

I. Pequena análise estatística da língua portuguesa: Machado de Assis e Pero Vaz de Caminha

Geraldo A. Barbosa

*QuantaSec Ltda. – Consultoria, Pesquisa e Projetos em Criptografia Quântica e Estatística,
Av. Portugal 1558, Belo Horizonte MG 31550-000, Brasil **

(Dated: 19 de janeiro de 2006)

Resumo: Este trabalho mostra a estatística de ocorrência de letras na língua portuguesa em textos de Machado de Assis, contabilizando mais de um milhão de letras. Apresenta-se também o bigrama de ocorrência de letras para a obra “Memórias Póstumas de Brás Cubas”. Foram calculados os valores da entropia para vários textos de Machado e feita uma breve discussão de flutuações estatísticas de acordo com o tamanho das amostragens. De maneira similar analisa-se a “Carta” de Pero Vaz de Caminha, no original e em versão contemporânea. Algumas comparações são feitas e discutidas ao final do trabalho.

INTRODUÇÃO

Dos sons guturais à representação sonora mais elaborada das idéias estabeleceu-se em passado pré-histórico a linguagem Proto-Indo-Européia. Na evolução desta estabeleceu-se a maioria das linguagens atuais. A representação da linguagem sonora através da escrita define outro salto qualitativo e intrincado da mente humana que, felizmente, deixou muitos registros preciosos para a linguística.

A complexidade da evolução de cada linguagem escrita deixa assinaturas inequívocas que permite não somente a identificação das linguagens particulares assim como suas inter-relações e até suas evoluções históricas. O advento do computador trouxe ferramentas de análise das linguagens que permitem a linguística avanços importantes que antes eram praticamente impossíveis pelo volume de dados a serem tratados.

Dentre as assinaturas mais fundamentais de cada linguagem está a ocorrência de letras e suas combinações. Para exemplificar, observemos as palavras em inglês e português “pale” e

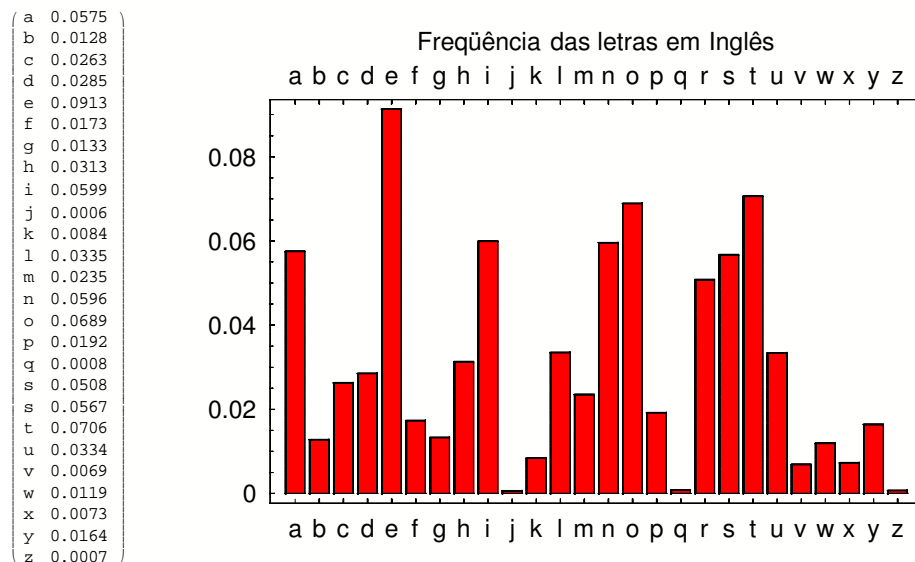


FIG. 1: Frequência das letras na língua inglesa. A tabela à esquerda da figura mostra cada letra utilizada e sua probabilidade de ocorrência no texto.

“pela” ou “train” e “intra”, ou ainda, “pious” e “pisou”, “goad” e “gado”, “blue” e “bule”. Cada uma destas palavras, em inglês e português, foram escritas com as mesmas letras. Entretanto, não é muito fácil encontrar palavras em duas línguas escritas com as mesmas letras e em igual número. Mais difícil ainda se torna quando estas palavras contém muitas letras, embora o número de permutações de letras cresça enormemente e seja igual ao fatorial do número de letras (Para N letras, existem $N!$ permutações possíveis).

Assim, a frequência da ocorrência das letras numa língua torna-se uma assinatura ou *impressão digital* – ou se utilizando a palavra em moda, o “DNA”, daquela língua. A simples existência de símbolos especiais numa língua já impossibilita a ocorrência de palavras com as mesmas letras. Assim as letras na palavra em inglês “tea” não poderão ser rearranjadas como “até” em português, pela inexistência do acento. A ocorrência de duas letras consecutivas ou bigramas (por exemplo, “pe”, “el”, “la” em “pela”), trigramas e outras combinações de ordem maior são também características determinantes de cada língua (As notações bigramas e dígramas ou dígrafos, etc. são entendidas como equivalentes).

A figura 1 e tabela ao lado mostram a ocorrência (percentual) das letras em inglês retirada de uma grande amostragem de textos. Observe-se que a letra “e” é a que mais ocorre na língua inglesa, seguida de “t” e “o” e assim por diante, nas proporções indicadas.

Qual é a frequência da ocorrência das letras em português? A resposta a esta pergunta

deve ser conhecida por muitas pessoas mas não está facilmente disponível na literatura eletrônica aberta pela Internet. A motivação para este trabalho deve-se a que o autor observou pessoas (além dele mesmo) buscando estas respostas na Internet e não terem encontrado estas estatísticas. Espera-se, portanto, que outros pesquisadores que já tenham feito estudos semelhantes sobre a língua portuguesa entendam as dificuldades para terceiros encontrarem estes trabalhos. Mais do que isto, estas dificuldades levam outros à “redescoberta da roda” e a não utilização de obras existentes. Assim, trabalhos preciosos deixam de ser compartilhados. Felizmente, isto também livra aqueles autores de quaisquer ônus nas análises aqui feitas. Estas, por sua vez, são análises comuns em estatística e linguística e podem ser reproduzidas com facilidade com “softwares” estatísticos adequados à linguagem simbólica.

Nota-se que as análises estatísticas linguísticas devem utilizar um razoável número de amostras para lhes conferirem confiabilidade como representativas daquela língua. *Qual é o número de amostras necessárias para se representar uma língua?* Este número deve ser grande e variado para não se utilizar repetições típicas de um determinado ramo do conhecimento. Uma obra como um conto genérico de maior porte talvez possa ser utilizada como representativa da língua numa determinada época. Igualmente, obras representativas de diferentes épocas podem ser utilizadas para se estudar a evolução da língua ao longo do tempo. Este trabalho pretende dar algumas respostas qualitativas e quantitativas sobre algumas destas questões.

As análises estatísticas linguísticas são tão importantes que elas fazem parte até do arsenal de grandes potências desde a Idade Média para o deciframento de mensagens dos concorrentes econômicos ou inimigos. A criptanálise tem tido papel extremamente relevante em muitos eventos de importância mundial [1]. A frequência de textos cifrados pode, em muitos casos, revelar se um ciframento é monoalfabético, a possível frequência ou dimensão da senha ou chave de ciframento e outras peculiaridades. A criptologia, juntamente com as técnicas de interceptação de mensagens, são parte do cotidiano de nações modernas [2]. Assim, as grandes potências conhecem segredos comerciais de nações antes de reuniões econômicas mundiais ou seus segredos militares tais como localizações de radares de outros países, grandes redes de transporte de energia elétrica e outras informações de importância estratégica de seu interesse. Assim, não somente os linguístas se interessam pelo estudo puro das línguas mas muitos outros vivem de seu estudo por diferentes motivações.

Este trabalho visa simplesmente fornecer um estudo de ocorrência de letras e bigra-

mas na língua portuguesa utilizando textos eletrônicos facilmente acessíveis no momento. Para uma análise estatística mais completa da língua um número bem maior e variado de textos deve ser analisado. Esta análise demanda tempo e recursos não compatíveis com o modesto objetivo do presente trabalho. Será mostrado, entretanto, que os resultados indicam convergência para índices possivelmente estacionários. Estes valores assintóticos seriam os índices representativos da língua portuguesa nesta época.

A ESCOLHA DE TEXTOS

Alguns textos foram escolhidos para análise e que pudessem mostrar alguns detalhes de diferentes épocas. Num primeiro grupo, o autor escolheu para análise alguns textos de Machado de Assis (1839/1908). Inicia-se com “Memórias Póstumas de Brás Cubas” (1881), disponível na Internet através da Fundação Biblioteca Nacional Departamento Nacional do Livro, do MEC (Ministério da Cultura – Brasil) (<http://www.machadodeassis.org.br/obras003.htm>). Para suporte das conclusões, utilizou-se os textos “Dom Casmurro” e “Quincas Borba”, ambos digitalizados pelo Núcleo de Pesquisas em Informática, Literatura e Linguística <http://www.cce.ufsc.br/alckmar/literatura/literat.html>, da Universidade Federal de Santa Catarina e, por último, deste mesmo autor utilizou-se os “Contos Fluminenses” (também disponibilizado pelo MEC, Brasil). O segundo grupo refere-se às versões moderna e antiga da “Carta” de Pero Vaz de Caminha (1450/1500), enviada a “El Rei D. Manuel”, por ocasião da descoberta do Brasil (1500). Este texto está disponível na Internet através da Biblioteca Nacional Digital (Portugal), no endereço eletrônico http://www.bnd.bn.pt/ed/viagens/brasil/obras/carta_pvcaminha/index.htm. As duas obras se distanciam em cerca de 400 anos. “Memórias” foi escrita em capítulos curtos e a “Carta” está dividida no original em páginas (frente e verso).

As estatísticas destes dois grupos serão apresentadas, com maior ou menor detalhe, de acordo com a necessidade da análise.

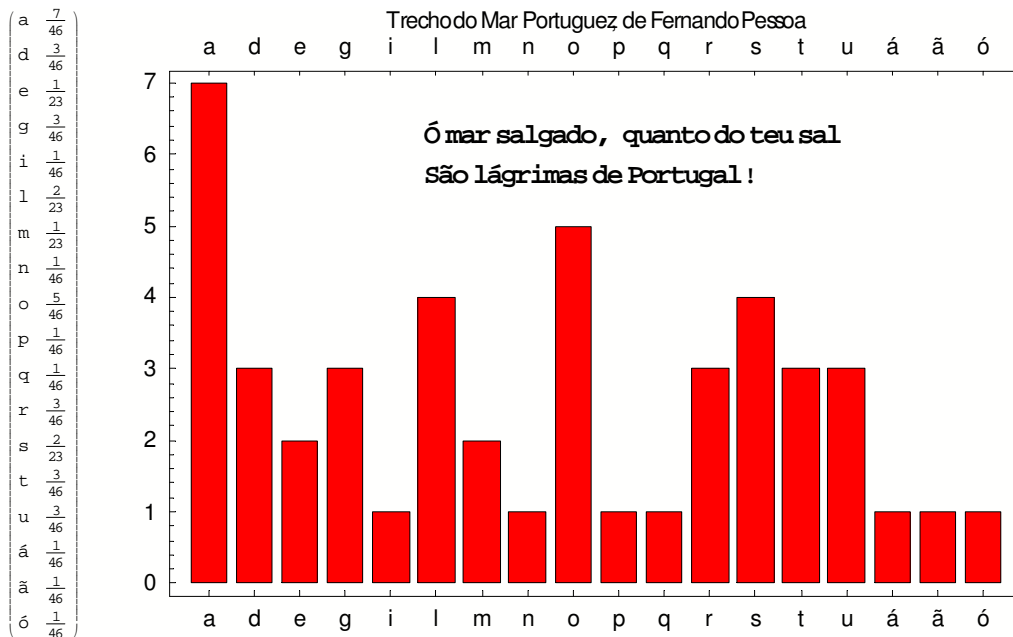


FIG. 2: Frequência das letras em trecho de poesia de Fernando Pessoa. A tabela à esquerda da figura mostra, ao lado de cada letra, sua probabilidade de ocorrência $P(x)$ em relação ao número total de letras no texto (46 letras).

METODOLOGIA

Para as análises estatísticas realizadas foram desenvolvidos programas de computação simbólica com suporte básico do “software” *Mathematica*, versão 5, da Wolfram Research. Os resultados serão espelhados na estatística de ocorrência de letras, na ocorrência de bigramas e nas entropias dos textos analisados. A entropia de cada amostragem analisada é definida pela Teoria da Informação de Shannon [3] como

$$H(X) = \sum_x P(x) \log_2 \frac{1}{P(x)}, \quad (1)$$

onde x especifica cada variável e X o conjunto destas variáveis. Aqui, X designa o alfabeto usado e x as letras do mesmo. $P(x)$ é a probabilidade de ocorrência da letra x na amostragem feita, isto é, a contagem de quantos letras x ocorreram dividida pelo número total de letras na amostragem utilizada. Muitas propriedades são associadas à entropia e se constituem em parte da Teoria da Informação. A base 2 utilizada para o logaritmo é usualmente adotada e particularmente prática para a comunicação binária (0,1). As entropias resultantes são valores em “bits”. Para exemplificar numericamente a construção

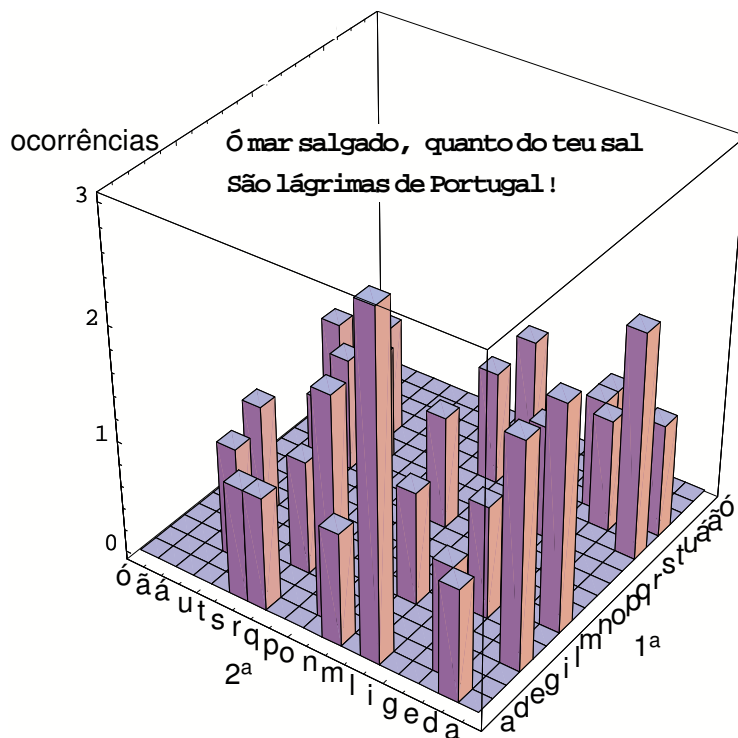


FIG. 3: Ocorrência de bigramas em trecho de poesia de Fernando Pessoa. Escolhe-se uma letra no eixo indicado 1^a e a letra seguinte no outro eixo 2^a.

das ocorrências e respectivas probabilidades, a figura 2 mostra um pequeno trecho de poesia de Fernando Pessoa com a contagem das letras. A probabilidade de cada letra está dada na lista ao lado. O mesmo trecho pode ser analisado em relação a bigramas, trigramas etc. A figura 3 mostra o bigrama do trecho escolhido. Dos bigramas pode-se retirar com facilidade as probabilidades adjuntas $P(x, y)$, assim como voltar as probabilidades $P(x)$, através de $P(x) = \sum_y P(x, y)$, etc.

MEMÓRIAS PÓSTUMAS DE BRÁS CUBAS

Frequências

As figuras e tabelas a seguir representam a frequência de ocorrência de letras do capítulo 1 ao capítulo 10 (Figura 4). As tabelas ao lado de cada figura representam numericamente o mesmo conteúdo da figura e são apresentadas para uso dos leitores que se interessam por aspectos numéricos mais detalhados.

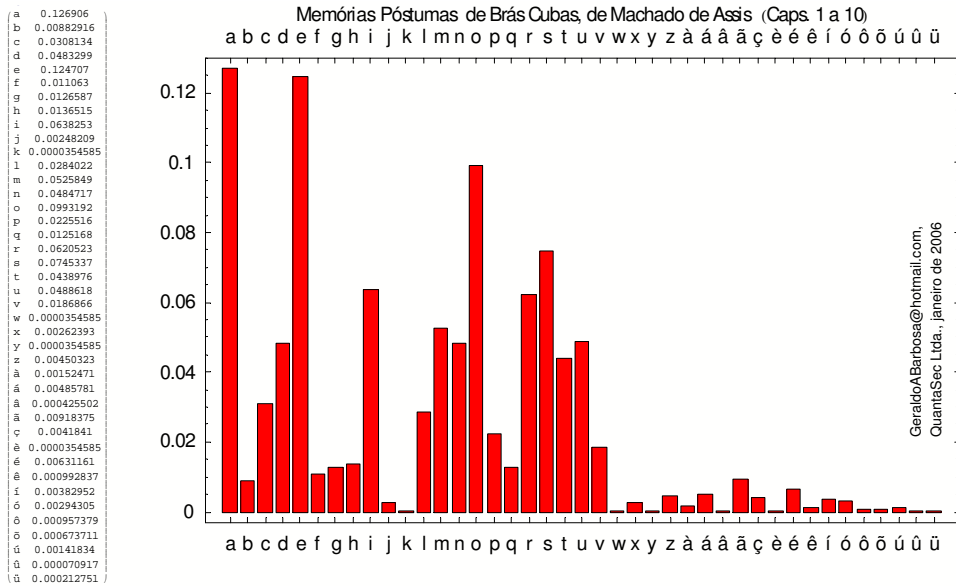


FIG. 4: Frequência das letras nos capítulos de 1 a 10 (28.202 letras) de Memórias Póstumas de Brás Cubas. A tabela à esquerda da figura apresenta numericamente a probabilidade de cada letra.

Observa-se que as letras “a”, “e” e “o” são as de maior ocorrência. A letra “a” ocorre com frequência maior do que a letra “e” do inglês. Será esta uma flutuação devida à pequena amostragem ou uma indicação de diferenças existentes entre estas línguas?. A resposta encontra-se na próxima figura. A figura 5 mostra as ocorrências de letras em todos os 160 capítulos. A maior ocorrência da letra “a” persiste assim como as frequências são distintas daquela apresentada na Fig. 1. Será que esta amostragem pode ser vista como característica da língua portuguesa? A seção a seguir responderá a esta pergunta. De fato, já foi possível perceber que a formação de palavras numa língua é de formação complexa e com baixíssima probabilidade de ocorrência mesmo nas palavras isoladas. Um texto completo numa língua é, portanto, muito característico daquela língua. Tal raciocínio já antevê a resposta qualitativa. Entretanto, pode ser dada uma resposta quantitativa? Passemos a próxima seção.

Entropias

A figura 6 mostra as entropias dos capítulos de 1 a 10, de 1 a 20, e assim sucessivamente até o último capítulo, de número 160. A esquerda estão os valores das entropias encontradas e, à direita, estes valores divididos pelo número de letras ou símbolos do alfabeto utilizado. Observe-se que o símbolo “ñ”, não pertencente à língua portuguesa foi incluído devido a

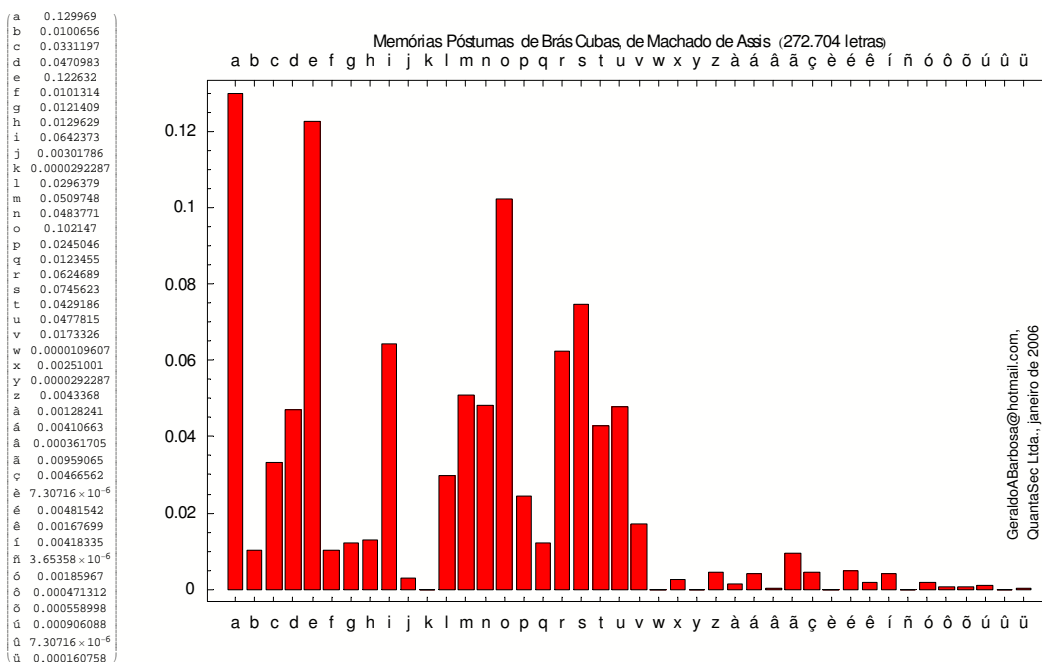


FIG. 5: Frequência das letras nos capítulos de 1 a 160 de Memórias Póstumas de Brás Cubas. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

citação de palavra espanhola no texto. Entretanto sua contribuição à entropia (valores à esquerda da figura) é desprezível. Os 35 símbolos considerados no texto são a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z, à, á, â, ã, ç, è, é, ê, í, ñ, ó, ô, õ, ú, û, ü. Desprezar-se o símbolo “ñ” traria uma correção por fator multiplicativo de $34/35 \simeq 0,97$ nos valores à direita. Deve-se notar que uma amostragem reduzida produz valores da entropia que variam consideravelmente. Ao se aumentar a amostragem considerada, com o acréscimo de capítulos, observa-se que estas flutuações são reduzidas ou suavizadas. A tabela I apresenta os valores numéricos das entropias obtidas. Espera-se que um limite assintótico seja atingido para números *finitos* amostrados para a entropia e a entropia por símbolo. Pode-se indagar se a amostragem utilizada já apresenta a estatística representativa da língua portuguesa. Entende-se também que pequenas flutuações poderão sempre ser notadas numa obra devido à particularidades daquele trabalho. Assim, num trabalho de caráter romântico, a palavra “amor” eventualmente poderá ocorrer com maior frequência do que a palavra “lobo”, com mesmo número de letras. Estas flutuações produzirão desvios em torno de uma média representativa da língua.

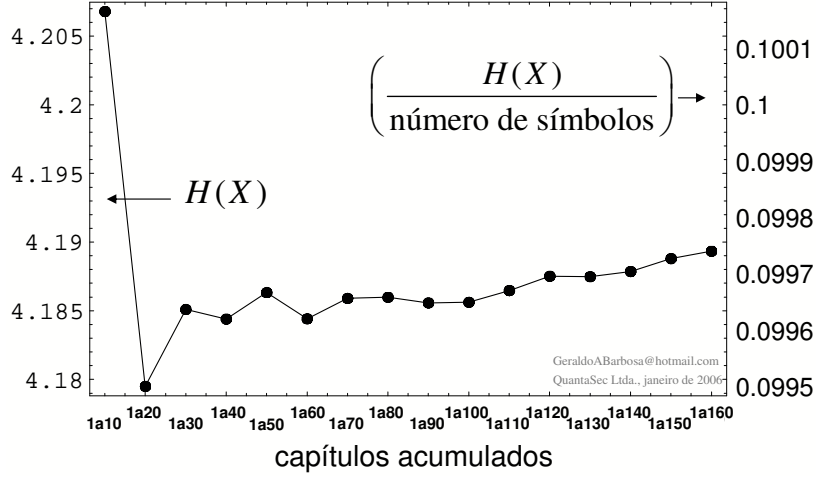


FIG. 6: Entropia de capítulos acumulados de 10 em 10. À esquerda estão os valores das entropias. À direita, a entropia por símbolo.

Redundância

Num alfabeto com L símbolos, a entropia é maximizada quando todos os símbolos têm a mesma probabilidade $P = 1/L$. Neste caso, a entropia seria (Ver Eq. 1)

$$H_{max} = \underbrace{L}_{\text{Número de símbolos}} \times \overbrace{\left(\frac{1}{L}\right) \log_2 \left(\frac{1}{L}\right)}^{\text{Entropia por símbolo}} = \log_2 L . \quad (2)$$

Este valor máximo sugere a medida da redundância r da língua em primeira ordem, isto é, a diferença relativa entre $H(X)$ para a ocorrência de letras e H_{max} :

$$r_1 = 1 - \frac{H(X)}{H_{max}} = 1 - \frac{H(X)}{\log_2 L} . \quad (3)$$

Se considerarmos que o valor r_1 representa a redundância da língua (o que não é uma boa aproximação, ver abaixo), seu valor sugere que a eficiência da língua possa ser melhorada neste percentual, ou seja, de alguma forma a língua poderia ser melhorada ou compactada para ser mais eficiente. A entropia encontrada para os textos analisados forneceu $H(X) \simeq 4.19$ com um alfabeto de 35 letras. Assim, $r_1 \simeq 0,223$ ou 22,3%. Note que uma “mensagem” arbitrariamente longa, consistindo da repetição de um mesmo símbolo teria entropia zero e, portanto, teria a máxima redundância possível ($r = 1$). Isto significa que a mensagem, a despeito de seu tamanho, não leva qualquer informação.

Como mencionado, a entropia pode ser calculada em aproximações mais representativas: assim como a Eq. (1) foi aplicada para a probabilidade de ocorrência de letras, ela

Caps.	Entropia	Entropia/símbolo
1 a 10	4.20679	0.100162
1 a 20	4.17949	0.0995118
1 a 30	4.18509	0.099645
1 a 40	4.18439	0.0996284
1 a 50	4.18632	0.0996743
1 a 60	4.18441	0.0996289
1 a 70	4.18590	0.0996643
1 a 80	4.18598	0.0996661
1 a 90	4.18557	0.0996564
1 a 100	4.18562	0.0996575
1 a 110	4.18647	0.0996779
1 a 120	4.18751	0.0997027
1 a 130	4.18748	0.0997019
1 a 140	4.18786	0.0997109
1 a 150	4.18880	0.0997333
1 a 160	4.18933	0.0997461

TABLE I: Entropia e entropia por símbolo.

poderia também ser estendida para os bigramas, trigramas e assim por diante. Como as maiores palavras numa língua tem um número de termos relativamente pequeno (“anti-constitucionalíssimamente?”), estes cálculos precisariam ser estendidos somente até estes valores.

Sendo N_{max} o comprimento do maior n -grama ocorrendo na língua, em princípio poderíamos então calcular entropias em várias aproximações:

$$\begin{aligned}
 H_1(X) &= \sum_x P(x) \log_2 \frac{1}{P(x)}; & H_2(X) &= \frac{1}{2} \sum_{x,y} P(x,y) \log_2 \frac{1}{P(x,y)}; \dots \\
 H_{N_{max}}(X) &= \frac{1}{N_{max}} \sum_{\underbrace{x,y,\dots}_{N_{max}}} P(\underbrace{x,y,\dots}_{N_{max}}) \log_2 \frac{1}{P(\underbrace{x,y,\dots}_{N_{max}})}. & (4)
 \end{aligned}$$

As entropias em ordem mais alta fornecem valores mais representativos para a língua do

que a entropia calculada simplesmente em primeira ordem. De fato, as entropias de ordem menor podem ser calculadas a partir das de ordem maior. Por exemplo, a probabilidade para acerto de um dado símbolo β numa língua, conhecendo-se um símbolo inicial α , será dada por

$$P(\beta|\alpha) = \frac{P(\alpha, \beta)}{P(\alpha)}, \quad (5)$$

e assim por diante. Entretanto, os cálculos passam a demandar muito tempo. Pode-se também optar pela estatística de ocorrência das próprias palavras, sejam como unidades independentes uma das outras ou até mesmo se considerando suas ocorrências vinculadas, isto é, a probabilidade de ocorrência de uma dada palavra condicionada ao aparecimento anterior de uma ou outras palavras.

Apesar de ser uma aproximação “pobre” a entropia calculada em primeira ordem é, entretanto, uma *característica* da língua e suficiente para os objetivos deste trabalho.

Deve-se observar que Shannon [4] se utilizou também de outros procedimentos bastante pragmáticos para se calcular a redundância de uma língua, como experimentos com textos “ocultos”, dos quais somente se conheciam uma certa percentagem das palavras. Por exemplo, apresenta-se o texto com um conjunto de espaços a serem preenchidos, igual ao número de letras e espaços do texto faltante. Passa-se a tentar adivinhar as letras ou espaços deste texto. O número das tentativas incorretas são anotadas para cada posição. Cada acerto é preenchido e age como uma informação para se diminuir o número de tentativas das letras ou espaços restantes. A partir destas contagem de tentativas constroem-se as probabilidades envolvidas para se acertar as letras ou espaços e, destas, constroem-se as entropias para cada letra (Um experimento deste tipo está sendo realizado atualmente através do “site” <http://www.numaboa.com.br/index.php>, desenvolvido por Viktoria Tkotz). Neste tipo de entropia encontram-se elementos de difícil estimativa numérica a-priori pois as probabilidades para acerto de cada símbolo sucessivo dependerá de forma correlacionada com todos elementos anteriores.

Neste trabalho, conforme utilizado na figura 6 e tabela I, as entropias por símbolo são determinadas automaticamente através das probabilidades e uso da equação 1 e da divisão do resultado pelo número de símbolos do alfabeto.

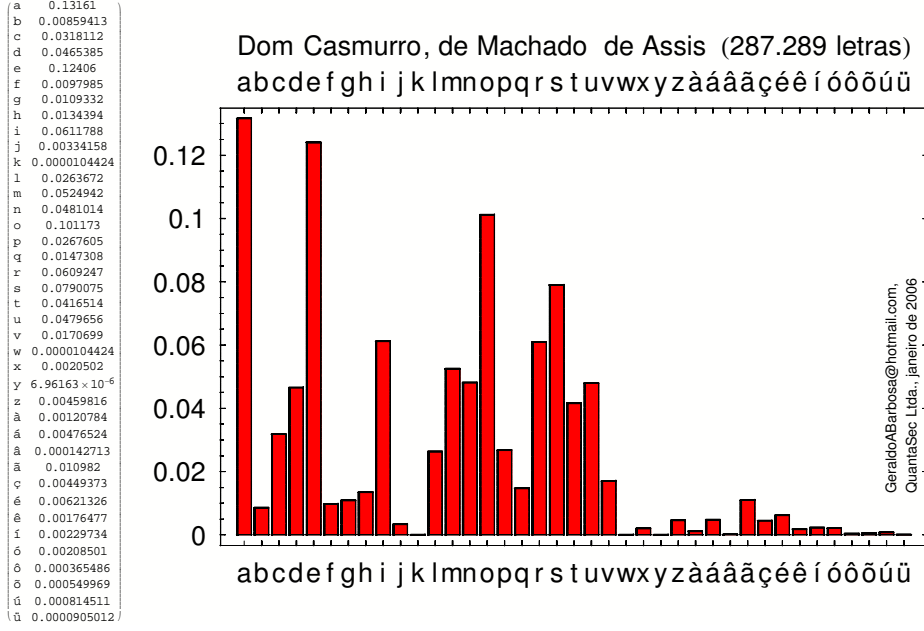


FIG. 7: Frequência das letras em Dom Casmurro. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

Outros textos de Machado de Assis

As figuras 7, 8, 9 mostram as frequências em “Dom Casmurro”, “Quincas Borba” e “Contos Fluminenses”.

Entropias

As entropias obtidas para os textos analisados de Machado de Assis, juntamente com a contagem das letras (entre parenteses) em cada texto são: $X_{Bras\ Cubas} = 4,208 (273.704)$; $X_{Dom\ Casmurro} = 4,180 (287.289)$; $X_{Quincas\ Borba} = 4,173 (342.190)$; $X_{Contos\ Fluminenses} = 4,167 (251.196)$. O total de letras considerado é acima de um milhão de letras: 1.154.379. A entropia média obtida para os textos, de forma ponderada, é

$$\overline{H}_W = \frac{273.704X_{BC} + 287.289X_{DC} + 342.190X_{QB} + 251.196X_{CF}}{1.154.379} \simeq 4,182 . \quad (6)$$

A entropia média, não ponderada, variância σ^2 e o desvio padrão σ são, respectivamente

$$\overline{H} = \frac{X_{BC} + X_{DC} + X_{QB} + X_{CF}}{4} = 4,187 \quad (7)$$

$$\sigma^2 = 0,000289 \quad (8)$$

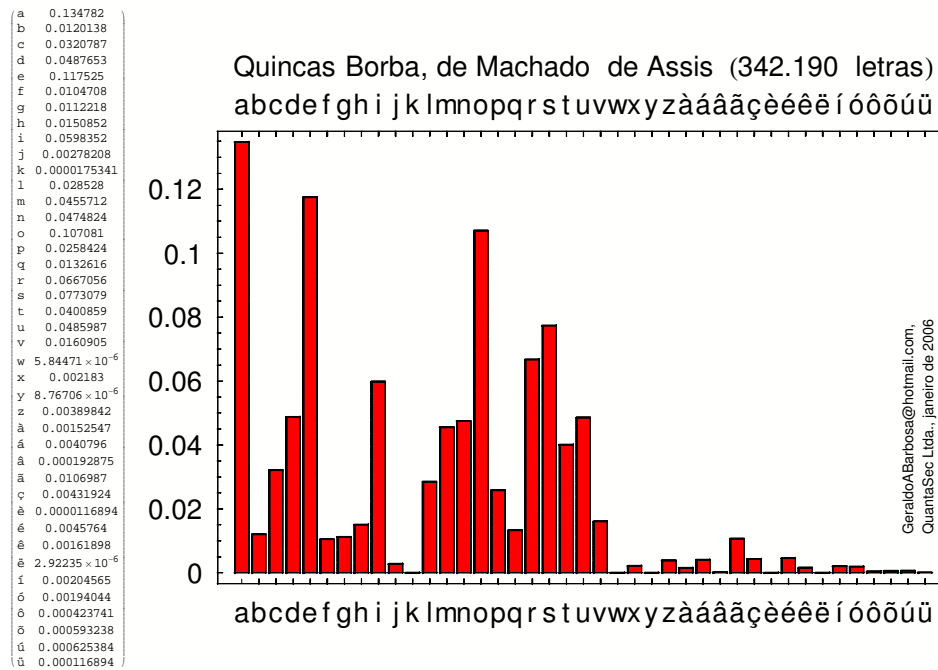


FIG. 8: Frequência das letras em Quincas Borba. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

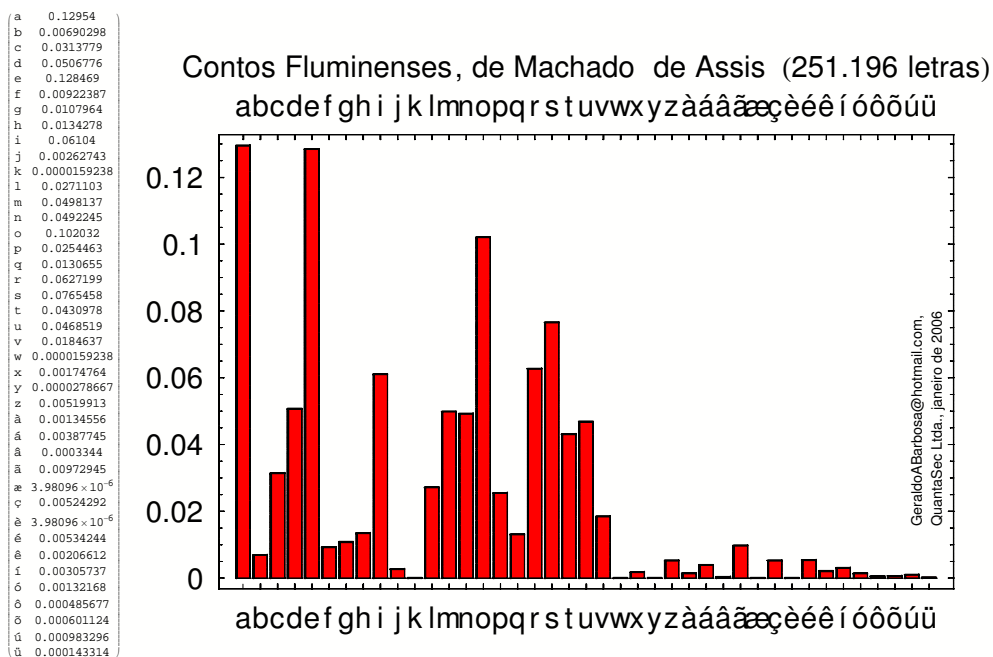


FIG. 9: Frequência das letras nos Contos Fluminenses. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

$$\sigma = 0,0170 . \quad (9)$$

Estes resultados indicam que, em amostragens bem maiores, muito possivelmente os valores a serem obtidos estarão em torno de $H \simeq 4.18 \pm 0,02$.

Bigrama de Mémoires Póstumas de Brás Cubas

O bigrama em barras do texto analisado mais detalhadamente, “Mémoires Póstumas de Brás Cubas” está mostrado na figura 10 e deve ser analisado como na figura 3. A figura 11 mostra a matriz de valores do bigrama de barras da figura 10. Deve-se observar que o alfabeto utilizado ou mostrado em cada análise pode variar, devido a utilização ou não de certas letras. De outra forma, pode-se também mostrar um alfabeto maior do que o utilizado para, eventualmente, se acomodar letras de algum alfabeto diferente (como o “ñ”). Por exemplo, a figura 10 foi construída com posições para o alfabeto de 45 letras: a1,b2,c3,d4,e5,f6,g7,h8,i9,j10,k11,l12,m13,n14,o15,p16,q17,r18,s19,t20,u21,v22,w23,x24,y25,z26,à27,á28,â29,ã30,ç31,è32,é33,ê34,ì35,í36,ñ37,ò38,ó39,ô40,õ41,ù42,ú43,û44,ü45. Neste caso, muitas letras jamais foram citadas e o bigrama mostrará o valor zero para as mesmas. A altura das barras mostra o número de ocorrências de cada bigrama. A divisão do número de ocorrências de cada bigrama, dividido pelo número total de letras com ocorrências fornece a probabilidade adjunta $P(x, y)$. A figura 11 mostra numericamente as ocorrências de cada bigrama.

CARTA DE CAMINHA

Esta seção analisa a “Carta de Caminha” nas versões contemporânea e original. Dentre os resultados buscados pergunta-se: 1) Qual é a diferença estatística na ocorrência das letras nas duas versões? 2) Qual é a entropia das duas versões e a entropia por símbolo?

Carta, versão contemporânea

A figura 12 mostra a estatística de ocorrência das letras na versão contemporânea da “Carta de Caminha” A figura 13 mostra a diferença de ocorrência de letras entre a versão contemporânea da “Carta de Caminha” e “Memórias Póstumas de Brás Cubas”.

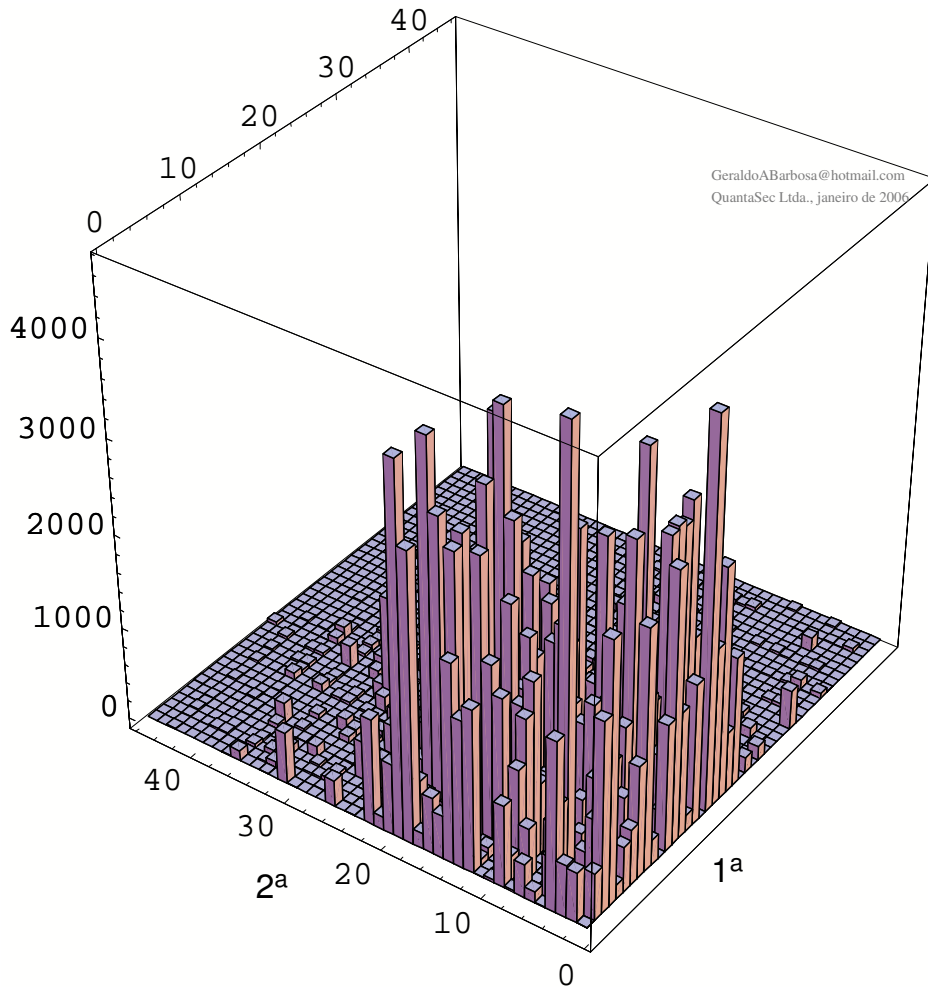


FIG. 10: Bigrama tridimensional do texto completo “Mémorias Póstumas de Brás Cubas”. Por simplicidade, deve-se utilizar os números associados ao alfabeto de 45 letras utilizado.

A entropia da versão contemporânea é $H_{Cont} = 4,117$ e sua entropia por símbolo é $H_{Cont}/\text{símbolo} = 0,114$. A redundância da Carta em versão contemporânea é $r_{Cont} = 0,204$.

Carta, original

A figura 14 mostra a estatística de ocorrência das letras na “Carta de Caminha” original. A figura 15 mostra a diferença de ocorrência percentual de letras entre a “Carta de Caminha” original e “Memórias Póstumas de Brás Cubas”.

A entropia da Carta original é $H_{Original} = 4,082$ e a entropia por símbolo é $H_{Original}/\text{símbolo} = 0,136$. A redundância da Carta original é $r_o = 0,168$.

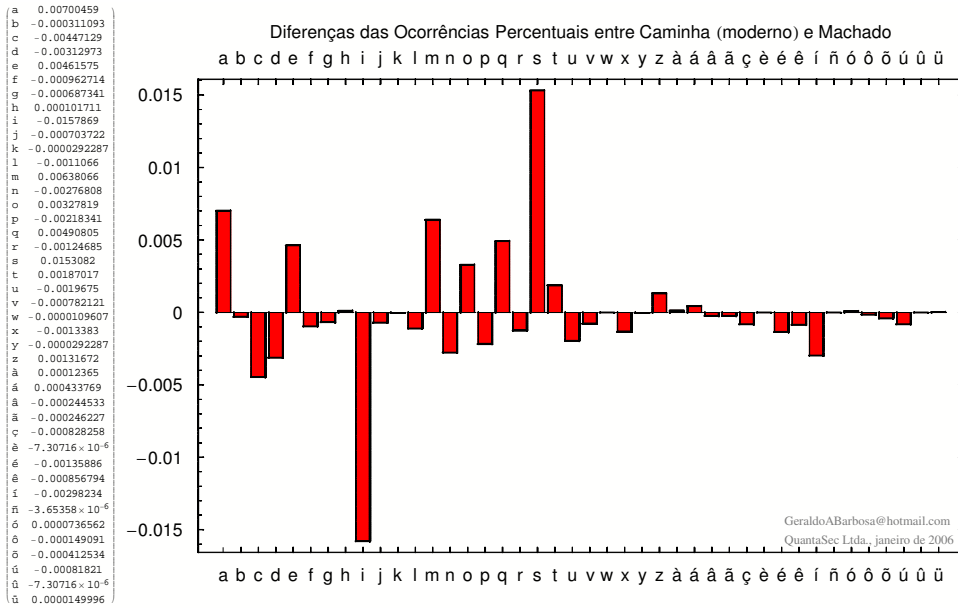


FIG. 13: Estatística da diferença de ocorrência de letras entre a versão contemporânea da “Carta de Caminha” e “Memórias Póstumas de Brás Cubas”. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

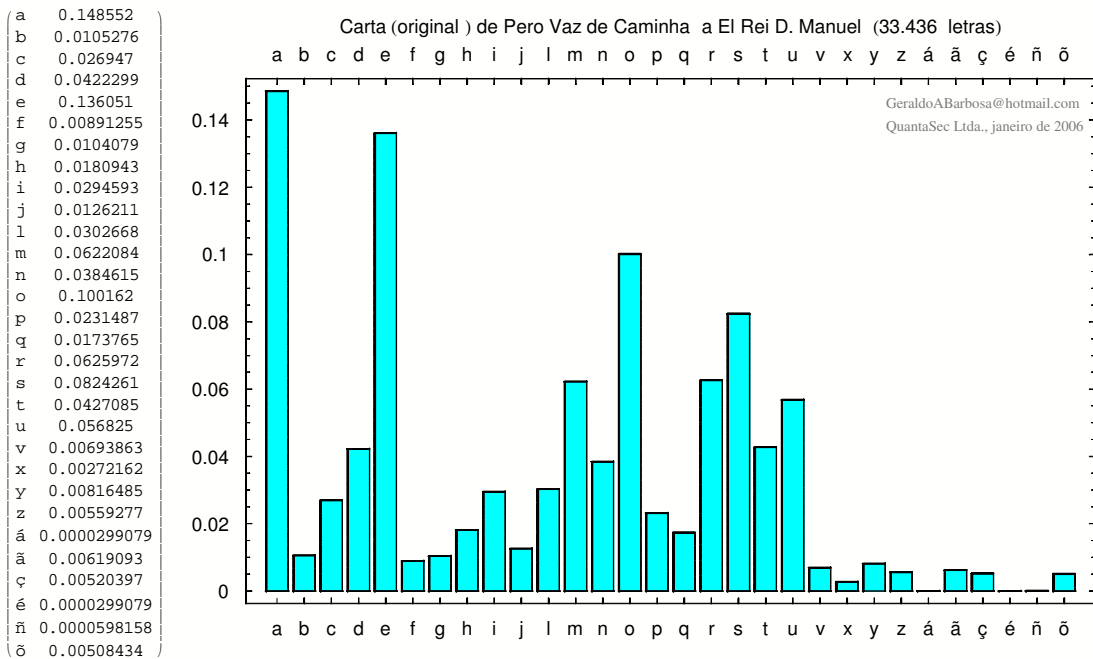


FIG. 14: Estatística de ocorrência das letras na “Carta de Caminha” original. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

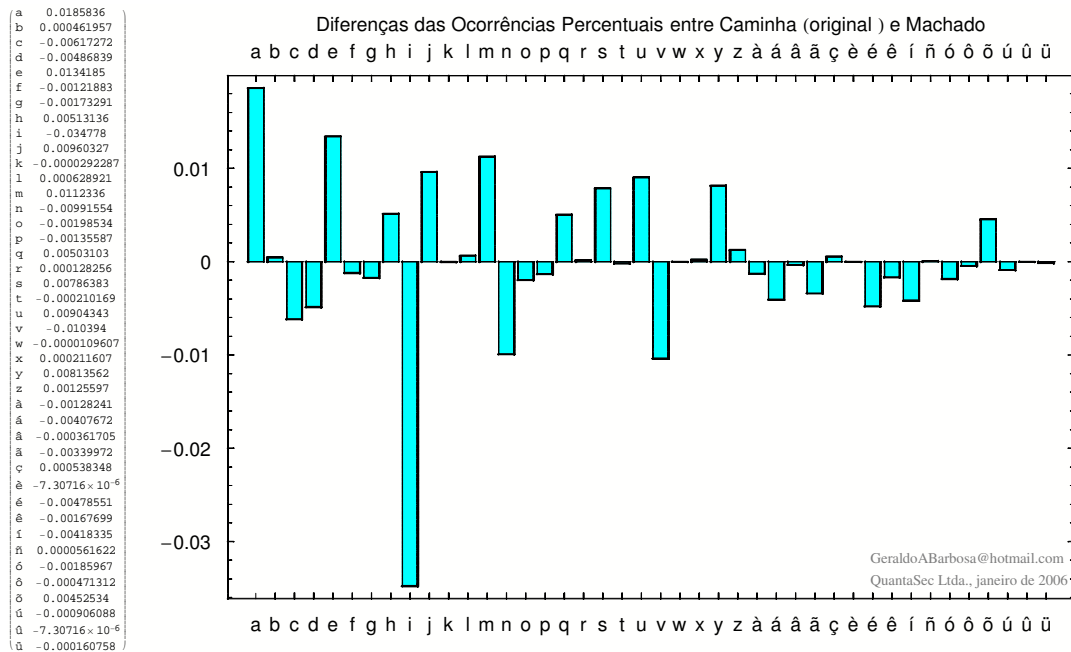


FIG. 15: Diferença de ocorrência percentual de letras entre a “Carta de Caminha” original e “Memórias Póstumas de Brás Cubas”. A tabela à esquerda da figura apresenta numericamente as letras encontradas e sua probabilidade de ocorrência.

Algumas comparações

É interessante notar que as versões contemporânea e original mostram as entropias por símbolo 0.114 e 0.136, respectivamente. A redundância da versão contemporânea é maior do que a do original, $r_{Cont} = 0,204 > r_o = 0,168$, mostrando que a mesma mensagem foi transmitida por Caminha de forma mais econômica do que a versão moderna. De fato, basta observar que a mensagem original (33.436 letras) utiliza cerca de 1.000 símbolos *menos* do que a contemporânea (34.138 letras). Será que este fato deve-se mais aos redatores da forma contemporânea, em uso redundante da língua portuguesa do que uma característica da própria língua?

A figura 16 mostra as diferenças percentuais relativas na ocorrência de letras entre a “Carta de Caminha” original e sua versão moderna. Os dados foram calculados como

$$\frac{\text{Ocorrência fracionária de letras no original} - \text{Ocorrência fracionária de letras na versão contemporânea}}{\text{Ocorrência fracionária de letras no original} + \text{Ocorrência fracionária de letras na versão contemporânea}} \times 100 . \quad (10)$$

Assim, os valores 100% indicam letras existentes no original mas não na versão contemporânea e -100% , o oposto.

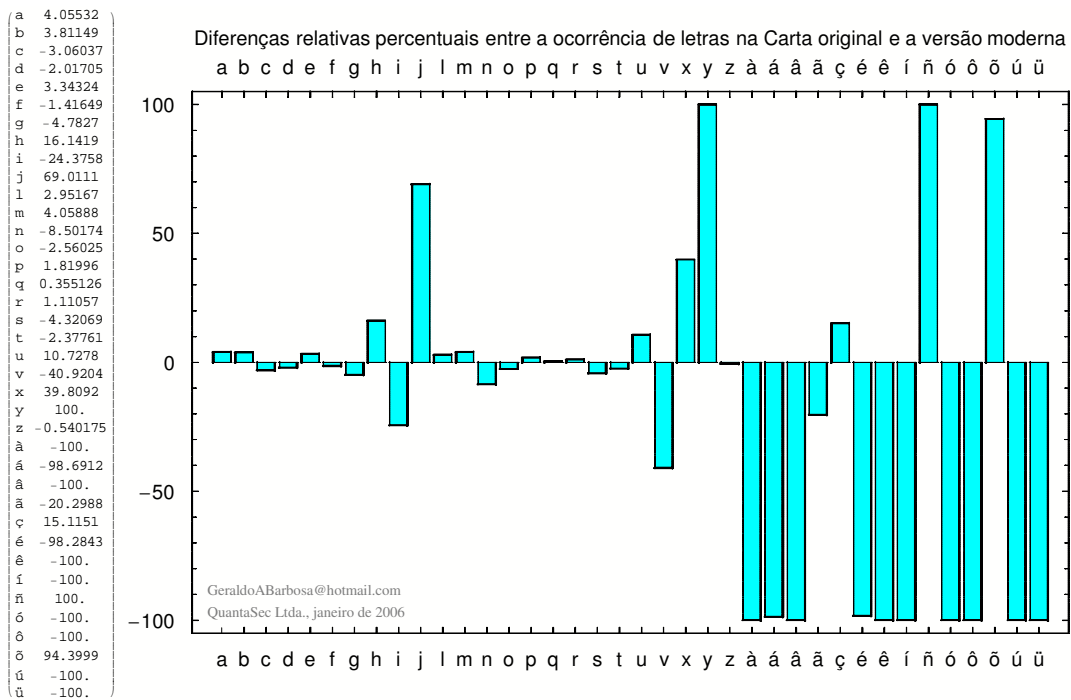


FIG. 16: Diferenças relativas percentuais na ocorrência de letras entre a “Carta de Caminha” original e sua versão contemporânea. A tabela à esquerda da figura apresenta numericamente as letras encontradas e as diferenças percentuais relativas para cada letra.

CONCLUSÕES

Neste trabalho examinou-se a estatística de ocorrência de letras na língua portuguesa utilizando-se textos de Machado de Assis acessíveis na Internet e contabilizando mais de um milhão de palavras. Para a obra “Memórias Póstumas de Brás Cubas” foi apresentado também seu bigrama. Um valor numérico representativo da *entropia* da língua portuguesa –pelo menos à época deste autor– foi calculado. A existência de apreciáveis flutuações estatísticas para amostragens em baixo número foi evidenciada assim como seu amortecimento com o aumento deste número. Foi indicado que pequenas flutuações da entropia estarão sempre presentes em obras de qualquer porte mas estas flutuações ocorrerão em torno de um valor que se pode dizer característico da língua portuguesa. Foram comparados também os resultados de Machado de Assis com os obtidos da Carta de Caminha, nas versões contemporânea e original. Mostrou-se também diferenças estatísticas entre o original e a referida versão contemporânea da Carta.

Seria muito interessante um estudo de maior número de textos visando corroborar ou

não os resultados aqui obtidos e se estabelecer de forma definitiva alguns parâmetros significativos da língua portuguesa contemporânea e mesmo a evolução da língua. A obtenção de entropia em melhores aproximações poderia ser também um dos objetivos nesta análise maior, para fornecer um estudo bem mais completo.

O advento da computação possibilitando o tratamento de grandes volumes de dados indica que mudanças introduzidas na língua por propostas conscientes deveriam ser acompanhadas por estudos estatísticos que apoiassem estas propostas.

Mais do que os resultados aqui apresentados pretende-se estimular os jovens interessados na língua portuguesa a se utilizarem de ferramentas estatísticas que se tornam indissociáveis do mundo moderno e que permitem análises impraticáveis de serem feitas manualmente.

*O autor é Ph.D. pela University of Southern California (1974) e foi Professor Titular da Universidade Federal de Minas Gerais, Brasil. Atualmente, é Professor no Center for Photonic Communication and Computing, ECE Department, Northwestern University, Evanston, IL 60208-3118, US

* Electronic address: `Email:GeraldoABarbosa@hotmail.com,g-barbosa@northwestern.edu`

- [1] D. Kahn, *The Code-Breakers, The Story of Secret Writing* (Scribner, New York 1996). S. Singh, *The Code Book, The Science of Secrecy from Ancient Egypt to Quantum Cryptography* (Anchor Books, New York 1999).
- [2] J. Bamford, *Body of Secrets, Anatomy of the Ultra-Secret National Security Agency* (Anchor Books, New York 2002).
- [3] D. J. C. MacKay, *“Information Theory, Inference, and Learning Algorithms”* (Cambridge 2003).
- [4] C. E. Shannon, *The Bell System Technical Journal*, Vol. **27**, pp. 379423, 623656, July, October, 1948.