

蒙古文拼写形式多样化问题研究¹

白双成¹²³ S·苏雅拉图²³张劲松¹

¹北京语言大学 信息科学学院, 北京 100190;

²内蒙古社会科学院 蒙古语信息处理研究所, 呼和浩特 010020;

³内蒙古蒙科立软件有限责任公司, 呼和浩特 010019)

摘要: 蒙古文文本中存在一个有别于多数其他文字的特别现象——看到的单词字形正确但其内码序列不正确, 或者说单词“变形显现字形”序列正确但“名义字符”序列不正确的现象, 我们称其为蒙古文的拼写形式多样化现象。本文先定义该现象及相关概念, 再通过简单图示、例词拼写形式穷举、新闻语料统计分析和基于整篇文章标注统计等多种方式、多角度论证这一现象的事实性和严重性, 分析导致这一现象的深层原因并指出拼写形式多样化对蒙古文信息处理和应用方面的严重影响, 最后提出通过推广普及录入规范和标准提高用户意识、使用智能输入法避免误录、使用校对纠错工具后纠正、基于生语料的统计学习方法为补充等多途径解决方法。本文对蒙古文标准编码的推广普及具有较好的参考价值。

关键词: 蒙古文文本; 拼写形式多样化; 读音错误; 字形错误; 智能输入

Study of Mongolian Spelling Diversity Phenomena

BAI Shuangcheng¹²³ S·Soyolt²³ ZHANG Jinsong¹

¹Beijing Language and Culture University, Beijing 100190, China;

²MIT Center, Inner Mongolia Academy of Social Science, Hohhot, 010020, China;

³Inner Mongolia Menksoft Software Co., Ltd., Hohhot, 010019, China)

Abstract: Unlike most other script, there is a special phenomenon in Mongolian text -- word shape is correct but the code sequence is maybe wrong, or the words, presentation form is correct but nominal character sequence is wrong, we call it Mongolian Spelling Diversity Phenomenon. This article firstly defined the phenomenon and explained related concepts, than demonstrated this phenomenon by exhausting spelling style of example word, raw corpus statistics and whole text annotation statistics etc. We analyzes the deep causes of this phenomenon and points out that MSDP's serious impact on Mongolian information processing and application. Finally, we promote serious of methods to solve MSDP, it includes popularization of standardized input by raise awareness of the user, using intelligent IME to avoid input mistake, using the spelling error correction tools, exploration of raw corpus based statistical learning methods etc.

Keywords: Mongolian text; Spelling Diversity Phenomena; intelligent IME; Spelling Error

引言

蒙古文编码国家标准²将蒙古文字符分为“名义字符”(Nominal Character)集和“变形显现字形”(Presentation Form/Character)集, 并规定前者用于信息存储、传输和计算, 有明确的码位, 后者(即蒙古文各类纷繁复杂的字母变形形式)仅用于信息的输出(显示和打印), 不需要码位^[1]。把一种文字的字符分为“名义字符”和“变形显现字形”进行编码,

¹**基金项目:** 国家电子发展基金 2010 年度、2011 年度蒙古文专项; 国家自然科学基金(61163020); 内蒙古自治区自然科基金项目 2011MS0918 资助项目

作者简介: 白双成(1974—), 男, 博士研究生, 研究员, 主要研究方向为语言工程; 苏雅拉图(1956—), 男, 研究员, 主要研究方向为蒙古文信息处理; 张劲松(1968—), 男, 博士, 教授, 博士生导师, 主要研究方向为基于语音语言处理技术的计算机辅助汉语教学技术的研究工作。

²蒙古文编码标准有国际标准和国家标准, 目前两者即相关又有所区别。国际标准只有大体框架, 国家标准是在其基础上的细化和升级, 是标准、用户协定和转换规则三者的统一体。符合国家标准必然符合国际标准, 而符合国际标准的未必能符合国际标准, 所以本文所述“标准编码”特指国家标准。

并把信息传输与信息输出截然分开来处理,是一个完全不同于英法德等西方文字和汉日韩等东方文字通用方式的较特殊的编码结构与编码方式。原本在其他文字中特别简单的字符 C_i 到字形 T_i 的“一对一”“映射关系”,变成了字符序列 $C_1C_2\cdots C_n$ 到字形序列 $T_1T_2\cdots T_m$ 的“多对多”的复杂“转换过程”。由于蒙古文的这一转换过程远复杂于阿拉伯文,仅仅依靠词首、词中、词尾等词内位置进行调形的阿拉伯文成功经验已无法满足蒙古文复杂变形需求。幸好目前已有部分操作系统、浏览器及通用基础软件环境具备了标准编码实现条件,技术层面上基本攻克了这一复杂转换过程的技术难题。只待“用户协定”和转换规则^[2]进一步查漏补缺,各机构共同执行即有希望实现对于任意一个名义字符序列 $C_1C_2\cdots C_n$ 转换为唯一一个变形显现字形序列 $T_1T_2\cdots T_m$,从而真正实现编码统一。为解决同形词(Homograph),完全有可能存在另一个名义字符序列 $C_1'C_2'\cdots C_k'$ (k 与 n 不一定相等)同样被转换为 $T_1T_2\cdots T_m$ 。Unicode 标准 Core Specification 中 13.4 节给出了如图 1 的一个案例。

表 1 Unicode 标准提供的蒙古文词形趋同案例

字符序列 C			词形	词意	字符序列 C'			词形	词意
ᠠ	1824	u	ᠠᠯᠠᠭᠤ	长度		1823	o	ᠠᠯᠠᠭᠤ	宫殿
	1837	r				1837	r		
	1832	t				1833	d		
	1824	u				1824	u		

但是, C' 是否可能是一个完全不合法的字符序列或根本不是用户想要的字符序列呢? 这种我们所不希望的事情发生概率如何呢? 非常不幸的是, 由于蒙古文的“一字一音”、“一字多形”和“一形多字”等自然属性, 这种可被映射为相同 T 的 C 的组合非常庞大且非常常见(见下文分析)。蒙古文文本中存在的, 这种单词输出(屏幕显示或打印)的“变形显现字形”所表现的字形正确, 但其存储的“名义字符”序列(文本内码序列)不正确现象, 我们称为**蒙古文拼写形式多样化现象**(Mongolian Spelling Diversity Phenomena), 简称为**拼写多样化**。也有文章称其为同形异码词^[3], 是对同一事实的不同层次命名。

正由于蒙古文的拼写多样化现象, 蒙古文拼写错误(spelling error)^[4]有别于其他文种, 可细分为“读音错误”和“词形错误”。依据字符序列 C 转换结果 T 的唯一性特性可知, 当一个单词的词形 T 错误时对应的 C 也一定错误, 所以词形错误时读音一定错误(详见 2.5)。相反, 词形正确却无法保证其读音正确。所以拼写多样化也可以叫“字形正确, 读音错误³”现象。从这个意义上讲, 蒙古文的读音错误才是其他文字中所述拼写错误, 但从其他文字经验和直觉而言, 很容易理解为字形错误是拼写错误。所以我们也可以说蒙古文的拼写错误有“字形正确但读音错误”和“字形读音都错误”两个层次。本文中“词”是指单词名义字符序列, 而“词形”是指其变形显现字形序列, 所以同形词即指同形异码词。本文生语料统计计算中没有进行字形纠错和归并。

另需要说明的是, 蒙科立编码^[5]作为一种“全字符编码”方式^[4], 本身就是基于标准编码框架的一种变形显现字形方案。所以, 不管是按名义字符形式保存的标准编码还是按变形显现字形形式保存的蒙科立编码, 只要是基于“音”的编码形式就必然存在拼写形式多样化现象, 本文分析结果通用于所有此类编码方式的文本。文章^[3]提到蒙科立编码文本语料库中存在同形异码词, 但此文集中于用同形异码字符替换和符合字符拆分、组合方式归并同形词上, 与本文目标具有较大差异。蒙古文自动校对^{[6][7]}等也提及蒙古文同形异码, 但都没有对此进行特别深入分析。

1 拼写多样化情况

1.1. 简单的拼写多样化案例

我们先看一个简单且容易理解的拼写形式多样化现象。因 \ddot{o}/\ddot{u} 、 x/g 、 d/t 、 a/e 四对字母在相同词内位置和相同阴阳性()条件下经常表现为同形, 人们很容易就“发现”

³ 同形词也属于拼写形式多样化, 但读音不错。

数据的原因主要有：

[1]. **可获得性**：选择网络资源的首要原因是它有便利的可获取性，便于其他研究人员也可以获取对照。

[2]. **可靠性**：正规新闻媒体机构主办，稿件经过编辑、审核等多道编审流程发布。具有内容相对可靠、术语相对规范统一、干扰因素相对低等优势。

[3]. **时效性**：作为新闻类网站，具有较强时效性，可基本反映新词术语和蒙古文使用现状。当然，时政类新闻为主的新闻内容词汇量必然比不上文学作品，文风也相对拘谨。

[4]. **代表性**：虽然是正规新闻稿件，但依然存在较为严重的读音错误，也不乏字形拼写错误，具有普遍代表性。

[5]. **可验证性**：因工作便利，可对此三个网站爬取内容进行正确性验证，确保网页爬取、网页模板分析、格式转换、行序回复等工作正确。

[6]. **结构性**：可额外获得结构化数据（Structured Data），便于进行按文档种类各自分类训练，便于进行关键字抽取（Keyword Extraction）、摘要生成（Summary Generation）等有监督学习（Supervised Learning）的后续研究工作。搜索引擎中结构化搜索（Structural Search）就是基于 MNN 的结构化数据（实际用 TREC 标记标示）进行了充分训练和验证。

[7]. **可延续性**：这三个新闻网站每日稳定更新。通过前期试验，语料搜集工具成熟后，可以定期更新扩充语料。

将 MGLNews 语料所有单词按字形进行归类后获得了所有词形统计数据，表 4 节选展示了部分典型数据。

表 4 拼写形式统计节选数据

拼写种类	所属词形频度	所属词形数	平均频度	所属词例 (top 5)
273	90485	1	90485	□□□□□□
179	5954	1	5954	
118	11472	2	5736	
.....				
50	122088	29	4209	
49	181698	31	5861	
.....				
1	310317	141006	2	

表 4 第一行表示，词形 □□□□□□ 在语料中共出现 90485 次，有 273 种⁵不同拼写方式，是拼写形式最多的词形。第三行表示 □□□□□□ 和 □□□□□□ 这两个词形各有 118 种拼写方式，共出现 11472 次，每个词形平均出现频度为 5736 次。最后一行表示 141006 个词形只出现了一种拼写形式。如前所述，语料中不仅存在读音错误，也有字形错误，表 3 统计过程中我们并没有进行错误字形纠错。例如，表 5 展示了 MGLNews 中单词 ᠶᠢᠯᠠᠭᠠᠨ 的部分词形相似词，其中只有一个是形近字，其余都是拼写错误。

表 5 词 ᠶᠢᠯᠠᠭᠠᠨ 的部分字形相似词

单词	类型	原因
ᠶᠢᠯᠠᠭᠠᠨ	原词	名词“要求”
ᠶᠢᠯᠠᠭᠠᠨ	形近	动词“要求”
ᠶᠢᠯᠠᠭᠠᠨ	拼错	替换 辅音“ ” 替换为辅音“ ”
ᠶᠢᠯᠠᠭᠠᠨ	拼错	替换 辅音“ ” 替换为“ ”
ᠶᠢᠯᠠᠭᠠᠨ	拼错	替换 辅音“ ” 替换为元音“ ”
ᠶᠢᠯᠠᠭᠠᠨ	拼错	替换 辅音“ ” 替换为辅音“ ”

⁵因篇幅所限，在此略去拼写方式数据。

a -yie 和 jin……jie jia等十几种错误拼写方法。

2.3. 控制字符误用

因控制字符是不可见字符，目前操作系统和编辑器又缺少控制字符查重（Duplication Checking）或过滤（Filtering）功能，乱用、误用情况在所难免。尤其是生僻且写法特殊的外来词，录入者可能反复交替试用几个控制字符后最终获得所需字形，但有可能录入了多余控制符而浑然不知。即使是常用词也有可能中间插入多余控制字符而表面上看不出来。例如 ġ 后放置 FVS1（U+180B）后变成

写词缀和前导词。由于部分独立单词字形与分写词缀同形，自动纠错难度较大，从而，统计准确率很难得到保障。工作量统计中分写词缀算作一个单词，也许有利于编辑人员稿费统计而故意为之。

检索、排序和统计等最基本的应用需求都难于得到满足，我们该怎么办呢？

4 拼写形式多样化问题的解决方案

既然有这么多问题，我们该如何解决呢？大致有两种解决思路。一种是要完全正确录入或同时对被搜索内容和搜索关键字进行拼写纠错，区分同形异音字母和词，达到精准搜索、精准排序、精准统计。另一种是通过额外算法解决同形异音字母的字形模糊匹配，达到模糊搜索和模糊统计，但这是否违反了标准编码制定初衷，纵容用户拼写错误了呢？对此，我们建议：

4.1. 推广普及录入规范和标准，提高用户意识

由蒙古文的自然属性、编码特性及拼写多样现象可知，让用户接受蒙古文标准编码的录入规范，形成良好习惯并不是简单的书写媒介更换问题，可以说是一次重大的变革。把握好了，我们可以借此机会将标准音与信息化同时推进。把握不好，将严重阻碍蒙古文信息化进程。为此我们认为，必须同时注重现有使用者培训和未来使用者培养，且后者更重要。

未来使用者培养方面，将蒙古文标准编码的录入规范纳入中小学课程，结合标准音推广工作，从入学儿童开始培养起良好的习惯，待他们毕业走上工作岗位时，将按标准录入视为很正常的必然事件。

现有使用者培训方面，我们先重点培训编辑、记者、相关蒙古文工作人员，再逐步扩大到普通用户步骤。单从报纸、期刊和图书出版角度考虑，要想让长期习惯于纸质出版做法的相关人员按规范录入具有一定难度。兼顾字形和读音会额外耗费精力，这明显有悖于他们追求绩效考核的初衷。那么如何让他们意识到这样做的好处并愿意付出这份努力也许是个不简单的系统工程问题。将出版资源升级为“语料库”，作为商品授权给研究机构及商业公司，获取经济效益等都不一定能凑效。如果不能形成长效机制，很难长久维持。所以从他们自身使用便利角度考虑，让他们意识到这样做了可以便利利用以往资料等可能是更合理的方式。

4.2. 使用智能输入法避免误录

在标准编码 OpenType/AAT 字库实现中，一般都会附带一个键盘映射(Keyboard Mapping)输入法。因复杂文本引擎和字体规则担负了名义字符到变形显现字形的映射转换，输入法本身一般只需从键盘字母映射为蒙古文字母，一般不用做额外处理。因蒙古文是个同形异音字符较多的文字，这种键盘映射输入法没有避免同形字符输入错误避免措施，从输出的字形又不易察觉错误，所以，即使用户了解标准编码框架、懂得输入规范，也意识到规范录入的必要性，但指望所有用户都能按标准录入是不现实的。

为标准编码的推广普及，不让终端用户陷入迷茫的一个有效途径就是推广智能化程度较高的输入法做预防性处理。鼓励使用完全符合规范和标准的智能输入法，从录入源头避免错误^[9]。智能化输入法确保录入字形和读音正确的同时给用户简单易用的用户体验，让用户不再感觉遵循规范是个负担。此处所述输入法不仅局限于全键盘录入，也包括智能终端的虚拟全键盘、数字键盘及 OCR 识别录入、语音识别录入等所有输入方式。这些输入方式上必须加以监督机制，尽量避免用户录入错误。不管输入法做到什么程度，总是无法避免 OOV，而这部分的正确录入只能依靠用户自律或通过网络协同等方式作为弥补。

4.3. 使用校对纠错工具后纠正

虽然通过前两项可以解决今后的问题，但对于历史数据或字形扫描、手写录入等场合，我们需要依赖自动校对和自动纠错。目前“词典+规则”是实现蒙古文文本校对常用方法^[6]。不管使用不确定有限状态自动机(NFSA)数据结构获取较高计算效率、使用词干/词缀和生

成规则来节省存储空间或是使用最一般的字符串匹配的库结构，其本质无非都是依赖词库，词库中有的词认为是正确词，词典中没有词（OOV）就认为是错别词。再进一步用搭配库或规则对部分同形异音词（例如 $\Pi\mathcal{P}'1 \sim \omega$ 的 $\Pi\mathcal{P}$ 录入为 $H' \ominus \text{H} \ominus$ 的 $H' \ominus$ 进行甄别。从公开资料来看，未登录词、同形多音词处理还不够成熟，句法和语义层面错误基本未能触及，甚至词法层面的形态变化分析^[10]还有待提高，校对和纠错效果基本取决于词典词汇量。字形拼写错误的纠错也不尽如人意，所以有待进一步完善和改进校对和纠错。

4.4. 探索基于生语料的统计学习方法

我们有了确保单词读音正确的熟语料，可以顺利开展一些统计建模的科研工作^[11]，但目前我们所能获得量还难以支撑实际应用需求，更无法满足需要大数据支撑的个别模型。虽说可以采取各种手段缓解数据稀疏（Data Sparse），但归根结底还得需要足够量数据支撑统计建模^[12]。很显然，深度机器学习（Deep Machine Learning）使用的词向量表示（Word Vector Presentation）来说，语料量越大，低维空间（Low Dimensional Space）上的词向量越趋于精准^[13]。更何况熟语料没有读音错误只是一种假设，而我们日常产生的原始数据又有如此严重的拼写多样化，我们所能采取的防范措施又不能解决所有问题，所以我们不能单纯等待和依赖加工足够量的熟语料后再开展相关研究工作。另一方面，新词术语研究、语言动态监测、舆情分析等工作总不能还要依赖加工的熟语料。直接利用生语料的研究工作是熟语料建设的必要补充和回旋途径，相辅相成，互为补充，也是解决拼写多样化的一个重要途径。

5 总结

综上所述，蒙古文文本中大量存在拼写多样化现象，严重影响着蒙古文文本的日常应用及科研工作。各种解决方式都无法独自满足需求，需加以综合利用。各项工作的开展势必依赖语料库建设、知识库建设及相应大数据、机器学习等方面的突破，所以我们的研究工作还任重道远。

参考文献

- [1]. The Unicode Consortium[EB]. <http://www.Unicode.org>.
- [2]. 确精扎布. 确精扎布蒙古文信息处理专辑[M]. 内蒙古教育出版社, 2014.
- [3]. 敖敏, 熊子瑜, 呼和. 基于蒙科立输入法的蒙古语同形异码词研究[C]. 第十一届全国人机语音通讯学术会议, 2011, 10.
- [4]. Aminul Islam and Diana Inkpen. Real-Word Spelling Correction using GoogleWeb 1T 3-grams . Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, pages 1241 - 1249, Singapore, 6-7 August 2009. c2009 ACL and AFNLP
- [5]. 白双成, 张劲松, 呼斯勒. 蒙古文输入法输入码方案研究[J]. 中文信息学报 2013(06):169-174.
- [6]. 斯·劳格劳. 基于不确定有限自动机的蒙古文校对算法[J]. 中文信息学报, 2009, (06) .
- [7]. 苏传捷, 侯宏旭, 杨萍等. 基于统计翻译框架的蒙古文自动拼写校对方法[J]. 中文信息学报, 2013, (06)
- [8]. 国家质量监督检验检疫总局, 国家标准化管理委员会. GB 25914-2010. 信息技术传统蒙古文名义字符、变形显现字符和控制字符使用规则[S]. 北京. 中国标准出版社, 2011. 11.
- [9]. S·苏雅拉图. 蒙古文整词计算机生成理论研究[J]. 中文信息学报, 2001(04):59-65.
- [10]. Deniz Yuret, Ergun Bici. Modeling Morphologically Rich Languages Using Split Words and Unstructured Dependencies[j]. ACL-IJCNLP 2009 Conference.
- [11]. 赵伟, 侯宏旭, 从伟, 宋美娜. 基于条件随机场的蒙古语词切分研究. 中文信息学报. 2010, 24(5).
- [12]. Daniel Jurafsky, James Martin. Speech and Language Processing[M] (英文版第2版 人民邮电出版社. 2010).
- [13]. Jacob Devlin; Rabih Zbib. Fast and Robust Neural Network Joint Models for Statistical Machine Translation. ACL2014. 1370—1380.