

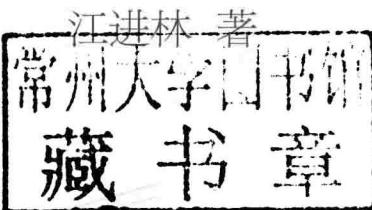
# 中国学生英译汉机器评分 模型的研究和构建

江进林 著

高等教育出版社

ZHONGGUO XUESHENG YINGYIHAN JIQI PINGFEN  
MOXING DE YANJIU HE GOUJIAN

# 中国学生英译汉机器评分 模型的研究和构建



高等教育出版社·北京

## 图书在版编目（C I P）数据

中国学生英译汉机器评分模型的研究和构建 / 江进

林著. -- 北京 : 高等教育出版社, 2016.1

ISBN 978-7-04-044316-5

I. ①中… II. ①江… III. ①英语—翻译—评分—机器识别—识别模型—研究 IV. ①H315.9-39

中国版本图书馆 CIP 数据核字 (2015) 第29714号

策划编辑 贾巍巍

版式设计 魏亮

项目编辑 王代军

责任校对 王春玲

责任编辑 刘瑾

责任印制 赵义民

封面设计 王洋

出版发行 高等教育出版社

社址 北京市西城区德外大街4号

邮政编码 100120

印 刷 北京市密东印刷有限公司

开 本 787mm×1092mm 1/16

印 张 12.25

字 数 224千字

购书热线 010-58581118

咨询电话 400-810-0598

网 址 <http://www.hep.edu.cn>

<http://www.hep.com.cn>

网上订购 <http://www.hepmall.com.cn>

<http://www.hepmall.com>

<http://www.hepmall.cn>

版 次 2016年1月第1版

印 次 2016年1月第1次印刷

定 价 28.00元

本书如有缺页、倒页、脱页等质量问题, 请到所购图书销售部门联系调换

版权所有 侵权必究

物 料 号 44316-00

本研究受到北京市哲学社会科学基金项目（批准号：15WYC064）、对外经济贸易大学学科建设专项经费（批准号：324-811005120501）和教学实验研究课题（批准号：X14520）的资助，以及对外经济贸易大学英语学院学术出版基金的资助。

# 序 言

已有的人工译文自动评分系统主要评价汉译英，极少数针对人工英译汉的自动评分研究还处于初等阶段。本研究旨在构建中国学生英译汉的机器评分模型，在人工评分标准、特征提取、语料和人机评分差异的分析方面与已有研究不同。首先，人工语义评分以原文的“翻译单位”为单元，翻译单位是符合搭配规则、意义单一、完整并具有一定区分度的多词单位，有利于评价译文的语义正误、语法性、连贯性、地道性等。人工形式评分增加了“风格切合度”标准，因为英译汉的目的语是学生的母语，译文的语言形式需要采用更高的评价标准。其次，本研究提取了翻译单位对齐数量等一些新特征。再次，以往研究采用单一文体，本研究对三种文体的译文分别建模，并确定了对每种译文评分最有效的质量预测因子。此外，本研究对人机评分差异较大的译文进行了质性分析，发现主要原因是机器依据的词表欠完美、语料预处理欠仔细。此发现为完善初始模型指明了方向。

创建评分模型需要五个步骤：语料收集、人工评分、特征提取、模型构建、模型验证。本研究收集了说明文、记叙文、叙议混合文译文各300多篇，译文评分按句进行，包括细致型和简化型两种。细致评分从语义和形式两个方面评价译文质量，简化评分仅对区分度较高的评分点进行语义评价。两种评分结果分别用于构建诊断目的和选拔目的的评分模型。评分结束后，提取多个语义和形式方面的文本特征，并在一半译文中进行多元线性回归分析，确定文本特征对人工评分预测力最大、共线性最小的方程。之后，利用方程计算另一半译文的得分，并考查机器与人工评分的相关性和一致性。最后，分析人机评分差异较大的译文，究其原因并提出改进措施。

本研究的主要结果如下：

1. 在三组语料的诊断性语义、形式评分模型及选拔性模型中，既有相同特征也有不同特征。其中，翻译单位对齐数量的贡献最突出。
2. 在诊断性模型中，篇章译文语义、形式评分模块的信度良好。交叉检验结果显示，三组语料篇章译文机器语义评分与人工语义评分的相关系数均值分别为0.846\*\*、0.881\*\*、0.925\*\*，机器形式评分与人工形式评分的相关系数均值分别为0.625\*\*、0.690\*\*、0.774\*\*。机器与人工评分比较一致。此外，与单句机器评分相加的篇章译文分数相比，机器对整体篇章的评分更稳定，更接近人工评分。

3. 以50、100、130、150、180篇训练集译文构建的选拔性评分模型都能较好地预测译文成绩。说明文和记叙文译文中130篇训练集、叙议混合文译文中100篇训练集所构建模型的机器评分与人工评分非常接近，选择此类数量的训练集不仅能够节约成本，还能满足大型评分的需要。

本书的主体部分是笔者的博士论文。衷心感谢我的导师文秋芳教授的精心指导。她踏实严谨的治学态度和批判性的思维能力一直鞭策着我，使我受益匪浅。感谢王立非教授、梁茂成教授对我论文撰写提出的宝贵意见。感谢参加我博士论文开题报告和预答辩的冯志伟教授、宋柔教授、陈国华教授、王克非教授、罗选民教授、王小捷教授，他们对本研究提出了很多精辟见解。感谢秦颖博士提供技术上的帮助。晏小琴、马晓雷博士、王金铨教授帮助评判学生译文，胡德香教授、熊兵教授、范娜、王东志博士提供专家译文，特此感谢。张莎、马玉学、陈功、蒋岳春、杨志红、刘长珍博士帮忙校对论文，一并致谢。

江进林

2015年12月7日

# 目 录

<b>绪 论</b>	<b>1</b>
0.1 引言	2
0.2 本研究的理论和实践意义	2
0.2.1 理论意义	2
0.2.2 实践意义	3
0.3 本研究概述	4
0.4 全书结构	5
0.5 小结	5
<b>第一部分 文献综述</b>	<b>6</b>
<b>第一章 自动评分系统综述</b>	<b>7</b>
1.1 现有自动评分系统	7
1.1.1 作文自动评分系统	7
1.1.2 翻译自动评价系统	9
1.1.2.1 机器翻译评价系统	9
1.1.2.2 人工译文评价系统	10
1.2 本研究与现有自动评分系统的区别	15
1.3 小结	17
<b>第二章 人工评分方法综述</b>	<b>18</b>
2.1 语言运用测试的人工评分方法和信度评价方法	18
2.1.1 人工评分方法	18

2.1.2 信度评价方法 .....	20
2.2 翻译质量评价方法 .....	21
2.2.1 翻译质量评价标准 .....	21
2.2.2 已有的翻译质量评价量化方法 .....	23
2.2.2.1 “信”的量化方法 .....	23
2.2.2.2 “达”的量化方法 .....	24
2.2.2.3 “切”的量化方法 .....	27
2.2.3 本研究补充的翻译质量评价量化方法 .....	28
2.2.3.1 “信”的补充量化方法 .....	28
2.2.3.2 “达”的补充量化方法 .....	28
2.2.4 本研究的量化方法小结 .....	29
2.3 小结 .....	31
<b>第二部分 研究设计 .....</b>	<b>32</b>
<b>第三章 研究问题与人工评分 .....</b>	
3.1 研究问题 .....	33
3.2 语料来源 .....	34
3.3 评分标准 .....	36
3.3.1 第一次评分标准 .....	36
3.3.1.1 语义评分标准 .....	37
3.3.1.2 形式评分标准 .....	42
3.3.2 第二次评分标准 .....	45
3.4 评分过程 .....	46
3.4.1 评分员选择 .....	46
3.4.2 评分员培训 .....	46
3.4.3 评分 .....	47
3.5 评分信度 .....	47
3.5.1 第一次评分信度 .....	47
3.5.2 第二次评分信度 .....	52
3.6 训练集、验证集、最佳译文集的形成 .....	53

3.7 小结 .....	54
<b>第四章 文本分析与数据分析 .....</b> 56	
4.1 研究流程概述 .....	56
4.2 研究工具 .....	57
4.2.1 文本预处理工具 .....	57
4.2.2 文本分析工具 .....	58
4.2.3 数据分析工具 .....	60
4.3 文本分析 .....	60
4.3.1 语义特征提取 .....	60
4.3.2 形式特征提取 .....	66
4.4 数据分析 .....	68
4.4.1 相关分析 .....	68
4.4.2 多元线性回归分析 .....	68
4.5 小结 .....	68
<b>第三部分 结果与讨论 .....</b> 69	
<b>第五章 译文质量预测因子 .....</b> 70	
5.1 译文语义质量预测因子 .....	70
5.2 译文形式质量预测因子 .....	71
5.3 文体与译文质量预测因子 .....	75
5.4 小结 .....	80
<b>第六章 诊断性英译汉评分模型 .....</b> 82	
6.1 诊断性英译汉评分模型的构建 .....	82
6.1.1 诊断性语义评分模型的构建 .....	82
6.1.2 诊断性形式评分模型的构建 .....	85
6.2 诊断性英译汉评分模型的验证 .....	90
6.3 以篇章和以单句为单位的篇章译文评分模型比较 .....	96
6.4 小结 .....	98

<b>第七章 选拔性英译汉评分模型</b>	<b>99</b>
7.1 选拔性英译汉评分模型概述	99
7.2 选拔性英译汉评分模型的构建	99
7.3 选拔性英译汉评分模型的验证	102
7.3.1 选拔性模型的自动评分信度	102
7.3.2 选拔性模型的自动排序信度	104
7.4 小结	106
<b>第四部分 结论</b>	<b>108</b>
<b>第八章 研究发现和价值</b>	<b>109</b>
8.1 研究发现	109
8.1.1 英译汉评分模型中的变量	109
8.1.1.1 研究证实的变量	109
8.1.1.2 研究改进的变量及提取的新变量	110
8.1.2 诊断性英译汉评分模型的性能	112
8.1.3 选拔性英译汉评分模型的性能	113
8.2 研究价值	114
8.3 不足和今后研究方向	116
<b>参考文献</b>	<b>118</b>
<b>附录</b>	<b>128</b>
附录一 记叙文、叙议混合文翻译题目和要求	128
附录二 记叙文、叙议混合文的翻译单位和连接词划分	131
附录三 说明文翻译单位和连接词的简化赋分	133
附录四 说明文翻译单位译文正误等级表	134
附录五 记叙文翻译单位译文正误等级表	138
附录六 叙议混合文翻译单位译文正误等级表	144
附录七 语言不合语法和不地道示例	148
附录八 三组语料的风格切合度示例	150

附录九 记叙文、叙议混合文评分点划分 .....	153
附录十 英汉词典示例 .....	154
附录十一 汉语高频字词表示例 .....	155
附录十二 关联词表示例 .....	156
附录十三 提取词类 .....	158
附录十四 语义变量与单句译文语义分数的相关关系 .....	158
附录十五 形式变量与单句译文形式分数的相关关系 .....	160
附录十六 单句译文语义评分模型 .....	165
附录十七 单句译文形式评分模型 .....	169
附录十八 篇章译文评分模型两次验证结果 .....	173
附录十九 单句译文评分模型两次验证结果 .....	173

## 表 目

### 第一章

表1-1 主要作文自动评分系统的特点 .....	7
表1-2 汉译英自动评分研究中的人工形式评分标准 .....	11
表1-3 汉译英自动评分研究中的变量 .....	12
表1-4 已有英译汉自动评分研究中的主要变量 .....	14

### 第二章

表2-1 整体评分法与分析评分法比较 .....	18
表2-2 拟提取的文本特征 .....	30

### 第三章

表3-1 说明文翻译题目和要求 .....	34
表3-2 语义评分标准 .....	37
表3-3 说明文的翻译单位和连接词划分 .....	40
表3-4 形式评分标准 .....	43
表3-5 说明文评分点划分 .....	45
表3-6 单句译文语义评分信度 .....	47
表3-7 单句译文形式评分信度 .....	48
表3-8 篇章译文语义、形式评分信度 .....	49
表3-9 评分员严厉度差异 .....	52
表3-10 评分员拟合数据 .....	52
表3-11 篇章译文第二次评分信度 .....	53
表3-12 两次平均评分的相关系数 .....	53

### 第四章

表4-1 形式特征及其提取工具 .....	66
-----------------------	----

### 第五章

表5-1 语义变量与篇章译文语义分数的相关关系 .....	70
-------------------------------	----

表5-2 词性分布与篇章译文形式分数的相关系数 .....	72
表5-3 词汇密度、词频广度、流利度、多样性与篇章译文形式分数的 相关系数 .....	72
表5-4 难易度、整体性与篇章译文形式分数的相关系数 .....	73
表5-5 句子层面的形式特征与篇章译文形式分数的相关系数 .....	74
表5-6 篇章层面的形式特征与篇章译文形式分数的相关系数 .....	75
表5-7 与篇章译文形式成绩相关的文本形式特征 .....	76

## 第六章

表6-1 篇章译文语义评分模型概况 .....	83
表6-2 篇章译文语义评分模型各系数的t检验结果 .....	83
表6-3 篇章译文语义评分模型的共线性数据 .....	85
表6-4 篇章译文语义评分回归方程 .....	85
表6-5 篇章译文形式评分的初次模型概况 .....	86
表6-6 篇章译文形式评分初次模型各系数的t检验结果 .....	86
表6-7 篇章译文形式评分初次模型的共线性数据 .....	87
表6-8 语料1篇章译文形式评分的第二次模型概况 .....	88
表6-9 语料1篇章译文形式评分第二次模型各系数的t检验结果 .....	88
表6-10 语料1篇章译文形式评分第二次模型的共线性数据 .....	89
表6-11 篇章译文形式评分回归方程 .....	89
表6-12 篇章译文机器语义、形式评分信度均值 .....	90
表6-13 机器/人工总分等级绝对一致的译文数量及其百分比 .....	91
表6-14 机器/人工总分等级绝对一致和相邻一致的译文数量及其百分比 .....	91
表6-15 机器/人工评分的差异显著性 .....	92
表6-16 机器/人工总分差异较大的译文数量 .....	93
表6-17 人机评分差异较大的原因及其改进方式 .....	93
表6-18 两种建模单位的语义、形式、总分评分信度比较 .....	96
表6-19 单句相加的篇章译文机器评分与人工评分的差异显著性 .....	97

## 第七章

表7-1 篇章译文选拔性模型概况（50篇训练集） .....	100
表7-2 篇章译文选拔性模型各系数的t检验结果（50篇训练集） .....	100
表7-3 篇章译文选拔性模型的共线性数据（50篇训练集） .....	101
表7-4 不同训练集的选拔性模型概况 .....	101
表7-5 不同训练集选拔性模型机器评分与人工评分的相关系数 .....	102
表7-6 不同训练集选拔性模型机器评分与人工评分的alpha值 .....	102
表7-7 不同训练集选拔性模型机器评分与 第二次人工评分的差异显著性 .....	103
表7-8 不同训练集选拔性模型机器排序与人工排序的相关系数 .....	105
表7-9 不同训练集选拔性模型机器排序与人工排序的alpha值 .....	105
表7-10 最佳选拔性评分模型 .....	106

## 第八章

表8-1 本研究改进的变量及提取的新变量 .....	110
----------------------------	-----

## 图 目

## 第三章

图3-1 翻译评分标准体系 .....	36
图3-2 各面的测量值概况 .....	51

## 第四章

图4-1 模型构建流程 .....	56
图4-2 模型验证流程 .....	57
图4-3 翻译单位词典示例 .....	63
图4-4 翻译单位对齐结果示例 .....	64
图4-5 初始矩阵 .....	65

# 绪 论

## 0.1 引言

自20世纪60年代以来，国内外已相继开发出多个英语作文自动评分系统，相关研究日臻成熟（Burstein 2003; Landauer et al. 1998, 2003; Page 1968, 2003; 梁茂成 2005）。近年来，少数研究者对中国学生译文的自动评分进行了尝试。其中，汉译英自动评分系统采用两种人工评分分别构建诊断性（diagnostic）和选拔性（selective）评分模型，前者可以对译文的语义、形式质量进行细致评分并提出反馈，后者满足大规模测试中汉译英自动评分的需要。两类模型的评分都比较接近人工评分，具有良好的效果（王金铨 2008）。不过，该研究的文体仅限于记叙文，且采用保留样本法（hold-out method）<sup>1</sup>验证模型，结果受到验证集译文的影响。在英译汉自动评分方面，已有研究挖掘了词对齐数量等特征，采用10折交叉检验（10-fold cross validation）<sup>2</sup>来验证模型，具有一定的优势。不过，该系统还处于初级阶段，比如语料仅涉及一个广告类段落；人工评分比较粗略；变量基本上属于词汇层面，且参考译文很少（王立欣 2007）。

基于以上不足，本研究将构建稳定、可靠的中国学生英译汉自动评分模型。研究采用细致型和简化型两种人工评分分别构建诊断性和选拔性评分模型。诊断性模型包括篇章、单句译文的语义、形式评分模块，构建单句模型的目的是对单句译文机器评分相加的篇章评分模型和整体篇章译文评分模型进行比较，选择篇章译文的最佳评分方法；选拔性模型包括篇章译文的语义评分模块。本研究在四个方面与已有研究不同。第一，人工评分在“忠实、通顺”的基础上，进一步考查译文的风格是否符合原文，并以“翻译单位”为单元对译文的语义质量进行穷尽性评价。第二，改进以往研究中的一些变量，并提取一些新的变量，比如翻译单位对齐数量。第三，为三种文体的译文分别构建评分模型。第四，采用2折交叉检验法验证模型，并深入分析人机评分差异较大的译文，究其原因并提出一系列减少机器评分偏差的措施。

总之，本研究将改进已有英译汉自动评分研究的诸多不足，并对汉译英自动评分系统的完善提供方法论指导，具有重要的理论意义。此外，研究构建的评分模型可以提高评分效率，帮助学生自学，具有很高的实践价值。

## 0.2 本研究的理论和实践意义

本节包括两个部分，分别概述本研究的理论和实践意义。

### 0.2.1 理论意义

本研究的理论意义主要有三点：

- 1 该方法将译文随机分为两半，一半进行训练，一半进行验证。
- 2 K折交叉检验是衡量机器学习结果推广性的常用方法，5折或10折交叉检验的平均结果通常最接近机器的真实性能（Breiman & Spector 1992）。