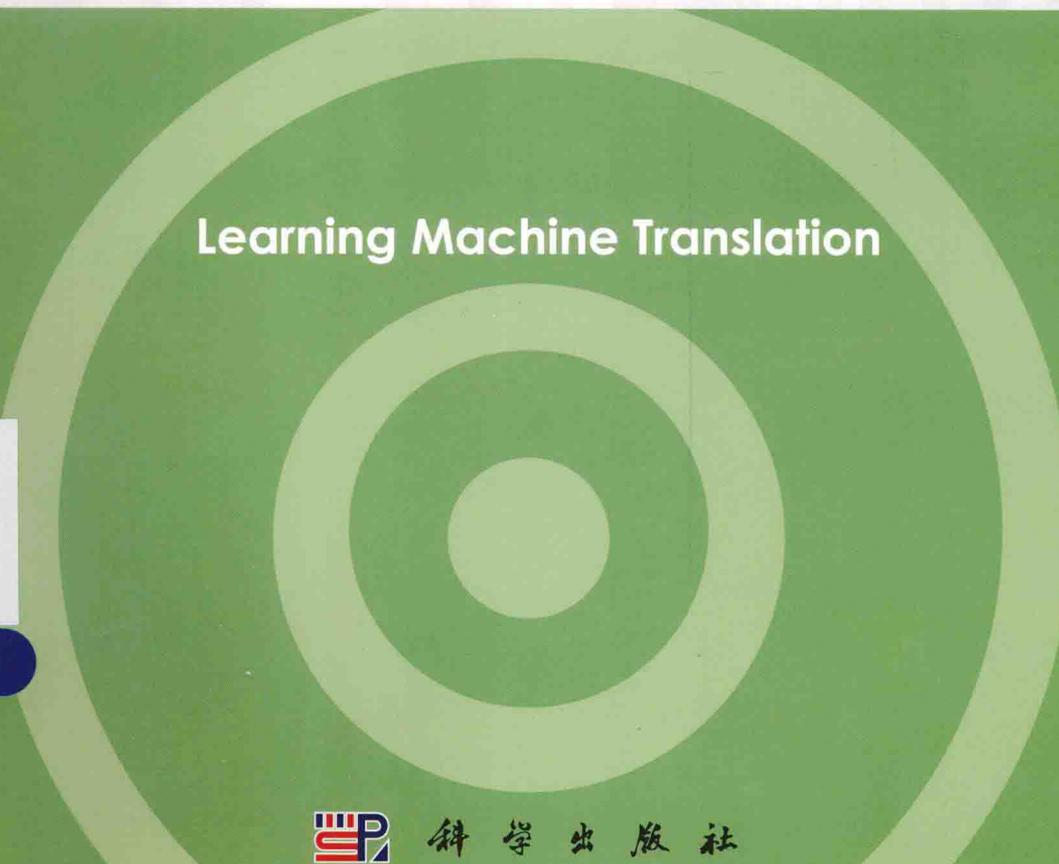




学习机器翻译

Cyril Goutte Nicola Cancedda 著
Marc Dymetman George Foster

曹海龙 赵铁军 译
朱聪慧 杨沐昀



Learning Machine Translation



科学出版社

学习机器翻译

Cyril Goutte Nicola Cancedda 著
Marc Dymetman George Foster

曹海龙 赵铁军 译
朱聪慧 杨沐昀

科学出版社

北京

图字:01-2013-5035

内 容 简 介

本书首先讨论若干使能技术,即解决那些不是机器翻译本身但却是与机器翻译系统开发相关的技术,其中包括从可比语料中获取双语句子对齐数据、多语名称词典的自动构造、词对齐技术等。随后介绍若干新的、改进的统计机器翻译技术,包括利用句法信息的判别式训练框架、半监督学习方法和基于核的学习方法的应用以及多机器翻译译文的组合以改进整个翻译系统的质量。

本书适合于从事机器翻译的研究者和研究生阅读,读者应具备统计机器翻译的入门知识。

Learning Machine Translation/Cyril Goutte, et al.

© 2009 Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means, including photocopying, recording, or by any information storage and retrieval system, without permission in writing from the publisher.

© 2014 Chinese Translation Publishing Science Press

图书在版编目(CIP)数据

学习机器翻译/(加)古特(Goutte, C.)等著;曹海龙等译. —北京:科学出版社,2014.10

书名原文:Learning Machine Translation

ISBN 978-7-03-042297-2

I. ①学… II. ①古…②曹… III. ①机器翻译-翻译机-基本知识
IV. ①H085②TP391.2

中国版本图书馆 CIP 数据核字(2014)第 250852 号

责任编辑:张艳芬 范 勃 / 责任校对:张小霞

责任印制:肖 兴 / 封面设计:陈 敬

科学出版社出版

北京东黄城根北街16号

邮政编码:100717

<http://www.sciencep.com>

北京佳信达欣艺术印刷有限公司印刷

科学出版社发行 各地新华书店经销

*

2014年10月第 一 版 开本:720×1000 1/16

2014年10月第一次印刷 印张:19

字数:362 000

定价:98.00 元

(如有印装质量问题,我社负责调换)

译者序

当今互联网时代对于机器翻译的需求越来越迫切,因为人们随时都可能浏览自己根本不懂的语言的网页。机器翻译这个极具商业和社会价值的研究领域受到了全球同行们的极大关注。改进技术使之更好地服务于人类需求,一直是研究者们孜孜追求的目标。在国际计算语言学学会的年度盛会(ACL)上,机器翻译始终是近年来发表文章最多的主题。

当前机器翻译的主流技术是统计机器翻译技术,其采用了大量机器学习特别是统计学习方法。统计机器翻译以大规模的双语数据为驱动,通过统计学习的手段挖掘出其中的翻译规律。这就是本书题名为“学习机器翻译”的缘故:计算机以“学习”为手段完成“机器翻译”。

本书是几年前一个研讨会论文的扩充,其中所涉及的问题虽有待继续解决,但所介绍的方法仍具有参考价值。本书除了第一章关于统计机器翻译的概述外,其余章节分为两部分:第一部分称为“使能”技术,这些技术虽然不直接服务于机器翻译系统的开发,却与统计机器翻译系统所依赖的基础资源相关,其中包括双语对齐语料的构建、多语名称词典的构建、命名实体识别、词对齐和语言模型构建技术等;第二部分介绍统计机器翻译在训练、调参、解码各个实现阶段的技术,包括判别式训练、翻译短语选择、基于核的翻译模型、句子重构翻译模型、半监督学习方法应用、翻译结果重排序和多机器翻译系统的译文融合等,多方实验表明这些技术对于改进整个统计翻译系统的输出质量有很大帮助。

本书的翻译分工如下:赵铁军负责第1、4、12章,曹海龙负责第2、6、7、9、10章,朱聪慧博士负责第3、5、11章,杨沐昀负责第8、13章。在本书翻译过程中,哈尔滨工业大学计算机学院机器智能与翻译实验室的多位研究生参与了初稿的翻译,他们是:刘淋(第2章)、崔一鸣(第3章)、李婷婷(第4章)、朱晓宁和崔一鸣(第5章)、赵弈欧和李晓倩(第7章)、张宇(第8章)、张文文(第9章)、张捷鑫(第10章)、史华兴(第12章)、朱俊国(第13章)等。

本书的出版得益于多方面的支持。感谢译者团队的成员侯亚楠同学,她为本书的校对、排版等付出了辛勤劳动。感谢哈尔滨工业大学教育部-微软语言语音重点实验室为本书出版提供了资助。在翻译过程中,译者承担了机器翻译方向的多项国家科研课题,包括国家863计划重大项目“互联网语言翻译系统研制”(2011AA01A207)、国家自然科学基金项目(61173073、61100093、61272384),这些研究工作无疑为翻译工作提供了坚实的知识基础。

原书前言

外语总是围绕在我们身边。现代通信技术使得人们能够获取自身并不完全理解的语言信息。数以亿计的互联网用户任何时间都可以在线,问题是语言障碍使得人们不能相互沟通。自动翻译的梦想激发着人们对自动或半自动翻译方法的持续兴趣。尽管自动翻译在有限领域和应用上取得了成功,千百万网页每天都被自动翻译,但是机器翻译却经常遭受那些并不欣赏产生全自动翻译这种挑战的人们的怀疑。但是,至少在现阶段,这是一个统计方法占据主流的快速发展和异常丰富的领域。

本书是2006年12月“神经信息处理系统会议”组织的用于多语言信息获取的机器学习研讨会的一个续篇。本书的几位作者在2006年的会议上就其工作作了介绍,但涉及的内容并不全面,而本书中约有半数内容是未公开发表的。

与以往的研讨会比较,本书也是坚持把统计机器翻译作为重点。本书分为两部分。第一部分涉及使能技术,即那些解决并不是机器翻译本身但与开发机器翻译系统紧密相连的问题的技术。例如,第2章涉及从可比语料库中获取双语句子对齐数据,这对于那些没有平行语料的领域或语言对是一个至关重要的任务。第3章和第4章探讨各种命名实体的多语言等价体的识别问题。该项技术的一个应用就是改进跨语言命名实体的翻译。第5章涉及词对齐,对于大多数统计机器翻译系统来说是一个基本的使能技术。它展示如何利用多预处理机制来改进对齐的质量。第6章展示词序核如何能被用于不同类型的语言学信息,并且指出这种方法可以改进判别式语言模型。

本书第二部分介绍新的统计机器翻译技术以及依靠统计方法或机器学习在现有技术基础上的改进。第7章探讨在翻译模型中纳入句法信息的问题,在一个判别式训练框架下进行参数估计。第8章提出一种在极大规模语料库上训练的统计机器翻译系统的假设输出重排序新方法。第9章介绍一种机器翻译新方法,它避开传统的特征函数对数线性组合,代之以基于核的方法(据作者所知,这是机器翻译领域涉及该方向的第一个研究)。第10章和第11章重点放在翻译模型的词或短语选择的改进上。在第10章中,基于源语言的全局信息,一个判别式过程决定目标语言词汇表中一个词是否出现在目标语言句子中,同时使用语言模型合作的加权转换机制为所选择的语言词汇进行排序。在第11章中,一个判别式短语选择模型被集成在一个基于短语的统计机器翻译系统中。这一章同时给出对大量的自动机器翻译评价标准的有趣分析和比较。第12章探索半监督学习的使

用,通过利用源语言中大量未翻译的材料来改进机器翻译输出。第 13 章展示多机器翻译系统的输出如何能够被组合起来,以便改进整体翻译质量。这种方法允许若干合作者贡献不同的机器翻译系统进行合作。系统融合目前在国际评测中产生了较好的结果。

作者期待本书对机器学习和统计机器翻译两个领域的研究者有所帮助。希望机器学习研究者将对统计机器翻译各个前沿方向有一个很好的整体了解,对于学习方法不同形式的直接影响力有所了解,并对一个具有挑战性的重要问题进行探索。也希望统计机器翻译研究者能够对本书产生兴趣,特别是对导论级统计机器翻译教材中可能没有覆盖的某些先进题目的介绍,对某些机器学习所激发的、具有产生本领域新方向潜力的新方法的介绍给予关注。

特别感谢 Susan Buckley 和 Robert Prior 在书稿撰写方面给予的支持。特别感谢为提高本书质量而提出意见和建议的以下人员: Caroline Brun, Marine Carpuat, Mauro Cettolo, Hal Daumé III, Hervé Déjean, Andreas Eisele, Alex Fraser, Patrick Haffner, Xiaodong He, Pierre Isabelle, Roland Kuhn, Dragos Munteanu, Miles Osborne, Jean-Michel Renders, Antti-Veikko Rosti, Craig Saunders, Libin Shen, Michel Simard, Sandor Szedmak 和 Dan Tufis。

Cyril Goutte

Nicola Cancedda

Marc Dymetman

George Foster

目 录

译者序

原书前言

第 1 章 统计机器翻译初步	1
1.1 背景	1
1.2 机器翻译的评价	3
1.2.1 基于编辑距离的方法	4
1.2.2 基于 n 元文法的方法	5
1.2.3 召回率的重要性	6
1.2.4 使用句法的方法	6
1.2.5 评价方法的评价与融合	7
1.2.6 统计显著性检验	7
1.3 基于词的机器翻译	7
1.3.1 模型 1、模型 2 和隐马尔可夫模型	8
1.3.2 模型 3、模型 4 和模型 5	9
1.3.3 搜索	9
1.3.4 现状	10
1.4 语言模型	10
1.4.1 n 元文法模型和平滑技术	11
1.4.2 最大熵模型	13
1.4.3 若干最新研究趋势	14
1.5 基于短语的机器翻译	16
1.5.1 对数线性模型	17
1.5.2 基于短语的翻译模型	17
1.5.3 最小错误率训练	19
1.5.4 搜索	20
1.5.5 重打分	22
1.5.6 现状	23
1.6 基于句法的统计机器翻译	23
1.6.1 无需句法分析的方法	24
1.6.2 目标语言端进行句法分析	25

1.6.3	源语言端进行句法分析	25
1.6.4	源语言端和目标语言端都进行句法分析	26
1.7	其他一些重要方向	27
1.7.1	因子化模型	27
1.7.2	模型自适应	27
1.7.3	系统融合	28
1.7.4	用于机器翻译的核方法	28
1.8	用于统计机器翻译的机器学习	28
1.8.1	翻译作为一个学习问题	29
1.8.2	使用不精确损失函数的学习	30
1.8.3	用于统计机器翻译的端到端学习	31
1.9	结论	32
1.10	附录	32

第一部分:使能技术

第2章	挖掘专利构建平行语料库	35
2.1	引言	35
2.2	相关工作	36
2.3	资源	37
2.4	对齐过程	38
2.4.1	句子对齐打分	38
2.4.2	降低句对齐中的噪声	40
2.5	专利平行语料库的数据统计	41
2.5.1	全集和源数据集的比较	41
2.5.2	基本的统计数据	42
2.5.3	关于机器翻译的统计数据	43
2.6	机器翻译实验	44
2.6.1	机器翻译系统	44
2.6.2	比较重排序限制	45
2.6.3	跨板块的机器翻译实验	46
2.6.4	对原始对齐数据的基于任务的评估	49
2.7	结论	51
第3章	多语言名称词典的自动创建	52
3.1	引言和动机	52
3.1.1	内容	53

3.1.2	专有名称和机器翻译	54
3.1.3	多语种名称实体词典与其他文本分析应用的相关性	54
3.1.4	存在名称变体的原因	55
3.2	相关工作	57
3.2.1	现有的名称词典或建立词典的相关探索	57
3.2.2	命名实体识别	58
3.2.3	名称变体的匹配	59
3.3	新名称的多语言识别	60
3.3.1	背景:多语言的新闻数据	60
3.3.2	一个允许多语言的轻量级识别过程	61
3.3.3	用维基百科扩充名称数据库	62
3.4	查找已知名称和其形态变体	62
3.4.1	处理词形变化	62
3.4.2	查找过程	63
3.5	人名识别的评价	65
3.6	名称变体的识别和合并	66
3.6.1	非罗马字符构成名称的音译	66
3.6.2	名称变体的“标准化”	67
3.6.3	(标准化)名称变体的近似匹配	68
3.7	总结与展望	69
第4章	多语料库中命名实体的音译和发现	71
4.1	引言	71
4.2	前人工作	73
4.3	协同排序:命名实体发现的一个算法	74
4.3.1	时间序列生成和匹配	76
4.3.2	音译模型	76
4.4	实验性研究	77
4.4.1	命名实体发现	78
4.4.2	初始例子集合规模	81
4.4.3	时间序列打分函数的比较	81
4.5	结论	82
4.6	未来工作	82
第5章	基于多预处理机制的统计词对齐融合	84
5.1	引言	84
5.2	相关工作	84

5.3	阿拉伯语的预处理机制	85
5.4	对齐的预处理机制	86
5.4.1	Giza++ 对齐	86
5.4.2	对齐重映射	87
5.5	对齐融合	87
5.6	评价	89
5.6.1	实验数据和评价指标	89
5.6.2	对齐重映射的贡献	90
5.6.3	融合特征的贡献	91
5.6.4	每个单一特征的作用	91
5.6.5	对齐合并实验	92
5.6.6	测试集评估	93
5.6.7	对齐规则分析	94
5.6.8	错误分析	95
5.7	后记:机器翻译和词对齐的改进	96
5.7.1	实验设置	97
5.7.2	结果	97
5.8	结论	99
第6章	用于判别式语言建模的语言学增强的词序列核	100
6.1	动机	100
6.2	增加语言学知识的词序列核方法	101
6.2.1	词序列核方法	101
6.2.2	因子化表示方法和核组合	103
6.2.3	因子化的核	103
6.2.4	实例说明	105
6.2.5	有理数核的解释	106
6.3	实验验证	107
6.3.1	各个因子上的核	108
6.3.2	因子的整合	109
6.3.3	与 n 元模型的比较	111
6.4	结论和未来的工作	113
6.5	附录	114

第二部分:机器翻译

第7章	走向树结构翻译模型的纯粹判别式训练	119
------------	--------------------------------	------------

7.1	引言	119
7.2	相关工作	120
7.3	学习方法	121
7.3.1	问题表征	122
7.3.2	目标函数	122
7.3.3	风险最小化	123
7.4	实验	127
7.4.1	数据	127
7.4.2	词转录	128
7.4.3	词包转录	131
7.4.4	树转录	133
7.5	结论	135
第 8 章	大规模统计机器翻译重排序	137
8.1	引言	137
8.2	背景	138
8.3	相关工作	138
8.4	我们的方法	140
8.5	实验 1: 汉译英系统的重排序	141
8.5.1	重排序器的训练	142
8.5.2	实验结果	142
8.6	实验 2: 法译英系统的重排序	145
8.6.1	实验结果	146
8.7	讨论	149
8.8	结论	150
8.9	附录	150
第 9 章	基于核的机器翻译	155
9.1	引言	155
9.2	统计机器翻译中的回归模型	156
9.2.1	岭回归	156
9.2.2	n 元语法字符串核	157
9.2.3	大规模训练	158
9.2.4	基于检索的稀疏近似法	158
9.3	解码	160
9.3.1	原像问题	160
9.3.2	柱搜索	160

9.3.3	复杂性分析	161
9.4	实验	162
9.4.1	语料	162
9.4.2	系统配置	163
9.4.3	岭回归实验	163
9.4.4	稀疏近似实验	165
9.4.5	搜索错误	166
9.5	进一步讨论	166
9.5.1	语言模型	166
9.5.2	语言学知识	167
9.6	小结	167
第 10 章	通过全局词汇选择和句子重构实现统计机器翻译	169
10.1	简介	169
10.2	SFST 训练和解码	170
10.2.1	单词对齐	170
10.2.2	双语言表示法	171
10.2.3	双语短语获取和局部重排序	172
10.2.4	SFST 模型	173
10.2.5	解码	173
10.2.6	单词插入模型	174
10.2.7	全局重排序	174
10.3	词汇选择判别模型	175
10.3.1	连续词汇选择模型	176
10.3.2	词袋词汇选择模型	177
10.4	选择分类器	177
10.4.1	多元与二元分类器对比	178
10.4.2	几何与概率解释	178
10.4.3	L1 与 L2 正则化	179
10.5	数据和实验	180
10.5.1	联合国和英国国会议事录语料	182
10.6	讨论	183
10.7	结论	184
第 11 章	统计机器翻译的判别式短语选择	185
11.1	引言	185
11.2	专用词语选择方法	187

11.3 判别式短语翻译	188
11.3.1 问题的设定	189
11.3.2 学习	189
11.3.3 特征设置	190
11.4 局部短语翻译	192
11.4.1 数据集及设置	192
11.4.2 评价	193
11.4.3 参数调整	193
11.4.4 性能比较	194
11.4.5 整体性能	195
11.5 为全局任务使用局部判别式短语翻译模型	197
11.5.1 基准系统	197
11.5.2 软集成判别式短语翻译的预测结果	198
11.5.3 设置	200
11.5.4 评价	200
11.5.5 参数调整	205
11.5.6 结果	206
11.6 结论	211
第 12 章 用于机器翻译的半监督学习	214
12.1 引言	214
12.2 基线机器翻译系统	215
12.3 框架	216
12.3.1 Yarowsky 算法	216
12.3.2 用于统计机器翻译的半监督学习算法	218
12.3.3 过滤器函数	218
12.3.4 估计函数	219
12.3.5 评分函数	219
12.3.6 选择函数	220
12.4 实验结果	221
12.4.1 设置	221
12.4.2 汉英翻译结果	223
12.4.3 法英结果	226
12.4.4 翻译例子	228
12.5 先前工作	229
12.6 结论与展望	230

第 13 章 学习系统融合机器翻译系统	232
13.1 引言	232
13.2 词对齐	234
13.2.1 问题表示	234
13.2.2 词对齐估计	234
13.2.3 词汇调序	237
13.2.4 相关研究中的其他对齐方法	238
13.3 CN 的生成和评分	239
13.3.1 建立 CN	239
13.3.2 概率估计	241
13.3.3 带有 R-best 译文翻译系统融合	242
13.3.4 共识翻译的抽取	242
13.3.5 语言模型重评分	243
13.3.6 保留单词的大小写信息	244
13.3.7 系统融合的参数优化	244
13.4 实验	244
13.4.1 翻译任务及环境	244
13.4.2 评价标准	245
13.4.3 对比实验	245
13.4.4 最终结果	247
13.5 结论	248
参考文献	249
中英文术语	279

第 1 章 统计机器翻译初步

Nicola Cancedda, Marc Dymetman, George Foster, Cyril Goutte

第 1 章主要是关于统计机器翻译(statistical machine translation, SMT)方面的一个简短介绍。特别是,覆盖了机器翻译(machine translation, MT)输出的自动评价、语言建模、基于词和基于短语的翻译模型以及 MT 中句法的应用。也对一些较新的方向做了一个快速综述,这些方向在未来具有重要意义。将 SMT 置于机器学习研究的一般情况下,并将重点放在它与标准机器学习问题和实践的相似和差异之上。

1.1 背景

MT 带着雄心勃勃的目标和未完成的承诺开始了其长期发展的历史。早期的工作一如当时的认识,以自动的或者“机械式”的翻译面貌出现,至少可以追溯到 20 世纪 40 年代。尽管其发展过程中存在少数磕绊如美国的 ALPAC 报告,但是 MT 依然以各种方式伴随着计算机科学和人工智能的进步而迅速发展(Hutchins, 2003)。

越来越强大的计算能力的提供使得 MT 的获取和使用更加直接。MT 也通过诸如搜索引擎服务这样的专门服务给公众以更广泛的展示。许多互联网用户熟悉至少一种如巴别鱼^①、谷歌语言工具^②、视窗在线翻译^③一类的工具。这些服务中的大部分由基于规则的系统所支持,如 SYSTRAN 所开发的系统^④,可是其中的某些系统(如谷歌和微软)至少部分地使用了统计方法。

在本章和本书的其他部分,翻译被定义为这样一个任务:将一个现存的源语言书写的文本转换为一个等价的不同语言即目标语言文本。传统的 MT(在本章论中指“统计之前的翻译”)依赖于源语言端不同层次的语言学分析和目标语言端的语言生成(图 1.1)。

^① <http://babelfish.yahoo.com/> or <http://babelfish.altavista.com/>。

^② http://www.google.com/language_tools。

^③ <http://translator.live.com/>。

^④ <http://www.systransoft.com/>。

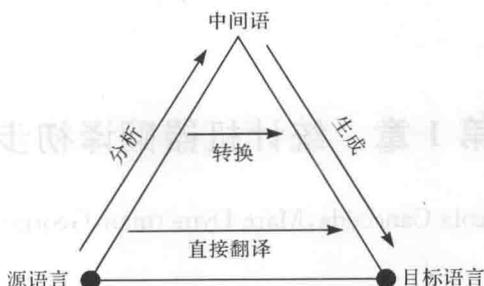


图 1.1 MT 金字塔

各种方法根据对分析和生成的依赖程度而不同：中间语言方法全部依赖分析和生成，而直接翻译方法最少依赖分析和生成，转换方法介于二者之间

在 20 世纪 80 年代后期，来自 IBM 的一个研究组开拓了第一个 MT 统计方法 (Brown et al, 1990)。事实上这可被看做是计算语言学总体潮流转移的一部分：在大约十年内，统计方法成为了这个领域中压倒性的支配方法。例如，在计算语言学协会 (Association for Computational Linguistics, ACL) 的年度会议论文集里的情况就是如此。

SMT 一般是指如何从源语言和目标语言的等价句对构成的大语料库中学习翻译。这是一个典型的机器学习框架：我们有一个输入 (源语言句子)，一个输出 (目标语言句子)，和一个试图为每一个给定输入产生正确输出的模型。

这里包含一系列关键问题，然而，其中一些关乎于 MT 应用。一个至关重要的问题是翻译质量的评价。机器学习技术典型地依赖于某些代价优化，以便学习输入和输出数据之间的关系。不过，自动地评价一个译文的质量或与某个给定的 MT 输出所相关的代价，是一个非常困难的问题。它可被归入于语言理解得更广泛问题，并且在很大程度上将停留在尚待解决的阶段。关于 MT 评价代价的定义与自动计算的困难，将在本章 1.2 节详述。

由 IBM 研究组倡导的 SMT 早期方法是信源通道方法。这基本上是两个模型组合的一个框架：一个基于词的翻译模型 (1.3 节) 和一个语言模型 (language model, LM) (1.4 节)。翻译模型保证了系统生成源语言句子对应的目标语言假设，而 LM 保证了输出的译文是尽可能符合语法和流畅的。

基于词的翻译模型取得了一些进展。然而，一个显著的突破却是采用对数线性模型和基于短语的翻译所带来的。这些内容将在 1.5 节中有更详细的描述。

虽然早期的 SMT 模型基本上忽略了语言学方面，但已有一系列努力尝试在翻译模型或者 LM 方面重新引入语言学思考。这些内容由 1.6 节和后续的某些章节所包含。此外，我们在 1.7 节给出了一个关于 SMT 目前趋势的概述，其中的一些内容在后面的章节中也会介绍。