto porque, no caso da medida ser puramente ordinal, as estatísticas e análises estatísticas a serem utilizadas são aquelas que os estatísticos chamam de estatísticas não-paramétricas; no caso da medida ser de intervalo, as estatísticas utilizadas são chamadas de paramétricas<sup>1</sup>.

---

<sup>1</sup> Os conceitos de paramétrico e não-paramétrico são discutidos no capítulo 11.

# Parte II

## A Descrição dos Dados da Pesquisa

A primeira coisa a fazer com os dados de uma pesquisa científica consiste na descrição dos mesmos. Essa descrição ou apresentação dos dados é estatisticamente efetuada de duas formas: (1) uma descrição algébrica, que utiliza os números, e (2) uma descrição geométrica, essa fazendo uso de desenhos ou gráficos.

Como ficou dito, a descrição algébrica dos dados de uma pesquisa é feita por meio de números. Essa descrição pode ser *completa* ou *parcial*. A descrição é completa se ela apresenta todos os dados individuais da pesquisa (distribuição de freqüências); ela é parcial ou sumarizada se apresenta somente sumários dos dados (medidas de tendência central e de variabilidade). Assim, na descrição dos dados, nós temos que falar de três tópicos, a saber:

- Distribuição de freqüências
- Medidas de tendência central
- Medidas de variabilidade.

### I – Distribuição de Freqüências

Imagine que apliquei um teste a 2.426 sujeitos e obtive os seguintes resultados:

27	26	22	27	27	21	24	26	27	11	24
28	26	26	27	25	25	26	28	23	22	26
26	25	25	26	26	28	25	29	28	27	23
24	25	25	30	27	26	24	23	16	25	16
22	19	24	13	28	17	22	19	13	30	7
16	29	26	26	15	19	24	28	26	24	24
11	15	24	23	28	21	25	25	13	17	25
29	14	13	23	26	28	26	28	25	22	26
25	18	26	25	26	9					

..... (mais de dois mil e quatrocentos deles!)

Com uma montanha de tais dados, fica praticamente impossível tirar qualquer informação sobre as habilidades dos sujeitos que o teste queria medir. Mesmo assim, eu posso ter a informação completa dos dados se puder dizer quantos dos sujeitos ganharam cada um dos possíveis escores contidos na listagem acima. De fato, o teste consta de 30 itens e o sujeito recebe um ponto por cada item acertado; de sorte que os escores possíveis vão de 0 a 30. Com essa informação, posso construir uma listagem dos escores, contendo o escore (X) e quantos dos sujeitos (frequência ou f) obtiveram tal escore, e apresentar tal informação numa tabela de duas colunas. Além disso, posso acrescentar uma terceira coluna na qual é representado o quanto por cento a ocorrência de cada escore do teste acontece nos 2.426 casos; é a percentagem (%). Essa tabela é chamada de distribuição de frequência (veja tabela 2-2).

Tabela 2-2. Distribuição de frequência de escores no TRAD de 2.426 sujeitos

Escore (X)	Frequência (f)	Percentagem (%)	Percentagem acumulada
0	0	0	0,0
1	9	0,4	0,4
2	19	0,8	1,2
3	55	2,3	3,4
4	77	3,2	6,6
5	79	3,3	9,9
6	66	2,7	12,6
7	74	3,1	15,6
8	58	2,4	18,0
9	66	2,7	20,7
10	60	2,5	23,2
11	42	1,7	24,9
12	60	2,5	27,4
13	45	1,9	29,3
14	65	2,7	31,9
15	54	2,2	34,2
16	54	2,2	36,4
17	55	2,3	38,7
18	59	2,4	41,1
19	66	2,7	43,8
20	93	3,8	47,7

21	90	3,7	51,4
22	112	4,6	56,0
23	95	3,9	59,9
24	145	6,0	65,9
25	191	7,9	73,7
26	217	8,9	82,7
27	178	7,3	90,0
28	143	5,9	95,9
29	76	3,1	99,1
30	23	0,9	100,0

Arranjando a lista dos dados segundo a tabela 2-2, já dá para ter uma noção bem mais clara do desempenho dos sujeitos no teste. Inclusive, vê-se que quase ninguém tem os escores mais baixos ou mais altos do teste. Ademais, se você representa essa distribuição de frequência num gráfico, a informação fica ainda mais clara. Veja a figura 2-2 para tal ilustração.

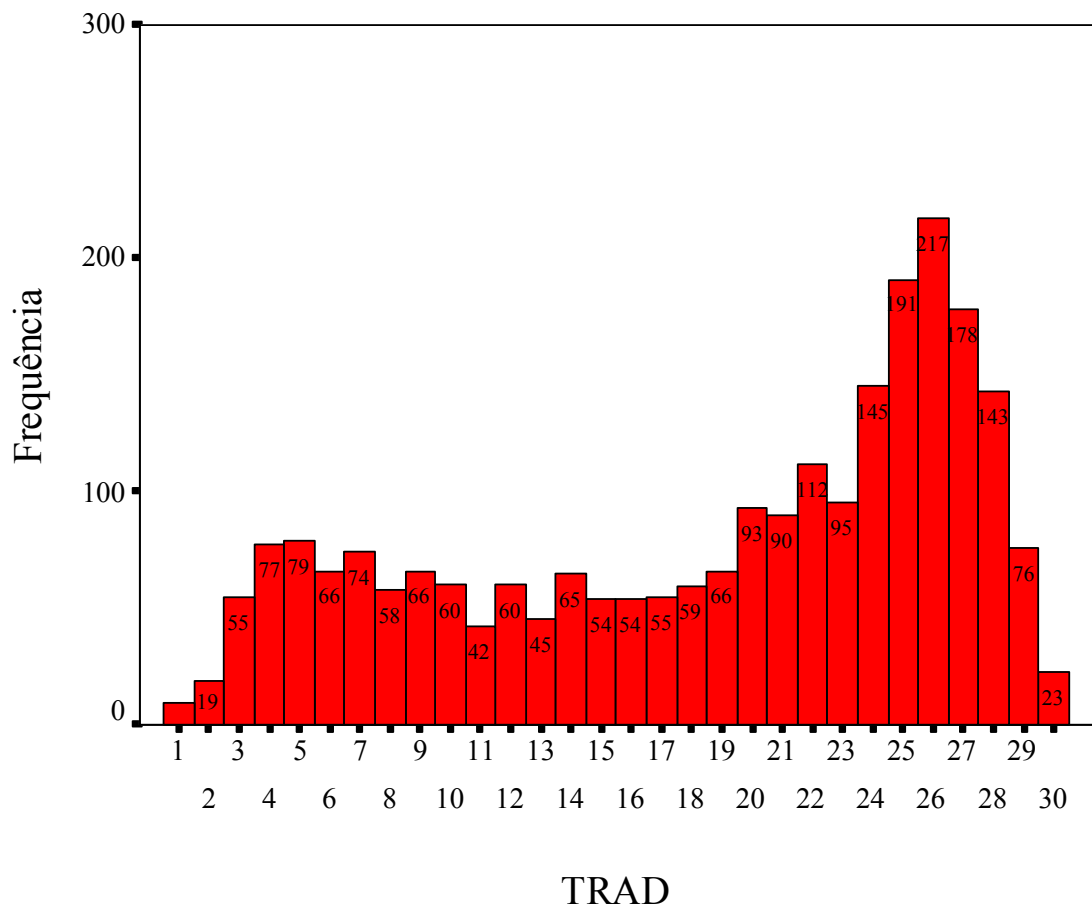


Figura 2-2. Histograma da distribuição dos dados da tabela 2-2

Você pode também ilustrar os dados com um polígono em lugar de um histograma, como faz a figura 2-3.

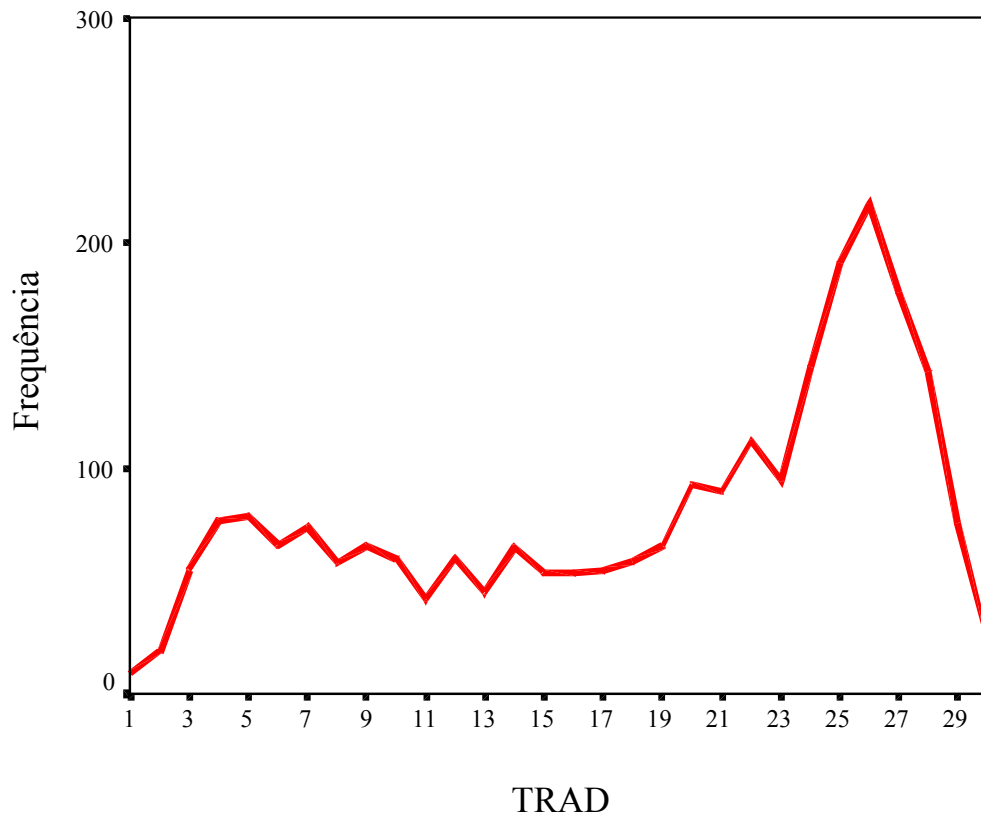


Figura 2-3. Polígono da distribuição dos dados da tabela 2-2

A representação da frequência acumulada aparece como na figura 2-4. A linha contínua vermelha representa essa curva para o TRAD. A curva pontilhada representa a mesma informação no caso em que os escores do teste se distribuíssem exatamente numa curva normal, com média de 18,50 (que é a média geral deste teste).

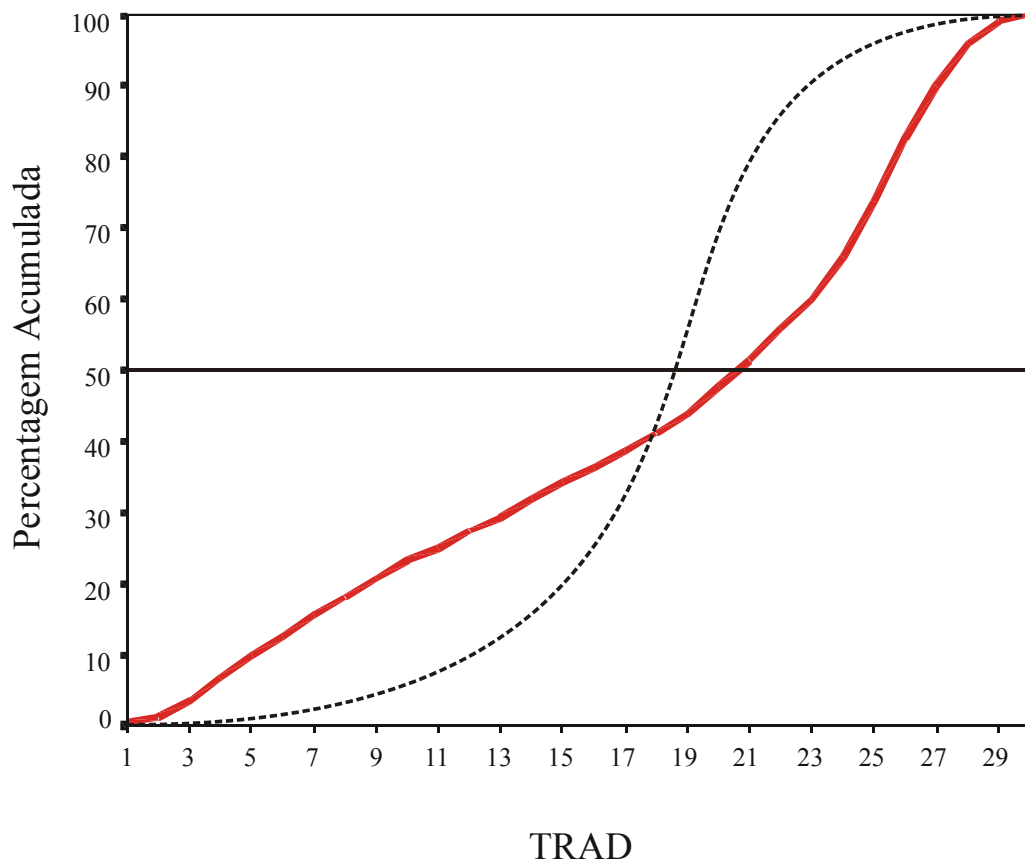


Figura 2-4. Curva da frequência acumulada dos dados da tabela 2-2

Se achar que 30 escores é uma amplitude muito grande, posso agrupá-los em categorias, digamos escores em intervalos de 5 em 5, e produzir uma nova distribuição de frequências, agora chamada de distribuição de frequência de dados agrupados, como está expresso na tabela 2-3.

Tabela 2-3. Distribuição de frequência de escores agrupados do TRAD de 2.426 sujeitos

Escore (X)	Frequência (f)	Porcentagem (%)	Porcentagem acumulada
1 - 5	239	9,9	9,9
6 - 10	324	13,4	23,2
11 - 15	266	11,0	34,2
16 - 20	327	13,5	47,7
21 - 25	633	26,1	73,7
26 - 30	637	26,3	100,0

Os gráficos que representam essa distribuição são os das figuras 2-5 (histograma) e 2-6 (polígono).

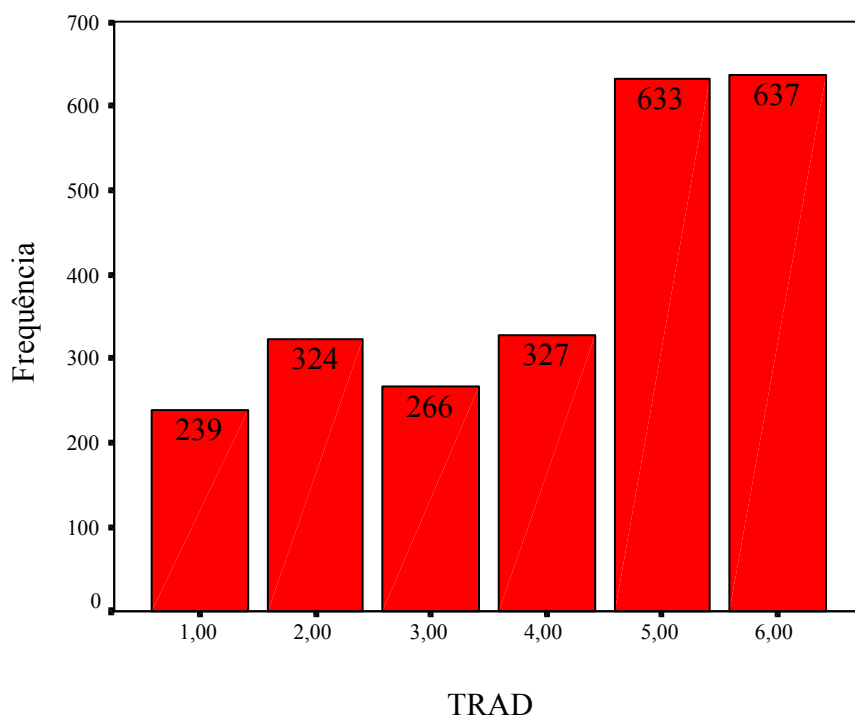


Figura 2-5. Histograma da distribuição dos dados da tabela 2-3

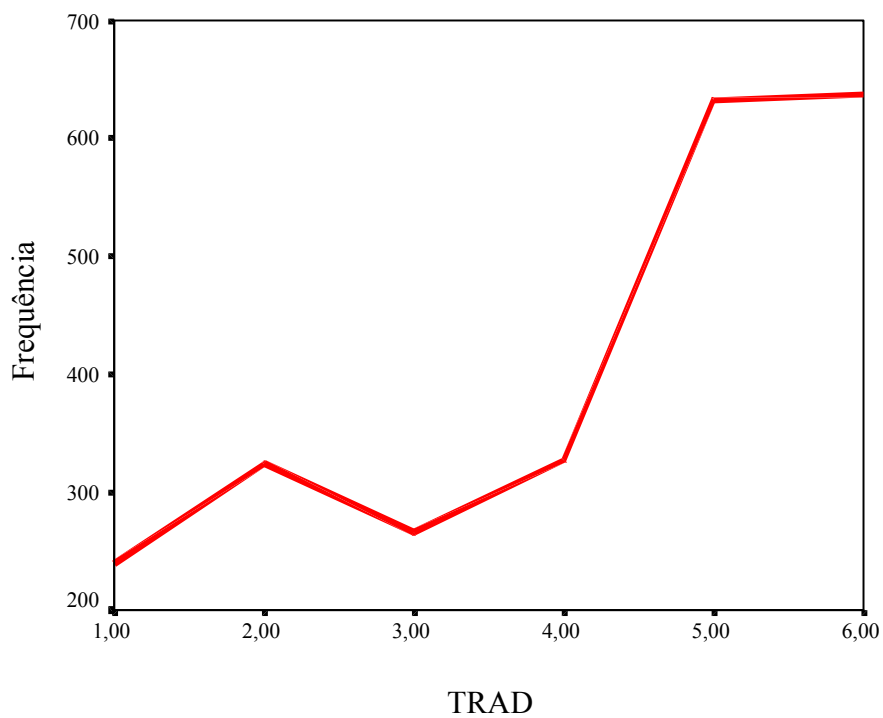


Figura 2-6. Polígono da distribuição dos dados da tabela 2-3

## II – Medidas de Tendência Central

A expressão dos dados de uma pesquisa por meio da distribuição de frequência já dá bastante informação sobre a mesma, mas ela é de manipulação muito pesada. Embora ela dê uma informação completa sobre os dados, fica difícil detalhar esta informação. Por exemplo, se você quer saber se este grupo de 2.426 sujeitos é forte ou fraco na aptidão que o TRAD mede ou se você quer saber se os homens são superiores ou inferiores às mulheres na mesma aptidão, como é que você faria para obter tal informação de uma distribuição de frequências? No caso da comparação entre homens e mulheres, você poderia fazer duas distribuições de frequências, uma para cada sexo, e comparar as distribuições. Mesmo assim a informação fica difícil de ser vista. Para casos como esses, seria de grande valia se eu pudesse expressar com uma única palavra a informação toda contida numa distribuição de frequência. Assim, seria muito mais fácil qualquer comparação que quisesse fazer com os dados da pesquisa. Claro, seria uma informação sumária chamada de parcial. Mesmo sendo parcial, a informação deve ser o mais representativa possível de todos os dados da pesquisa. Os estatísticos inventaram uma maneira para, precisamente, dar tal informação sumária e representativa, não com uma palavra, mas com duas, as quais chamaram de medidas de tendência central e medidas de variabilidade. As medidas de tendência central procuram dar um valor que seja o mais típico possível de toda a distribuição de valores e as medidas de variabilidade procuram dizer quão bem ou mal tal valor representa todos os valores.

Você já percebeu, também, que estamos falando de medidas e não de medida, porque há diferentes maneiras de expressar tanto o valor típico quanto a variabilidade dele. Esse fato está ligado à história que discutimos acima sobre o nível das escalas de medida (se lembra, escala nominal, ordinal etc.). Assim, há três tipos de medidas de tendência central: a média, a mediana e a moda. A média se usa com escalas que chegam a ser, pelo menos, do tipo intervalar; a mediana para escalas ordinais; e a moda para escalas nominais. Isso tem como consequência que, quanto mais baixo o nível da escala de medida, menor será a informação que a estatística de tendência central produz sobre os dados originais. Assim, a média é a medida parcial que



produz a maior informação; mas ela exige que você esteja medindo o seu objeto de estudo com escalas intervalares, pelo menos.

## 1 – A Média Aritmética<sup>2</sup>

Esta medida estatística consiste, então, em se procurar aquele valor que melhor representa todos os valores da pesquisa. Tal façanha é conseguida pela soma de todos os valores, dividindo a mesma pelo número de valores somados. Os estatísticos expressam tal processo de várias formas, a saber:

Se você tiver uma distribuição contínua dos dados:

$$M = \bar{X} = \frac{\Sigma X}{N} \quad (2.1)$$

ou, se você tiver uma distribuição agrupada de dados, como nas tabelas 2-2 e 2-3,

$$M = \bar{X} = \frac{\Sigma fX}{N} \quad (2.2)$$

onde

X = os valores empíricos (2.426 deles, no nosso exemplo);

N = número de sujeitos da pesquisa (2.426, no nosso exemplo);

f = frequência de ocorrência de um dado valor ou classe de valores.

Exemplificando com a tabela 2-3:

---

<sup>2</sup> Fala-se que a média é o primeiro momento estatístico e a variância, o segundo. Essa história vem do seguinte:

Número	Momento bruto	Momento central	Cumulativo
0	1	0	
1	$\mu$	0	$\mu$
2	$\mu^2 + \sigma^2$	$\sigma^2$	$\sigma^2$
3	$\mu^3 + 3\mu\sigma^2$	0	0
4	$\mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$	$3\sigma^4$	0

No caso de dados agrupados em classes, para o valor de X você escolhe o ponto médio do intervalo da respectiva classe. Assim, para a classe 26 – 30, o ponto médio é 28.

Escore (X)	Ponto médio (PM)	Frequência (f)	fPM
1 – 5	3	239	717
6 – 10	8	324	2.592
11 – 15	13	266	3.458
16 – 20	18	327	5.886
21–25	23	633	14.559
26 –30	28	637	17.836
Soma		2.426	45.048

Assim, a média dessa distribuição será:

$$M = \frac{\sum fX}{N} = \frac{45.048}{2.426} = 18,57.$$

Assim, o valor mais típico na TRAD é de 18 pontos e meio. Isso significa que todos os 2.426 sujeitos da amostra deveriam ter obtido este escore no teste.

**Uma digressão:** A média é chamada em estatística de primeiro momento. Vamos ver o que exatamente isso significa. Os estatísticos definiram momentos para todas as distribuições. Essa história está ligada ao difícil problema do teorema do limite central. Sem entrar na discussão desse teorema, entenda momentos estatísticos como análogos aos momentos em física, onde ele é definido como uma força multiplicada pela sua distância a partir do fulcro ou de um ponto central de referência, chamado de centróide. Também, em estatística os momentos são definidos em função de um ponto de referência, quais sejam o zero ou a média. Assim, o primeiro momento estatístico é a média, que consiste na *soma das distâncias* a partir do zero, vezes a probabilidade de estar naquela distância; ela é, assim, chamada de o valor esperado de X ou a *expectância* de X, ou seja, E(X). O segundo momento é a variância, porque ela é a *soma das distâncias ao quadrado* a partir da média, vezes a probabilidade de estar naquela distância. Momentos de ordem superior são a assimetria e a curtose, onde as distâncias com

respeito à média são elevadas à terceira e quarta potência, respectivamente. Assim, temos o seguinte:

Momento	Momento bruto	Momento central	Cumulativo
0	1	0	
1	$\mu$	0	$\mu$
2	$\mu^2 + \sigma^2$	$\sigma^2$	$\sigma^2$
3	$\mu^3 + 3\mu\sigma^2$	0	0
4	$\mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$	$3\sigma^4$	0

## 2 – A Mediana

Os autores não se entendem muito bem sobre o que seja exatamente a mediana de uma distribuição. Contudo, a definição mais comumente utilizada é a de que a mediana constitui naquele valor da distribuição de escores abaixo do qual caem 50% de todos os escores, isto é, ela é o valor que fica no meio da distribuição. Uma fórmula para estimar a mediana é a seguinte:

$$\text{Mediana} = L_s + \frac{fa - 0,5N}{f} \quad (2.3)$$

onde

$L_s$  = limite superior da classe onde caem os 50% dos casos;

$fa$  = frequência acumulada da classe onde cai a mediana;

$f$  = frequência dentro da classe em que cai a mediana.

Veja os dados do nosso exemplo da tabela 2-3:

Escore (X)	Frequência (f)	Porcentagem (%)	Porcentagem acumulada (fa)
1 - 5	239	9,9	9,9
6 - 10	324	13,4	23,2
11 - 15	266	11,0	34,2
16 - 20	327	13,5	47,7
21 - 25	633	26,1	73,7
26 - 30	637	26,3	100,0
Total	2.426		

A classe onde cai a mediana é a 16 – 20; a  $f_a = 47,70$ ; o  $f = 327$ . Assim,

$$\begin{aligned} \text{Mediana} &= L_s + \frac{f_a - 0,50N}{f} = 20 + \frac{47,7 - 0,50 \times 2.426}{327} = 20 + \frac{-1213}{327} = \\ &= 20 - 3,71 = 16,29 \end{aligned}$$

### 3 – A Moda

A moda simplesmente diz qual o escore ou qual a classe de escores que tem mais sujeitos. Ou seja, a moda é aquele valor que ocorre mais vezes numa dada distribuição; é o valor mais freqüente. Acontece muitas vezes que há mais de um escore ou classe em tal situação, como acontece com o nosso exemplo onde a moda pode ser tanto o escore 23 quanto o 28 ou, ainda, as classes 21-25 e 26-30; isso acontece porque os dados do nosso exemplo se apresentam como uma curva com dois picos na extrema direita. A informação que a moda dá é importante e, em caso de escalas nominais é a única possível, mas, como você vê, ela é uma informação bastante grosseira.

## III – Medidas de Variabilidade

As medidas de variabilidade ou de dispersão objetivam dar um significado mais preciso às medidas de tendência central; isto é, elas procuram caracterizar quão representativas dos dados originais estas medidas de tendência central o são; ela informa se os dados se concentram em torno da medida de tendência central ou se eles se espalham ou se afastam dela. Como previsível, as medidas de dispersão dependem do tipo de medida de tendência central que se utiliza para descrever os dados. Assim temos: (1) a amplitude (a única viável com escalas nominais, onde se usa a moda como medida de tendência central), (2) o intervalo semi-interquartil (para escalas ordinais, onde se usa a mediana) e (3) a variância (para escalas de intervalo, onde se usa a média).

## 1 – A Amplitude

A amplitude é a medida mais simples da dispersão de uma série de dados. Ela consiste simplesmente na diferença entre o escore mais alto ( $X_s$ ) e o mais baixo ( $X_i$ ):

$$\text{Amplitude} = X_s - X_i \quad (2.4)$$

No nosso exemplo (veja tabela 2-2), os escores vão de 1 a 30 (ninguém ganhou 0). Assim, a amplitude dos dados será de 29 ( $30 - 1 = 29$ ).

Então, você vê que essa medida é bastante precária, porque usa somente dois valores da série de dados, a saber, os extremos, que podem variar enormemente. Mas com escalas nominais não há outra maneira de dar a informação da variabilidade. Mesmo assim, existem estatísticas não-paramétricas que fazem uso da amplitude para tirar conclusões dos dados da amostra para a população.

## 2 – O Intervalo Semi-Interquartilico (Q)<sup>3</sup>

Esta estatística é baseada no cálculo dos percentis da distribuição de frequência dos dados da pesquisa. Ela consiste na metade da diferença entre os escores que representam o percentil 75 e o percentil 25, isto é, o quartil 3 ( $Q_3$ ) e o quartil 1 ( $Q_1$ ); entre esses dois quartis cai a metade de todos os escores da pesquisa. Assim, ela visa dar uma idéia da dispersão dos dados em torno da mediana, esta última sendo o quartil 2 ou percentil 50; você vê que essa medida não depende mais dos escores extremos como dependia a amplitude. Sua fórmula é:

$$Q = \frac{Q_3 - Q_1}{2} \quad (2.5)$$

Veja o caso do nosso exemplo com o TRAD da tabela 2-2. Os quartis aparecem ali como:

$$Q_1 = 12$$

$$Q_2 = 21$$

---

<sup>3</sup> Alguns autores abreviam esse escore como IQ, ou seja, intervalo interquartilico

$$Q_3 = 26$$

Assim,

$$Q = (26 - 12)/2 = 7.$$

Esse dado significa que 25% (o  $Q$  é a metade de 50%) dos escores se situam a 7 escores acima e abaixo da mediana, ou seja,  $21 \pm 7$ , isto é, entre os escores 14 e 28.

### 3 – A Variância

A variância leva em conta todos os dados da pesquisa e informa quanto cada um desses escores se distancia do escore típico da distribuição, isto é, da média da distribuição. De fato, ela não dá a distância individual de cada escore em relação à média, mas ela dá a *média de todas as distâncias* de todos os escores com relação à média da distribuição, como se dissesse: em média, os escores distam tanto e tanto da média da distribuição. Assim, nesse contexto precisamos falar de duas coisas, a saber, os desvios e os desvios médios ou a variância.

#### 3.1 – Os Desvios ( $x$ )

Este conceito de distância, chamado de *desvio*, é muito importante em estatística. Ele é tipicamente representado pela letra minúscula correspondente à letra maiúscula que expressa o escore bruto. Assim, se o escore bruto é representado por  $X$ , o desvio será expresso por  $x$ . A fórmula do desvio é:

$$x = X - M \tag{2.6}$$

Agora, a soma de todos os desvios em torno da média aritmética da distribuição dos escores vai dar zero (0), isto é,

$$\Sigma x = 0 \tag{2.7}$$

e, conseqüentemente, também a média dos desvios vai dar 0, ou seja,

$$MD = \frac{\Sigma(X - M)}{N} = 0 \quad (2.8)$$

Isso acontece porque você soma os valores algébricos dos desvios, onde uns são positivos e outros negativos, de sorte que, no final das contas, eles se anulam mutuamente. A questão, porém, é que um desvio positivo é desvio, como também um desvio negativo também é desvio; assim, os dois não podem se anular mutuamente. Para evitar esse desfecho, a soma é feita sobre os valores *absolutos* dos desvios, ou seja,

$$MD = \frac{\Sigma|X - M|}{N} \quad (2.9)$$

Entretanto, como o uso de valores absolutos impede outros cálculos estatísticos (restrição matemática), costuma-se trabalhar com o quadrado dos desvios, ou seja, a soma do quadrado dos desvios ou a soma dos desvios quadráticos, abreviada como SQ (soma dos quadrados ou soma quadrática) ou SS (do inglês, *sum of squares*). Como o quadrado de valores positivos ou negativos sempre dá um valor positivo, os desvios ao quadro evitam que eles se anulam no momento de serem somados. A fórmula dos desvios quadráticos é a seguinte:

$$SQ = \Sigma x^2 = \Sigma(X - M)^2 \quad (2.10)$$

Quando se trata de dados numa distribuição de frequência, então os desvios quadráticos devem ser multiplicados pela respectiva frequência em que ocorrem. Assim, a fórmula 2.10 se torna:

$$SQ = \Sigma fx^2 \quad (2.11)$$

Outras fórmulas que você vai encontrar e que dão os mesmos resultados são:

$$SQ = \Sigma X^2 - \frac{(\Sigma X)^2}{N}$$

$$SQ = \Sigma X^2 - NM^2$$

*Uma nota:* Os estatísticos se referem à quantidade  $\Sigma X^2$  como a soma não-corrigida dos quadrados e à quantidade SQ ou  $\Sigma x^2$  como a soma corrigida. A diferença aritmética entre a soma não-corrigida e a soma corrigida dos quadrados (isto é, o termo  $NM^2$ ) é referida como a correção para a média ou o fator de correção.

### 3.2 – A Variância ( $s^2$ )

A SQ tipicamente dá um número enorme e a gente não sabe bem o que fazer com ele para tirar alguma informação sobre a variabilidade dos escores em torno da média. Por isso, os estatísticos inventaram a variância, que consiste em expressar essa SQ em termos de todos os escores para obter uma estimativa média da variabilidade em torno da média da distribuição dos escores. Consegue-se tal efeito, dividindo a SQ pelo número de sujeitos (N). Sua fórmula é:

$$s^2 = \frac{SQ}{N} = \frac{\Sigma x^2}{N} = \frac{\Sigma(X - M)^2}{N} \quad (2.12)$$

Novamente, quando se trata de dados numa distribuição de freqüência, então os desvios quadráticos devem ser multiplicados pela respectiva freqüência em que ocorrem. Assim, a fórmula 2.12 se torna:

$$s^2 = \frac{\Sigma fx^2}{N} \quad (2.13)$$

Quando a amostra de sujeitos é pequena, costuma-se substituir o N por N-1 para tornar a estatística da variância não-enviesada, dizem os estatísticos (este N-1 é conhecido como graus de liberdade, que será tratado no capítulo 6).

Vamos dar um exemplo com os dados da tabela 2-2 (veja a tabela 2-4).



Tabela 2-4. Cálculo da variância com os escores no TRAD de 2.426 sujeitos

Escore (X)	Frequência (f)	fX	Desvios (x)	Desvios Quadráticos (x <sup>2</sup> )	fx <sup>2</sup>
1	9	9	-17,59	309,41	2784,67
2	19	38	-16,59	275,23	5229,33
3	55	165	-15,59	243,05	13367,65
4	77	308	-14,59	212,87	16390,84
5	79	395	-13,59	184,69	14590,36
6	66	396	-12,59	158,51	10461,53
7	74	518	-11,59	134,33	9940,28
8	58	464	-10,59	112,15	6504,59
9	66	594	-9,59	91,97	6069,89
10	60	600	-8,59	73,79	4427,29
11	42	462	-7,59	57,61	2419,54
12	60	720	-6,59	43,43	2605,69
13	45	585	-5,59	31,25	1406,16
14	65	910	-4,59	21,07	1369,43
15	54	810	-3,59	12,89	695,96
16	54	864	-2,59	6,71	362,24
17	55	935	-1,59	2,53	139,05
18	59	1.062	-0,59	0,35	20,54
19	66	1.254	0,41	0,17	11,09
20	93	1.860	1,41	1,99	184,89
21	90	1.890	2,41	5,81	522,73
22	112	2.464	3,41	11,63	1302,35
23	95	2.185	4,41	19,45	1847,57
24	145	3.480	5,41	29,27	4243,87
25	191	4.775	6,41	41,09	7847,83
26	217	5.642	7,41	54,91	11915,06
27	178	4.806	8,41	70,73	12589,60
28	143	4.004	9,41	88,55	12662,38
29	76	2.204	10,41	108,37	8235,98
30	23	690	11,41	130,19	2994,33
Total	2.426	45.089			163.142,71
Média	fX/N = 18,59				

Assim, a variância será:

$$s^2 = \frac{fx^2}{N} = \frac{163.142,71}{2.426} = 67,25$$

O valor 67,25 representa a distância média em que todos os escores se situam da média de 18,59. É difícil perceber o que exatamente esse enorme valor quer informar sobre os dados, inclusive, porque ele não representa uma distância e, sim, uma área (isso será explicado logo mais). Por isso, a variância é tipicamente transformada no desvio padrão (DP), que é simplesmente a raiz quadrada da variância, ou seja,

$$DP = \sqrt{s^2} = 8,20$$

O desvio padrão é mais fácil de entender, porque ele está expresso na mesma escala de medida dos escores do teste. Esse valor de 8,20 significa que a maior parte dos escores (na verdade 68,26% deles – veremos isto quando falarmos da curva normal) se situa entre a média e mais 1DP e menos 1DP, ou seja:  $M \pm DP = 18,59 \pm 8,2$ , que corresponde aos escores extremos de 10,39 e 26,79.

A variância é difícil de entender, porque ela representa uma área, enquanto o DP representa um intervalo. Veja a figura 2-7 para visualizar essa asserção.

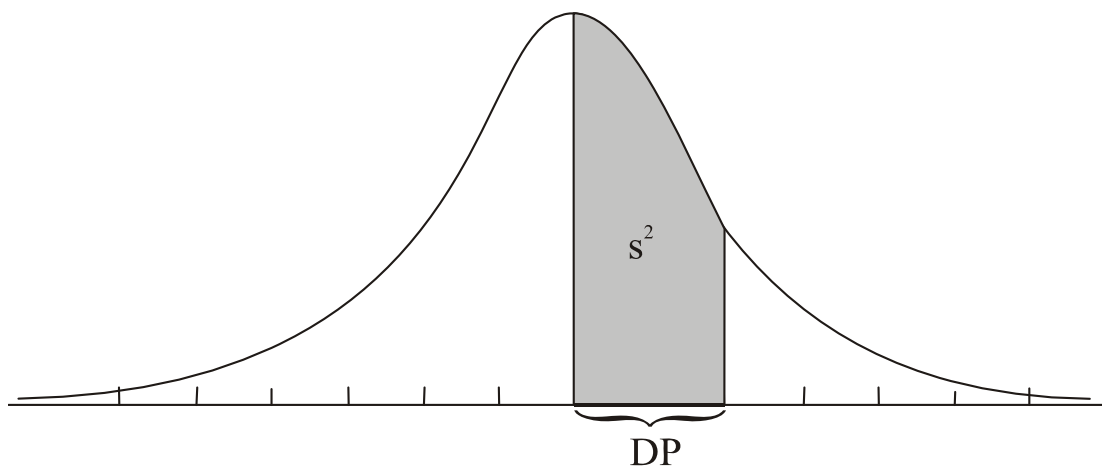


Figura 2-7. Ilustração da variância e do desvio padrão

## IV – Escore Padrão (z)

As medidas de tendência central, em particular a média, procuram expressar uma série de números num único valor. Essa série de números constitui um escala que tem uma origem, quase sempre

arbitrária (só a escala de razão tem origem natural), como é o caso do teste TRAD, que inicia no zero (este 0 é arbitrário) e vai até 30. Assim, um escore qualquer nessa escala é expresso em termos dos demais escores da escala. Mas, como temos a média que expressa todos os escores, pergunta-se por que não fazer dessa média a origem da escala? Assim, a média ficaria a origem e os escores todos se distribuiriam em torno dela, uns acima dela e outros abaixo dela. Pois esta é a idéia dos escores padronizados ou escores-padrão, a saber, expressar os escores da escala original, que inicia no 0 e acaba no 30, no caso do TRAD, iniciar os escores com a média, que no caso seria 18,59. Este último valor, então, seria a origem da escala, recebendo o valor 0. Qual a vantagem dessa manobra? A vantagem é enorme, porque qualquer que seja a escala original (deve ser de nível intervalar, pelo menos), sua transformação numa escala padrão torna todas as escalas diretamente comparáveis. Assim, escalas de qualquer tamanho, por exemplo, de 10, de 30, de 100, 523 etc., e que comecem não importa com que número, se transformadas em escalas padrão, todas elas serão idênticas, indo de  $-\infty$ , passando pela média (0), até  $+\infty$ , sendo a média 0 o ponto de referência para todas elas.

Esta transformação da escala de medida em escala padrão leva em conta a média e o desvio padrão da escala original, sendo sua fórmula a seguinte:

$$z = \frac{X - M}{DP} \quad (2.14)$$

Você vê, então, que o escore individual na escala original é expresso em relação à média da mesma escala (M) e em função do DP dela. Essa é uma transformação linear; por isso, o  $z$  terá as mesmas características da escala original, isto é, se esta escala era normal, a nova escala  $z$  também o será, mas se aquela não era normal, esta também não o será. Isso é importante observar, porque a distribuição dos escores na escala  $z$  não precisa ser normal, como erroneamente se pensa; o escore  $z$  não está ligado à curva normal padronizada, como veremos ao falarmos desta; ele é simplesmente uma transformação de uma escala, qualquer escala (de intervalo!), com a única intenção de tornar a média, a origem da escala.

Vamos exemplificar com o TRAD, observando a tabela 2-5.

Tabela 2-5. Escores brutos transformados em z

Escore (X)	Frequência (f)	fX	z	T	CEEB	Desvio QI
1	9	9	-2,15	28,55	285,49	67,82
2	19	38	-2,02	29,77	297,68	69,65
3	55	165	-1,90	30,99	309,88	71,48
4	77	308	-1,78	32,21	322,07	73,31
5	79	395	-1,66	33,43	334,27	75,14
6	66	396	-1,54	34,65	346,46	76,97
7	74	518	-1,41	35,87	358,66	78,80
8	58	464	-1,29	37,09	370,85	80,63
9	66	594	-1,17	38,30	383,05	82,46
10	60	600	-1,05	39,52	395,24	84,29
11	42	462	-0,93	40,74	407,44	86,12
12	60	720	-0,80	41,96	419,63	87,95
13	45	585	-0,68	43,18	431,83	89,77
14	65	910	-0,56	44,40	444,02	91,60
15	54	810	-0,44	45,62	456,22	93,43
16	54	864	-0,32	46,84	468,41	95,26
17	55	935	-0,19	48,06	480,61	97,09
18	59	1.062	-0,07	49,28	492,80	98,92
19	66	1.254	0,05	50,50	505,00	100,75
20	93	1.860	0,17	51,72	517,20	102,58
21	90	1.890	0,29	52,94	529,39	104,41
22	112	2.464	0,42	54,16	541,59	106,24
23	95	2.185	0,54	55,38	553,78	108,07
24	145	3.480	0,66	56,60	565,98	109,90
25	191	4.775	0,78	57,82	578,17	111,73
26	217	5.642	0,90	59,04	590,37	113,55
27	178	4.806	1,03	60,26	602,56	115,38
28	143	4.004	1,15	61,48	614,76	117,21
29	76	2.204	1,27	62,70	626,95	119,04
30	23	690	1,39	63,91	639,15	120,87
Total	2.426	45.089				

Média:  $M = 45.089/2.426 = 18,59$

Desvio padrão:  $DP = 8,2$  (veja tabela 2.4)

Assim, para o escore 30, o escore padrão será:

$$z = (30 - 18,59)/8,2 = 1,39; \text{ etc. (veja o restante na tabela 2-5).}$$

O que pode parecer estranho nesta transformação dos escores originais em escores-padrão é o fato de que esses, em parte, se apresentam com sinais negativos, isto é, aqueles escores originais que estão abaixo da média; além de aparecerem com decimais. Para superar essas dificuldades ou deselegâncias, tipicamente os escores-padrão são novamente transformados em uma escala padronizada mais elegante, que evita sinais negativos e, inclusive, as decimais. Dessas transformações cosméticas temos uma série delas, tais como a escala T, a escala CEEB (*Certificate of Entrance Examination Board*) e o desvio QI. Essas transformações seguem a seguinte equação linear:

$$\text{Escore transformado} = a + bz \tag{2.15}$$

Onde

$a$  e  $b$  são constantes quaisquer.

Algumas dessas transformações se tornaram clássicas, dando valores universalmente aceitos para o  $a$  e para o  $b$ . Veja algumas delas:

$$T = 50 + 10z$$

$$\text{CEE}B = 500 + 100z$$

$$\text{Desvio QI} = 100 + 15z.$$

Assim, para o escore bruto de 30, essas transformações do  $z$  serão (veja resultados na tabela 2-5):

$$T = 50 + 10 \times 1,98 = 69,8$$

$$\text{CEE}B = 500 + 100 \times 1,98 = 698$$

$$\text{Desvio QI} = 100 + 15 \times 1,98 = 144,70.$$

Se arredondar os valores dessas transformações, você terá valores sem decimais. Observe que o T e o desvio QI, para evitar as vírgulas devem ser arredondados com mais violência, enquanto o CEEB já evita as decimais mais facilmente. São pequenas conveniências!

**Nota:** todas essas transformações são lineares. Há também transformações não-lineares, quando o  $z$  é calculado por meio da curva normal e não pela fórmula 2.14. Veremos isso ao falarmos da curva normal.

## V – O Erro Padrão (EPM)

O conceito de erro padrão, também chamado de erro amostral, introduz um conceito muito importante em toda a estatística inferencial; por isso é útil ser discutido aqui, porque ele vai dizer algo relevante sobre a média da distribuição de casos, como viemos falando até agora.

Veja: no caso do TRAD, nós selecionamos uma amostra de 2.426 sujeitos e em cima dos dados obtidos deles fizemos todas as descrições acima relatadas. Agora, esses sujeitos representam apenas uma amostra possível de uma população de sujeitos similares. Eu bem que podia ter selecionado outros sujeitos dessa população e, assim, teria uma outra amostra, digamos de 3.000 casos. Com essa nova amostra, posso fazer todas as descrições dos dados novamente, inclusive da média e da variância e seus derivados. Provavelmente, essas descrições, a média e a variância, seriam um pouco diferentes das da amostra de 2.426 casos. Se eu for fazendo isso muitas vezes, terei no final uma série de amostras diferentes da mesma população e para todas elas posso ter também todas as estatísticas que discutimos neste capítulo (média, DP etc.). Agora, o curioso é que tendo muitas estatísticas, isto é, muitas médias, muitas variâncias etc. da mesma população, eu posso fazer uma distribuição de frequência e análises descritivas dessas estatísticas, como fiz com os escores individuais de cada amostra. Essa distribuição de estatísticas é chamada de *distribuição de amostragem (sampling distribution)*.

Assim, se eu fizer tal distribuição de médias de muitas amostras da mesma população e calcular o desvio padrão dessa distribuição, estou obtendo o desvio padrão de uma distribuição de amostras. Para distinguir este desvio padrão do desvio padrão da distribuição de casos, ele é chamado de *erro padrão*, e é simbolizado como EPM (erro padrão da média da distribuição de amostragem) ou  $\sigma_M$ , ou seja, o desvio padrão de médias.

O problema prático com tudo isso é que, para calcular o erro padrão, eu preciso do desvio padrão da população, o qual eu nunca tenho. Assim, tenho que estimar esse desvio padrão da população a partir de uma única amostra que possuo (no caso do TRAD, os 2.426 sujeitos). Felizmente, existem fórmulas para isso, sendo a seguinte a mais utilizada:

$$\sigma_M = \frac{DP}{\sqrt{N-1}} \quad (2.16)$$

isto é, para descobrir o erro padrão, basta dividir o desvio padrão da amostra de sujeitos pela raiz quadrada de  $N-1$  (ou seja, pelos graus de liberdade, veja capítulo 6 sobre os graus de liberdade).

Outra fórmula comumente utilizada para o cálculo do erro padrão é a seguinte:

$$\sigma_M = \sqrt{\frac{\Sigma x^2}{N(N-1)}} \quad (2.17)$$

O erro padrão dá a importante informação sobre quanto a média da minha amostra (de 2.426 sujeitos no caso do TRAD) se afasta da média da população da qual tirei a amostra. Veja o exemplo com o TRAD:

$$M = 18,59$$

$$DP = 8,2$$

$$N = 2.426$$

Assim, o erro padrão da média do TRAD será:

$$\sigma_M = \frac{8,2}{\sqrt{2.426-1}} = \frac{8,2}{49,24} = 0,1665.$$

Quer dizer, então, que a média da minha amostra ( $M = 18,59$ ) indica que a média da população muito provavelmente se encontra dentro do intervalo de  $18,59 \pm 0,17$  (arredondando), ou seja, entre 18,42 e 18,76. Disse muito provavelmente, porque estou utilizando 1

erro padrão em torno da média, o que corresponde a uma probabilidade de cerca de 68% (veja capítulo sobre a curva normal). Esse intervalo é chamado de intervalo de confiança da média, o qual me diz que a minha média de 18,59, de fato, pode ser qualquer outro valor contido entre os extremos desse intervalo (18,42 e 18,76).

**Nota:** O erro padrão da mediana é maior do que o da média.

De fato ele é  $\sigma_{Me} = \frac{1,253\sigma_M}{\sqrt{N}}$ . Por isso diz-se que a mediana é

menos eficiente que a média para descrever uma distribuição, porque está mais sujeita a flutuações de amostragem (veja Capítulo 5).



# Parte III

## A Descrição dos Dados da Pesquisa com o TRAD

Para responder todas as questões, levantadas no capítulo 1 sobre a pesquisa com respeito à aptidão medida pelo TRAD, é preciso ter as estatísticas descritivas de acordo com o delineamento ali proposto. Esse delineamento procurava responder questões sobre a influência de algumas variáveis sobre a aptidão medida pelo teste, a saber: (a) sexo (masculino e feminino), (b) idade (11 a 20 anos e 21 ou mais anos) e (c) nível escolar (Ensino Fundamental mais Ensino Médio e Ensino superior).

Assim, preciso dar as estatísticas descritivas dos escores no teste estratificados em termos dessas variáveis. A tabela 2-6 mostra como esses dados se apresentam.

Tabela 2-6. Descrição dos escores do TRAD por sexo, idade e nível escolar

Estatísticas	Sexo		Idade		Nível Escolar		Total
	M	F	≤ 20	≥21	F e M	Superior	
Média	18,02	18,91	16,64	21,37	15,72	22,52	18,59
DP	8,23	8,17	8,40	7,06	8,28	6,13	8,20
Erro Padrão	0,29	0,20	0,23	0,22	0,23	0,19	0,17
N	804	1.610	1.336	2.372	1.314	1.065	2.426

Nota: por causa dos missing nas diferentes variáveis, os N variam

### Nota importante (e prematura):

Nesta descrição dos dados do TRAD, aparecem muitas médias e muitas variâncias (ou DP). As diferentes variâncias que aparecem são particularmente importantes, porque elas representam o que vai se chamar de variância dentro dos grupos, variância entre os grupos e variância total. Isto é, a variância de uma população pode ser decomposta em vários pedaços. Esse dado vai ser a base do

método estatístico da análise de variância que será exposto no capítulo 8.

Esses dados distribuídos pelas variáveis de estudo podem ser representados em gráficos (polígonos), como mostram as figuras 2-8 a 2-10. Para obter esses gráficos, faça os seguintes comandos no SPSS:

GRAPHS

Line...

Line Charts: Multiple [Define]

Define Multiple Lines: *nesta janela marque*

N of Cases

Category Axis: *coloque o escore total no teste*

Define lines by: *coloque a variável a analisar (sexo, idade, escolaridade)*

OK

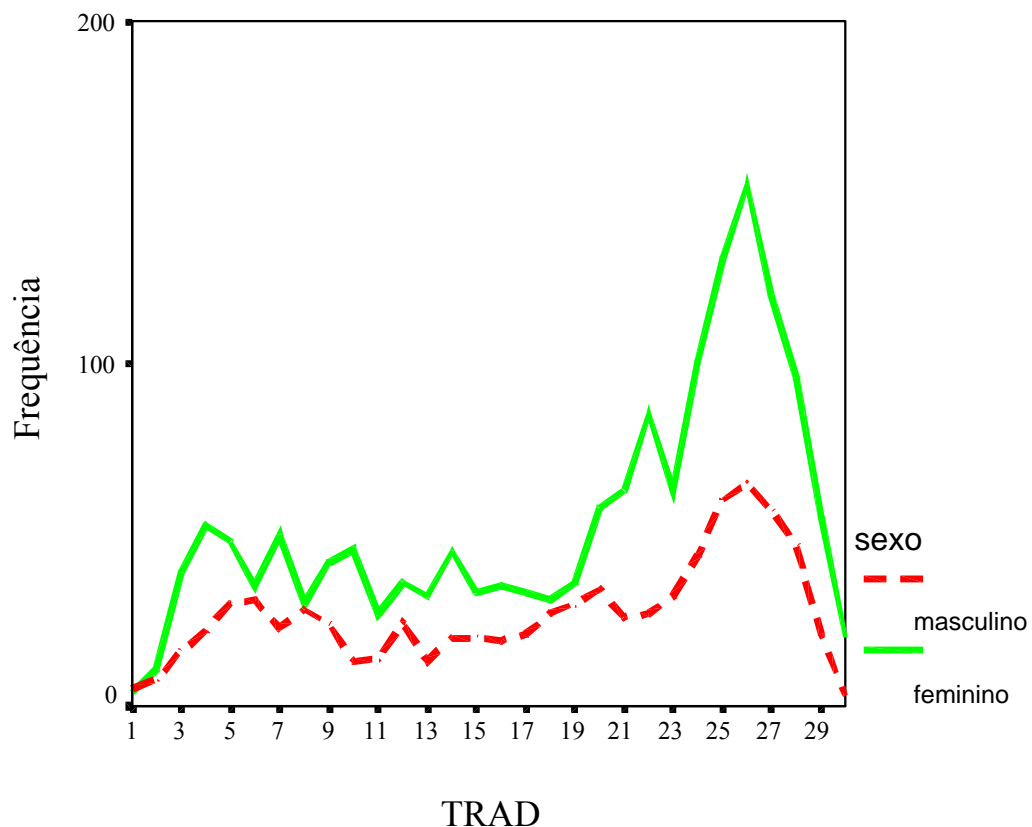


Figura 2-8. Escores no TRAD por sexo dos sujeitos

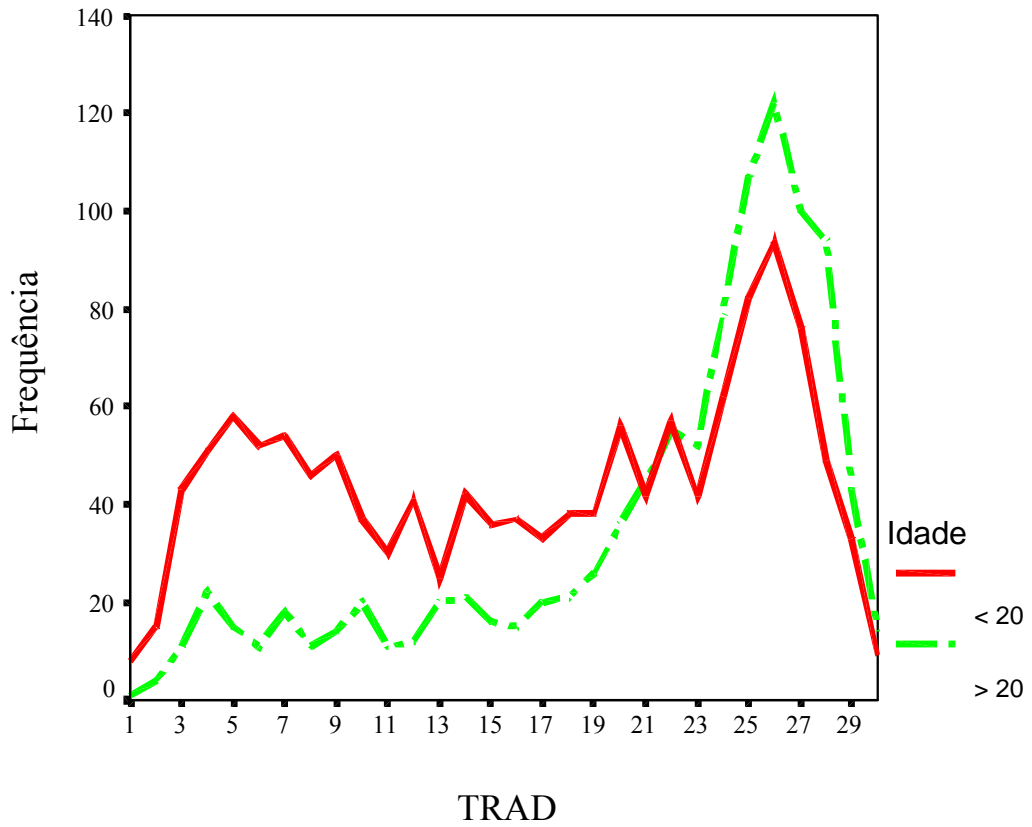


Figura 2-9. Escores no TRAD por idade dos sujeitos

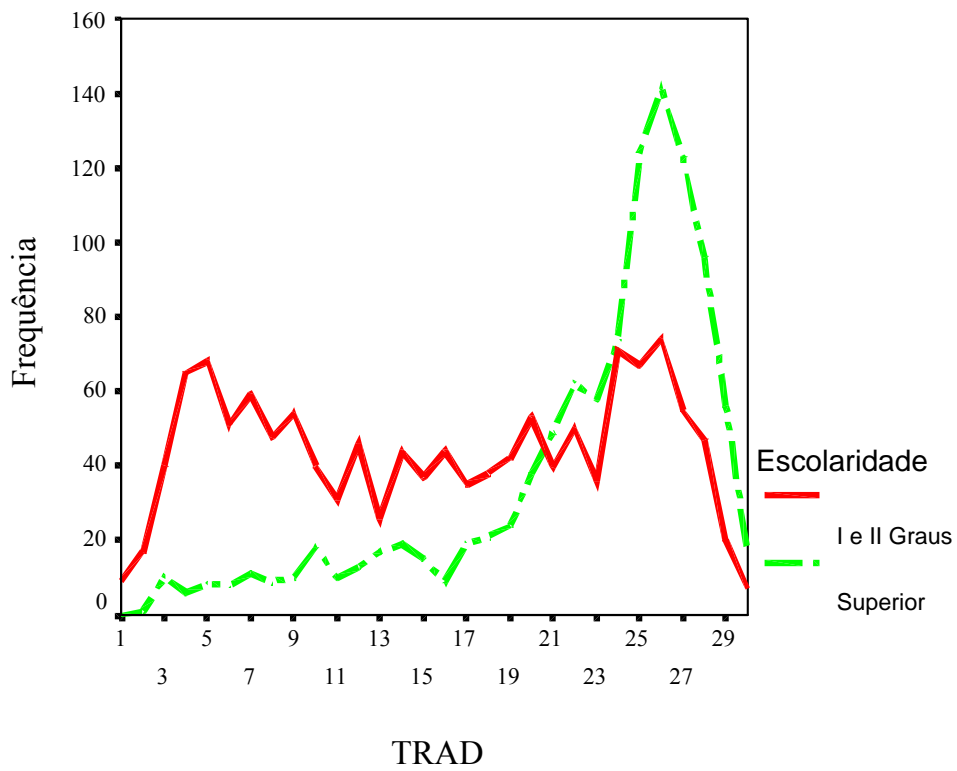


Figura 2-10. Escores no TRAD por nível escolar dos sujeitos

Você também pode ilustrar em gráficos as médias da tabela 2-6, utilizando o *Inserir Figura Gráfico* no Word, conforme mostram as figuras 2-11 a 2-13 (histogramas).

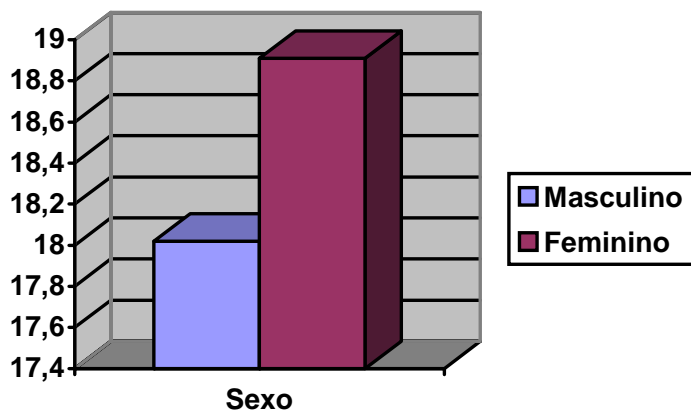


Figura 2-11. Escores médios no TRAD por sexo

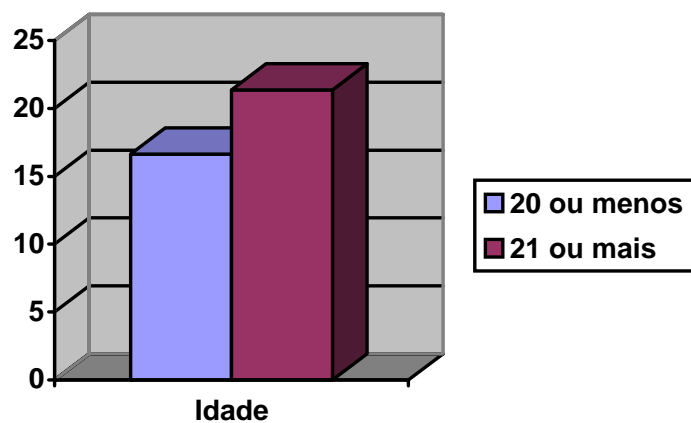


Figura 2-12. Escores médios no TRAD por idade

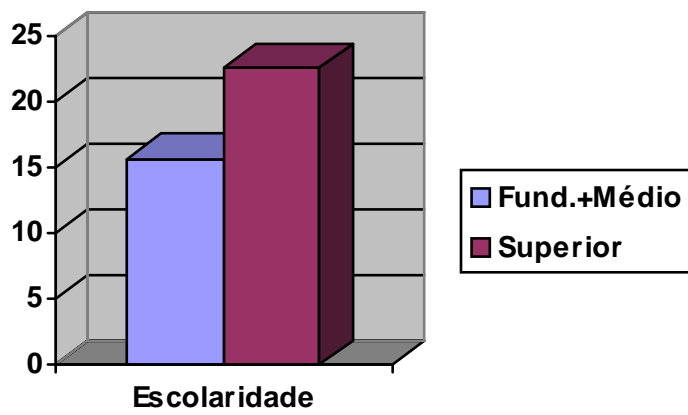


Figura 2-13. Escores médios no TRAD por escolaridade