

自然言語処理の基礎

工学博士 奥村 学 著

コロナ社

ま え が き

筆者が自然言語処理という言葉を知ったのは、大学の学部3年生のころであった。筆者が高校生のころはまだパソコンなどというものはなく、コンピュータというものはまだ世の中で一般的ではない時代に、私はなにを思ったか「コンピュータって面白そう」と勝手に思い込み、「情報工学科」に入学した。入学したものの、大型計算機でのプログラミングには必ずしもなじめず、四苦八苦し、卒業研究でもなにをテーマにしようか迷っていた。そのとき本で人工知能、とりわけわれわれの言語を理解するコンピュータの話を読んでこれだと思ったわけである。

筆者にとっての幸運は、4年生になった年に、当時出版された本で、自然言語処理の草分けとして名前のあった田中穂積先生が、東工大に着任されたことであった。いうまでもなく、卒業研究では田中先生の研究室を希望し、以後、現在までの26年間、自然言語処理というテーマで研究開発を行うことになった。

筆者の卒業研究は、(構文解析で用いる)日本語の文法を作成するというテーマであった。1980年代の自然言語処理はまだ小規模なシステムが多かったが、1章で説明するように、その後、筆者自身も自然言語処理技術の歴史的变化、大きな技術的な変革を直接経験することになる。筆者自身の研究の立場も、基礎的な自然言語を理解するという研究から、WWW上のテキストデータを対象とした応用研究の割合が少しずつ大きくなっていった。26年間の自然言語処理との付き合いの中では、研究開発した技術を基に会社を立ち上げるという貴重な経験もすることができた。いまや自然言語処理技術を基に会社だって設立できるのである!

また同時に、長年、教員として自然言語処理を若い学生さんに教えてきて、今回自然言語処理の教科書を書くという機会を得た。この教科書は、学部の3、

4年の学生さんが初めて自然言語処理について学ぶことを念頭に置いて書いたものである。また、半期の講義で扱える程度の内容を盛り込むように配慮したつもりでもある。伝統的に重要であると考えられている、人間の作成した知識を用いた自然言語処理手法とともに、コーパスから獲得した確率などを用いた、近年の自然言語処理手法についてもふれるように努力した。また、実用に近い自然言語処理技術の応用についても可能なかぎり説明するように努めた。

本書だけで自然言語処理のすべてがわかるということは、残念ながらないと思うが、本書の中で紹介している他の書籍なども参照して、さらに自然言語処理の理解を深めて下されば幸いである。

非常に残念だが、昨年7月われらが恩師である田中先生は、若くしてお亡くなりになった。本書を端緒として、自然言語処理って面白い! 自然言語処理って役に立つんだ! といったように、自然言語処理の面白さ、自然言語処理の可能性に若い読者が気づき、自然言語処理の世界にふれてくれるようになるなら、田中先生から自然言語処理の面白さ、難しさ、そして実用化の可能性について教えていただいた筆者にとっては、望外の喜びである。われわれの後をさらに引き継ぎ、自然言語処理の分野を引っ張ってってくれる若い研究者が、本書の読者から出てくれることを心から望んでいる。

この本を出版するにあたり、まずコロナ社に感謝します。なにかと遅れがちな原稿作成の過程で、つねに叱咤激励をいただきました。出版の日を迎えられたのは皆様の激励があったからこそともいえます。また、田中穂積先生がいらっしやらなければ今日の自分はなく、また本書はなかったといえます。ご冥福をお祈りするとともに、学部で4年生からずっとご指導、ご助言、ご支援下さったことに心から感謝いたします。ありがとうございました。また、本書の草稿に貴重かつ適切なコメントを下さいました、東京工業大学情報理工学研究所の徳永健伸教授、北陸先端科学技術大学院大学情報科学研究科の白井清昭准教授、広島市立大学情報科学研究科の難波英嗣准教授にも心から感謝いたします。

2010年8月

目 次

1. 自然言語処理概論

1.1 自然言語処理とは	1
1.2 自然言語処理の歴史	2
1.3 自然言語処理はなにに使えるか	3
1.4 自然言語処理の概要	7
1.5 本書の構成	9
章末問題	9

2. 辞書とコーパス

2.1 辞書	10
2.2 コーパス	13
2.3 言語の統計	16
2.4 機械学習を用いた自然言語処理	18
章末問題	21

3. 形態素解析

3.1 形態素解析とは	22
3.2 日本語の形態素解析	24
3.3 英語の形態素解析	33
章末問題	38

4. 構文解析

4.1 構文解析とは	39
4.2 文脈自由文法	41
4.3 構文解析アルゴリズム	43
4.4 CKY 法	48
4.5 チャート法	53
4.6 文脈自由文法の補強	61
4.7 構文解析における選好	63
4.8 まとめ	66
章末問題	67

5. 意味解析

5.1 意味解析とは	69
5.2 格フレームを用いた意味解析	73
5.3 コーパスを用いた語義曖昧性解消	78
5.4 まとめ	81
章末問題	82

6. 文脈解析

6.1 文脈解析とは	83
6.2 照応解析, 省略補完	84
6.2.1 照応, 省略に関する構文的, 意味的選好	88
6.2.2 焦点を用いた照応, 省略解析	89
6.2.3 照応解析の確率モデル	94

6.3 テキスト構造の構築	95
6.3.1 文間の意味的關係	96
6.3.2 表層的な情報を用いたテキスト構造の構築	101
6.3.3 テキストのセグメンテーション	102

7. 自然言語処理の応用

7.1 機 械 翻 訳	106
7.1.1 伝統的な機械翻訳方式	106
7.1.2 用例に基づく機械翻訳	109
7.1.3 統計的機械翻訳	110
7.2 自然言語処理とテキスト処理	111
7.3 情 報 検 索	113
7.3.1 いくつかの検索モデル	114
7.3.2 自動索引づけ	116
7.3.3 情報検索システムの評価	119
7.4 テキスト分類	122
7.5 情 報 抽 出	125
7.6 テキスト要約	133
7.7 質 問 応 答	138
7.8 ま と め	139
章 末 問 題	141
引用・参考文献	142
章末問題解答	146
索 引	153

1

自然言語処理概論

本章では、われわれ人間が日常的に話し聞き、読み書きしている言語をコンピュータ上で処理する技術である自然言語処理の概要について説明する。特に、自然言語処理の歴史、応用について説明する。

1.1 自然言語処理とは

自然言語 (natural language) とは、日本語、英語、フランス語など、われわれ人間が日常的に話し聞き、読み書きしている言語のことをいう。歴史的には、コンピュータにおける言語といえば、C, Ruby, Perl などのプログラミング言語、すなわち人工言語のほうが一般的であったため、人工言語と区別するため、ただ単に「言語」というのではなく「自然言語」と呼ぶこととなった。この自然言語をコンピュータ上で扱う技術を**自然言語処理** (natural language processing) という。

自然言語処理は、大きくつぎの二つの立場に分けることができる。

- コンピュータに「ことば」を理解させる。
- コンピュータ上で「ことば」を処理する。

次節以後で説明するように、自然言語処理の研究が始まった当初から、自然言語処理は、前者の立場に立って研究されてきたとあってよい。いわゆる伝統的な人工知能 (artificial intelligence) 的研究であり、コンピュータ上に知的な能力をもつものを実現したいという動機づけで研究されてきた。

一方で、1990年代に入ると、インターネットが急速に普及するのに伴い、

WWW上などに膨大な量のデータが蓄積されるようになり、その中には、テキストの形で蓄積されるデータもかなりの量を占めるようになった。このような形で、コンピュータ上には大量のテキストデータが蓄積されるようになり、このテキストデータをデータとして加工することで、なんらかの有用な情報を取り出したり、大量のデータを有効に活用する技術が開発されるようになってきた。これが後者の立場の自然言語処理といえるだろう。後者の立場の自然言語処理については、7章で詳しく説明する。

現在では、この二つの立場の研究が相補的に働き、自然言語処理の研究開発が活発に行われ、技術が深化するとともに、有用なシステムも続々と新たに生まれてきている。

1.2 自然言語処理の歴史

自然言語処理に関する研究開発は、コンピュータの開発とほぼ同時に開始されたといつてよい。コンピュータが1940年代に出現したのと^{とき}時を同じくして、機械翻訳システムの構想が生まれたというのは、今日から考えると驚くべき事実ではないだろうか。情報検索の歴史も同様に古く、1950年代には情報検索の研究が始まっている。日本では、1970年代に仮名漢字変換機能をもつワードプロセッサが登場し、以後、仮名漢字変換技術は、われわれにとって最もなじみのある自然言語処理の応用技術として用いられている。

1990年代に入ると、インターネットが急速に普及するのに伴い、WWW上には膨大な量のテキストデータが蓄積されるようになり、これらの膨大なテキストデータを有効に利用するための技術開発も進み、ユーザが欲しい情報を検索する情報検索技術や、逆にユーザが不要と考える情報をフィルタリングするフィルタリング技術、膨大なテキストデータの蓄積から、「おもしろい」情報を発掘しようとするテキストマイニング技術などが研究開発されるとともに、実用化されるに至っている。

一方で、自然言語処理の手法も同時期に大きな転機を迎えることになる。コ

コンピュータの HDD が安価になり、非常に大容量のテキストデータを保存できるようになり、また大量のテキストデータが容易に入手できるようになったことに伴い、それらの大量のテキストデータを「コーパス」として蓄積し、そのコーパスから自然言語処理に用いる知識として規則や確率を獲得し、それらを用いた自然言語処理の手法が生まれるようになってきている。本書では、2章でまずコーパスについて説明し、その後の章では、個別の解析技術においてどのようにコーパスが用いられるかについて説明する。

1.3 自然言語処理はなにに使えるか

本節では、自然言語処理の応用について、典型的なものをいくつか紹介しよう。前節で説明したように、自然言語処理の研究の当初から研究開発が行われていた機械翻訳 (machine translation) は、その代表例といってよいだろう。機械翻訳の手法については、7.1 節で簡単に説明する[†]。

また、これも前節で紹介したが、いまやわれわれにとってはなくてはならない入力手段である仮名漢字変換も、自然言語処理技術の典型的な応用の一つであろう。この技術のおかげで、キーボードから入力したローマ字あるいは仮名文字列は単語に分割されて漢字に変換され、われわれは現在何不自由なく日本語のテキストを入力できるようになっている。

もう一つ、わかりやすい自然言語処理の応用例を挙げよう。われわれが日常用いている言語でコンピュータやロボットとコミュニケーションができるようになると、よりそれらの機器はわれわれにとって親しみのあるものになる可能性がある。このようなコンピュータの機能は、対話システム、自然言語インタフェースなどと呼ばれる。これらのシステムの研究開発は、自然言語処理の歴

[†] 自然言語処理の書籍に、‘Time flies like an arrow.’ (「光陰矢の如し^{ごと}」) という文がよく例として登場するが、これは、古典的な機械翻訳システムが「時^{とき} 蠅^{ばえ}は矢を好む」と見事に誤訳する (「諺^{ことわざ}の翻訳というのは本来難しいのだが) ことにちなんていると思われる。

史の中でかなり初期の段階から行われている。

少し古い話になるが、ここで、歴史上に名前の残る二つの対話システムについて説明しておこう。一つは、1966年にワイゼンbaum (Joseph Weizenbaum) が開発した **ELIZA** (「イライザ」と発音する) である⁵⁶⁾†1。ELIZAは、ユーザがコンピュータ上でチャットで対話できるカウンセラ役をすることができる対話システムである。図 1.1 に ELIZA の対話例を示す。図では、全部大文字で表記されているのがシステムの発話である。英語なので難しいかもしれないが、図の対話例を見て、どのように感じられるだろうか。比較的自然的に対話が進んでいるように感じられた読者もいるのではないだろうか。

ELIZA は当時自然な対話が行え、チューリングテスト (Turing test)^{†2}にパスするという事で注目を集めたが、さらに驚かされたのは、このシステムが、きちんとした自然言語処理を行っておらず (1960年代であるから無理はないが)、単にユーザの発話に含まれる特定の単語列と照合するパターンを用意しておき、照合したパターンに応じて応答を出力するだけのものであったことである。例えば

ユーザの発話	システムの応答
i am X	how long have you been X?

というパターン (パターン中‘X’は変数を表す) をもっていれば、‘i am unhappy’ というユーザの発話とパターンの左側が照合するので、‘how long have you been unhappy?’ というシステムの応答を出力できるというわけである。

ELIZA に類する対話システムは現在でも広く世の中に存在する。例えば、Emacs/Mule から M-x doctor で実行可能なプログラムがそうであるし、日本

†1 肩付数字は、巻末の引用・参考文献の番号を表す。

†2 チューリングテストとは、コンピュータの原型を考えたチューリング (Alan Turing) が 1950年に提唱した、コンピュータの知能テストである。コンピュータに知性があるかどうかを判断するのに人間の被験者を用い、その被験者に端末を通して対話させる。人間の被験者が対話相手を人間と間違えるようならば、そのコンピュータは知性があると判断する。現在でも、チューリングテストを行うコンテストがレーブナー賞 (Loebner prize) という名前で行われている (<http://www.loebner.net/Prizef/loebner-prize.html>)。なお、本書に掲載する URL は執筆当時のものであり、変更される場合がある。

Men are all alike.
 IN WHAT WAY
 They're always bugging us about something or other.
 CAN YOU THINK OF A SPECIFIC EXAMPLE
 Well, my boyfriend made me come here.
 YOUR BOYFRIEND MADE YOU COME HERE
 He says I'm depressed much of the time.
 I AM SORRY TO HEAR YOU ARE DEPRESSED
 It's true. I am unhappy.
 DO YOU THINK COMING HERE WILL HELP YOU NOT TO BE UNHAPPY
 I need some help, that much seems certain.
 WHAT WOULD IT MEAN TO YOU IF YOU GOT SOME HELP
 Perhaps I could learn to get along with my mother.
 TELL ME MORE ABOUT YOUR FAMILY
 My mother takes care of me.
 WHO ELSE IN YOUR FAMILY TAKES CARE OF YOU
 My father.
 YOUR FATHER
 You are like my father in some ways.
 WHAT RESEMBLANCE DO YOU SEE
 You are not very aggressive but I think you don't want me to
 notice that.
 WHAT MAKES YOU THINK I AM NOT VERY AGGRESSIVE
 You don't argue with me.
 WHY DO YOU THINK I DON'T ARGUE WITH YOU
 You are afraid of me.
 DOES IT PLEASE YOU TO BELIEVE I AM AFRAID OF YOU
 My father is afraid of everybody.
 WHAT ELSE COMES TO MIND WHEN YOU THINK OF YOUR FATHER
 Bullies.
 DOES THAT HAVE ANYTHING TO DO WITH THE FACT THAT YOUR
 BOYFRIEND MADE YOU COME HERE

図 1.1 ELIZA の対話例 (文献 6) より引用)

で「人工無脳」と呼ばれていたプログラムは、この ELIZA が起源であるということができる。最近では、Twitter などに存在するボット (bot) と呼ばれるプログラムも同様のものである。

もう一つの対話システムは、1972 年にウィノグラッド (Terry Winograd) が開発した **SHRDLU** (シユルドウルと発音する) である⁵⁷⁾†。SHRDLU は図

† 英語の文字を頻度順に並べると

ETAOINSHRDLU...

となるところなどから、このシステムの名前が付けられたとされている。

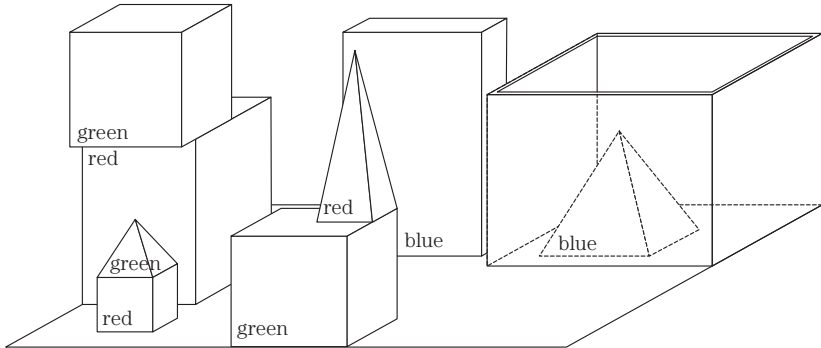


図 1.2 SHRDLU の積み木の世界

1.2 に示したような積み木の世界に存在するロボットアームに対して指示を出すことのできるシステムである。

仮想的につくられた積み木の世界にはロボット（アーム）があり、人間の自然言語による指示を理解し、積み木の世界を操作し、モニタ上で積み木の世界が変化する様子を表示する。例えば、図 1.2 の状況において

Pick up the blue box on the table.

と指示をすると、ロボットアームがテーブルの上の青の箱をもち上げ

Put it on the green box on the red box.

と指示すると、もっていた箱を（赤の箱の上の）緑の箱の上に置いてくれる。

指示を出した人間は、システムが言語による指示を理解したということ、モニタ上で世界の変化により（自分の出した指示のとおり世界が変化することで）確認することができる。SHRDLU は、積み木の世界という限定はあるものの、自然言語の指示を理解するコンピュータというものを明確に示した最初の例といえる。

最後に、自然言語処理に関連する、われわれにとってなじみ深いもう一つの応用として、Google, Yahoo! などに代表される検索エンジンを挙げておこう。検索エンジンの根幹をなす技術である情報検索は、自然言語処理とは別に研究開発が始まったが、いまやこの二つの技術は非常に近い位置にあるということ

索引

	<p>教師あり学習 18</p> <p style="text-align: center;">【く】</p> <p>句構造文法 41</p> <p>訓練データ 18</p> <p style="text-align: center;">【け】</p> <p>形態素 22</p> <p>形態素解析 7, 22</p> <p>形態素数最小法 29</p> <p>系列ラベリング 20</p> <p>結束性 103</p> <p>決定木学習 18</p> <p>原形の復元 23</p> <p>源言語 106</p> <p>言語モデル 111</p> <p>現代日本語書き言葉均衡 コーパス 15</p> <p style="text-align: center;">【こ】</p> <p>語彙的結束性 103</p> <p>語彙的トランスファ 107</p> <p>語彙的連鎖 104</p> <p>項 54</p> <p>構造的トランスファ 108</p> <p>構文解析 7, 39</p> <p>構文木 39</p> <p>構文構造 39</p> <p>構文的曖昧性 53</p> <p>後方照応 85</p> <p>語義曖昧性解消 71, 80</p> <p>コスト最小法 26, 29</p> <p>コーパス 3, 14</p> <p>固有名抽出 127, 139</p>	<p style="text-align: center;">【さ】</p> <p>再現率 119, 125</p> <p>再現率-精度グラフ 121</p> <p>最左導出 44</p> <p>最新性 88</p> <p>最長一致法 28</p> <p>索引語 114</p> <p>雑音のある通信路モデル 37, 110</p> <p>サポートベクトルマシン 18</p> <p>三角行列 48</p> <p style="text-align: center;">【し】</p> <p>指 示 103</p> <p>指示的 134</p> <p>自然言語 1</p> <p>自然言語処理 1</p> <p>シソーラス 12, 77, 103</p> <p>質問応答 112, 138</p> <p>自動索引づけ 114</p> <p>修辞関係 96</p> <p>修辞構造理論 96</p> <p>終端記号 41</p> <p>重要文抽出 135</p> <p>出現頻度 117</p> <p>照応解析 84</p> <p>照応詞 84</p> <p>焦 点 89</p> <p>情報検索 2, 112, 113</p> <p>情報抽出 112, 125</p> <p>省 略 103</p> <p>省略解析 84</p> <p>省略補完 85</p> <p>人工無脳 5</p>
<p style="text-align: center;">【あ】</p> <p>合図句 101</p> <p>曖昧性 8</p> <p>後戻り 43</p> <p>ア-リー法 61</p> <p style="text-align: center;">【い】</p> <p>依存構造解析 43</p> <p>一般化 LR 法 47</p> <p>意図の解析 84</p> <p>意味解析 7, 69</p> <p>意味素 77</p> <p>意味役割付与 73</p> <p style="text-align: center;">【え】</p> <p>衛 星 96</p> <p style="text-align: center;">【か】</p> <p>開始記号 41</p> <p>解 析 107</p> <p>概念階層 12, 77</p> <p>係り受け解析 43</p> <p>格 71</p> <p>核 96</p> <p>格フレーム 11, 73</p> <p>格文法 72</p> <p>確率文脈自由文法 64</p> <p>活性弧 55</p> <p>角川類語新辞典 13</p> <p>仮名漢字変換 2, 23</p> <p style="text-align: center;">【き】</p> <p>機械学習 18, 125</p> <p>機械翻訳 2, 3, 106</p>		

深層格	71				
		【す】	【ち】		【ね】
推論	132	茶 筌	33	根	40
		チャート	53		
		チャート法	53	【は】	
		チャンキング	20, 128	葉	40
正規文法	41	中間言語	109	バイグラム	16
生成	107	中間言語方式	106, 109	パーザ	39
生成規則	41	チューリングテスト	4	抜粋	135
精度	119, 125	著者同定	122	パラレルコーパス	14
制約	8	チョムスキー標準形	48		
接続	103			【ひ】	
遷移確率	37	【て】		非終端記号	41
選好	8	手がかり語	101	ビタビアルゴリズム	30, 36
先行詞	84	テキスト	22, 83	左再帰規則	45
センタ	89	テキストクラスタリング	112, 122	左隣構文解析	58
選択制限	77, 87	テキスト検索	139	必須格	71
前方照応	85	テキスト構造	84	ピボット方式	106
		テキスト構造解析	95	評価型ワークシヨップ	81, 139
【そ】		テキストセグメンテー		表層格	71
素性	62	ション	102	品詞付与	22, 33
素性構造	63	テキスト分類	112, 122		
		テキストマイニング	112	【ふ】	
【た】		テキスト要約	99, 112, 133	フィルタリング	122
代入	103			不活性弧	55
代表性	15	【と】		不要語リスト	116
タグ付コーパス	14	統計的機械翻訳	110	プーリアンモデル	114
縦型	43, 55	導出	44	文	22
単一化	63	動的計画法	30	文書類度	118
単一化文法	63	トップダウン	43, 55	文節数最小法	29
単語	22	トライグラム	16	文短縮	136
単語辞書	10, 24	トランスファ方式	106, 107	文法	40
単語出力確率	37			文脈	42
単語直接方式	106	【な】		文脈依存文法	41
単語分割	22, 24	ナイーブベイズ分類器	18	文脈解析	7, 83
談話	83	生コーパス	14	文脈自由文法	41
談話外照応	85			分類器	18, 125
談話セグメント	102	【に】		分類語彙表	13
談話単位	95	日本語語彙大系	13		
談話内照応	85	日本語 WordNet	13	【へ】	
談話マーカ	101	任意格	71	平均精度	122
				ベイズの定理	34

ベイズ分類	80		
ベクトル空間モデル	115	【や】	ラティス
変換	107	訳語選択	107
【ほ】		【よ】	【り】
報知的	134	要約率	133
補強文脈自由文法	62	用例に基づく機械翻訳	109
ボトムアップ	43, 55	用例に基づく手法	75
翻訳モデル	111	余弦類似度	116
【も】		横型	43, 55
目標言語	106	【ら】	類似度
		ラッパ	133
			【る】
			履歴リスト
			88
			【れ】
			類似度
			12, 75
			接続可能性行列
			24

【A】		ELIZA	4
abstract	135	extract	135
【B】		【F, G, H】	【R】
Balanced Contemporary		<i>F</i> 値	120
Corpus for Written		GETA	139
Japanese	15	Hyper Estraier	139
BCCWJ	15	【I, J, K】	right association
BNC	15	IDF	63
British National Corpus	15	JUMAN	13
Brown Corpus	15	<i>k</i> 近傍	96
【C】		KNP	67
Cabocha	67	<i>k</i> -NN	123
CFG	41	【L, M】	【S】
Charniak Parser	67	Lucene	5
CKY 法	48	MeCab	33
Collins Parser	67	【N】	SHRDLU
CSG	41	Namazu	90
【D】		N-best 探索	110
DF	118	n-gram	73
DP	30	【P】	SVM
【E】		PCFG	18, 125
EBMT	109	Penn Treebank	14
			【T】
			TF
			117
			TF 法
			117
			TF-IDF 法
			119
			【W】
			WordNet
			13
			WSD
			71
			【X】
			Xerox Tagger
			37

— 著者略歴 —

1984年 東京工業大学工学部情報工学科卒業
1989年 東京工業大学大学院博士課程修了（情報工学専攻）
工学博士
1989年 東京工業大学助手
1992年 北陸先端科学技術大学院大学助教授
2000年 東京工業大学助教授
2007年 東京工業大学准教授
2009年 東京工業大学教授
現在に至る

自然言語処理の基礎

Introduction to Natural Language Processing

© Manabu Okumura 2010

2010年10月28日 初版第1刷発行

検印省略

著者 ^{おく} 奥 ^{むら} 村 ^{まなぶ} 学
発行者 株式会社 コロナ社
代表者 牛来真也
印刷所 三美印刷株式会社

112-0011 東京都文京区千石 4-46-10

発行所 株式会社 コロナ社

CORONA PUBLISHING CO., LTD.

Tokyo Japan

振替 00140-8-14844・電話(03)3941-3131(代)

ホームページ <http://www.coronasha.co.jp>

ISBN 978-4-339-02451-7 (金) (製本：牧製本印刷)

Printed in Japan



無断複写・転載を禁ずる

落丁・乱丁本はお取替えいたします