

# 機械学習とは何か？

## 背景と概要

# ◆世界的企業によるAIの取り組み

## トップ5社のCEO発言

### ▶ トップ5社のCEO発言 図表02-5

<b>Google CEO</b> スンダー・ピチャイ (2017年5月 Google I/O '17)	モバイルファーストから AI ファーストにシフトする
<b>Amazon CEO</b> ジェフ・ベゾス (2016年株主への手紙)	我々は今まさに、明確な力強いトレンドのさなかにある。AI と機械学習だ
<b>Facebook CEO</b> マーク・ザッカーバーグ (2018年2月の Facebook 投稿)	AI を用いた我々のゴールは、よりよいサービスを作るために、Facebook 上のすべてのコンテンツを理解できるようにすることだ (図表02-6)
<b>Apple CEO</b> ティム・クック (2017年8月の決算説明会)	自律的なシステム (autonomous systems) の巨大なプロジェクトを動かしていて、莫大な投資をしている。乗り物もその使い方の1つだが、ほかにもさまざまな使い道のある、あらゆる AI プロジェクトの母親のようなものに取り組んでいる
<b>Microsoft CEO</b> サティア・ナデラ (2017年のアニュアルレポート)	あらゆる開発者が AI 開発者に、あらゆる企業が AI 企業になれるよう、我々の AI に関するケイパビリティを使って特徴あるポジショニングを取り、AI を民主化する

# ◆ トップ企業にみる機械学習の取り組み

- Google

- 🔗 自動運転の紹介動画

- ▶ <https://waymo.com/>

- 🔗 医療・科学・芸術ほか

- 医療 ▶ <https://ai.googleblog.com/2017/03/assisting-pathologists-in-detecting.html>

- 科学 ▶ <https://www.blog.google/topics/machine-learning/hunting-planets-machine-learning/>

- 芸術 ▶ <https://experiments.withgoogle.com/ai/ai-duet/view/>

- Amazon

- 🔗 無人配達

- ▶ <https://www.youtube.com/watch?v=vNySOrl2Ny8>

- 🔗 レジ無し店舗

- ▶ <https://www.youtube.com/watch?v=NrmMk1Myrxc>

- Facebook

- 🔗 SNSコンテンツの理解

- ▶ <https://research.fb.com/downloads/detectron/>

- Apple

- 🔗 ユーザ体験の向上

- Microsoft

- 🔗 顧客企業の支援

# ◆機械学習が注目される理由

## ・ 注目されるきっかけとなった事例

### (1) 画像認識

☞画像認識の世界的なコンペティション（ILSVRC：ImageNet Large Scale Visual Recognition Challenge）において、2012年にディープラーニング（深層学習）という方法を用いて大幅な画像認識の精度の改善が達成された。  
LSVRCでは任意の画像に何が写っているかを判別し、エラー（不正解）率の低さを、世界中の企業や大学の研究者たちが競い合う。  
ディープラーニングの手法を用いたチームは、それまで約50年の人工知能研究の成果で到達した26%のエラー率を、一気に15%にまで改善した。

### (2) 囲碁

☞2016年にGoogle傘下のDeepMindの開発した「AlphaGo」が、囲碁の世界チャンピオンであるイ・セドル氏を破った。

<https://www.youtube.com/watch?v=WXuK6gekU1Y>

AlphaGoには、ディープラーニングだけでなく、強化学習という機械学習の手法が用いられている。

<https://ja.wikipedia.org/wiki/AlphaGo>

## ・なぜ機械学習が成果を生むようになったのか？

### アルゴリズムの進化

- ・ディープラーニングや強化学習など、応用可能性と実行力の高い機械学習の手法の進化が続いているため

### データ量の増大

- ・インターネットの発展により、帯域の増強と共に画像、映像、音声、テキストなどさまざまなデータが増大しているため
- ・企業における業務のシステム化やセンサーの普及により、さまざまな種類のデータ量が増大しているため

### 計算資源の進化

- ・コンピューターの処理能力が格段に高まっているため（スパコンの処理能力の向上、GPU、TPU などの開発・普及）

### アルゴリズム、データ、 計算資源の利用 可能性の向上

- ・オープンソースのライブラリや TensorFlow などのツールにより、機械学習やディープラーニングのアルゴリズムが簡単に使えるようになっているため
- ・各種のデータソースが整備され、学習のためのデータが用意しやすくなっているため
- ・クラウドサービスの普及により、高性能の計算資源が安価に使えるようになっているため

# ◆機械学習ライブラリ・フレームワーク

## 1. NumPy

### 🔗NumPyの特徴

- 数値計算を行うための定番ライブラリの1つで、「ナンパイ」または「ナムパイ」と読む。
- 機械学習だけでなく、多言語配列や画像処理・音声処理にも活用できる利用頻度の高いライブラリ。
- 数値計算を得意とするライブラリだが、SciPyを追加する事で、さらに高度な科学計算を処理することが可能。

### 🔗NumPyの便利な点

- NumPyは数値計算の中で、特に配列処理能力に優れている。Pythonでも計算することは可能だが、インタプリタ型のプログラミング言語のため、実行速度が遅く、処理に時間がかかる。NumPyはC言語やFortranといったコンパイル型言語で実装されているため、処理速度が高速。
- データ分析において、データを受け渡すためのデータコンテナとしての役割もある。
- NumPyにC言語呼び出しのAPIがある事で、NumPyから外部ライブラリへデータを受け渡すことができたり、逆に外部ライブラリの計算結果をNumPyに戻すことが可能。
- この機能によって、NumPyは既存ライブラリを簡単に呼び出せる、動的インターフェースとしての役割を担ってる。

### 🔗NumPyの利用シーン

- ベクトルや行列といった多次元配列をPythonより早く効率的に処理する必要があるプログラムで利用される。
- 大規模のデータ処理に優れているだけでなく、他ライブラリと連結することがある。例えば「Pandas」「SciPy」「Scikit-learn」といったライブラリと連結して使用することが実際の業務で多数発生する。機械学習の世界では、かなり高い確率で使うことになる。

## 2. Pandas (パンドス)

### 🔑Pandasの特徴

- Pandasは、数表および時系列のデータ操作やデータ構造を変更するなど、テーブルデータを取り扱えるようにするライブラリ。
- 機械学習において、大量のデータをAIに学習させるが、最終的なモデル精度をより高めるためには、不要なデータを取り除き、必要なデータを精査する前処理をする必要がある。
- Pandasはその前処理の際に、データセット処理を効率化するために使用する。
- Pandas は重要なコードをC言語で実装しているため、NumPy同様、高速に処理することが可能。

### 🔑Pandasの便利な点

- NumPyで作成されたデータの加工や入出力が可能。
- 数値以外のデータ処理を行う場合にPandasを使うと便利。
- SQLと似た操作でデータ加工が行えるためデータベースを触っている人は使いやすい。
- 値に対するラベル付けを簡単に行えるため、機械学習の前処理などの工程を効率的に行うことが可能。
- 合計や平均だけでなく、分散や標準偏差といった機械学習において必須となる統計量も簡単に処理できる。
- 機械学習でよく行われる特定条件のデータを除外する処理も、SQLを書くように記述できる。
- 他にもソート、欠損値の補完、グラフの描画等、テーブルデータを取り扱う上で必須のメソッドを多く備えている。

### 🔑Pandasの利用シーン

- 一般的な表計算から、統計量の算出、データ整形、csv等のさまざまなフォーマットでの入出力といった、テーブルデータを扱う場合に、その豊富な機能面から利用されるデータ分析に必須のライブラリ。
- 金融データを取り扱うために最も適した時系列分析機能をもっているため、金融データ分析アプリケーションにも利用されている。
- 機械学習でデータ分析を行う前処理であるデータの読み込みやクリーニングといった作業でよく利用される。

### 3. SciPy (サイパイ)

#### 🔗SciPyの特徴

- SciPyは、信号処理や統計などの科学計算用のライブラリ。
- 語源は Science + Python に由来。
- SciPyではNumpyで行える配列や行列の演算はもちろん、さらに信号処理や統計といった計算ができるライブラリ。
- NumPyよりも高度な数値計算処理を行う場合は、SciPyを利用する。

#### 🔗SciPyの便利な点

- SciPyは統計、最適化、補完、積分、線形代数、フーリエ変換、信号処理、画像処理、遺伝的アルゴリズム、ODEソルバ、特殊関数といった、高度な科学技術計算処理をPythonを使って実行可能にする。

#### 🔗SciPyの利用シーン

cluster：階層的クラスタリング、ベクトル量子化、K平均 fftpack: 離散フーリエ変換

integrate：数値積分ルーチン、微分方程式ソルバ

linalg：線形代数ルーチン

io:データ入出力

maxentropy：エントロピー分布処理

ndimage:画像処理

sparse:疎行列処理

stats：統計処理、離散分布、連続分布

signal：信号処理ツール

tscipy:画像配列操作

weave：PythonにC言語やC++言語を組み込む処理



## 4. Matplotlib (マットプロットリブ)

### 🔗 Matplotlibの特徴

- ・データをグラフや画像データとして表示することができるブラフ描画のためのライブラリ。
- ・機械学習では、統計量の可視化や学習経過のグラフ化、画像の出力等の機能が多く利用されている。
- ・ヒストグラムや散布図を描いたり、JavaScriptを利用してインタラクティブなグラフを生成することも可能。
- ・Pandasでもデータの可視化は可能だが、Matplotlibを利用する事で更に複雑な表示が可能。

### 🔗 Matplotlibの便利な点

- ・出版用にも使えるほど高品質なグラフを作成することができる。
- ・画像をPDFやJPEG、GIF等であらゆる形式でエクスポート可能。

### 🔗 Matplotlibの利用シーン

- ・NumPyなど他ライブラリで前処理したデータを可視化する場合によく使われる最も一般的なPythonライブラリ。
- ・データを可視化する事で、異常値の検出や必要データの変形を発見することができるため、データ分析の現場ではよく使われる。
- ・機械学習でも、NumPyやPandasと組み合わせ、分析対象のデータにどんな傾向があるのかグラフで表示したり、スコアの変化を描画する等、データを図表で表示させる際に活用されている。

## 5. scikit-learn (サイキット・ラーン)

### 🔗scikit-learnの特徴

- ・機械学習全般のアルゴリズムが実装された機械学習の基盤となっている大人気のライブラリ。
- ・統計学、パターン認識、データ解析の技法が豊富に使うことができるので、特に研究者の間で人気がある。
- ・NumPyやSciPy、matplotlibと比較してもscikit-learnは、様々な機械学習の実装を簡単に試すことができる。

### 🔗scikit-learnの便利な点

- ・機械学習全般のアルゴリズムを簡単に実装できる。
- ・ただしどのアルゴリズムを選ぶべきか検討する際に、[チートシート](#)が参考になる。
- ・データの状況に応じて細かく分岐されたチートシートを活用する事でアルゴリズム選択が容易になる。

### 🔗scikit-learnの利用シーン

- ・主に分類・回帰・クラスタリング・次元削減という4つの目的に応じて利用される。

#### 分類：

ラベルとデータを学習し、ラベルがわからないデータにラベルを付けること。

例えば、身長体重とシワの数から性別を識別するというようなことに使われるイメージ。

#### 回帰：

数値をデータで学習し、その学習モデルを利用して、数値を予測すること。

例えば、身長体重かとシワの数から年齢を予測するといったことに使われるイメージ。

#### クラスタリング：

漠然としたデータから近い特徴を見つけて分類し、その学習モデルを利用して、同じようなデータからグループを見つけ分類すること。

例えば、星を光の強さや大きさ、色、距離、温度といったいくつかの特徴をもとにグループに分類するときに使われるイメージ。

#### 次元削減：

データの次元を削減しデータの可視化、圧縮化を行うこと。

イメージでいうと、3次元である現実世界から写真を撮ることで2次元に落とし込み、状況を判断するようなことをデータでも行うという感じ。

# ◆ディープラーニングフレームワーク

## 1. TensorFlow

- ・TensorFlowはGoogleのGoogle Brainチームによって2015年に開発。
- ・ディープラーニングフレームワークの中で現在最も人気のあるフレームワーク。
- ・Googleが開発したフレームワークなので、GmailやGoogle翻訳などに使われている。

主な特徴としては、

- ・C++とPythonで書かれている
- ・ドキュメンテーションが充実している
- ・iOSやAndroidなどのモバイルのプラットフォームでも利用できる
- ・Tensorboardと呼ばれる、訓練過程を観察できる機能がついている

🔗公式サイト：<https://www.tensorflow.org/>

## 2. PyTorch

- ・PyTorchはFacebookの人工知能グループによる開発からスタートし、2016年にリリース。
- ・上で述べたTensorFlowの人気に最も近い。

主な特徴としては、

- ・構文がわかりやすい
- ・研究者の間で人気

論文で発表されたモデルの実際のコードは、PyTorchで書かれたものが比較的に見つけやすい

- ・pdbやPyCharmなどのデバッガーが使える

動的な計算グラフ（学習速度は落ちるものの、柔軟なモデルの構築が可能）

🔗公式サイト：<https://pytorch.org/>

### 3. Keras

- ・ KerasはFrancois Chollet氏（現在はGoogleのエンジニア）を中心として2015年に開発。
- ・ 現在はTensorFlowに取り込まれ、tf.kerasの形で使われるのが一般的。
- ・ 開発者のFrancois Chollet氏自身が執筆したKerasの解説本も出版されている。

主な特徴としては、

- ・ モデルの構築がとても簡単で、初心者にとってもわかりやすい。
- ・ 素早く実装できるので、プロトタイプ作成時に便利。
- ・ TensorFlow、CNTK、Theanoという複数のバックエンドをサポート。

🔗公式サイト：<https://keras.io/>

### 4. Caffe

- ・ Yangqing Jia氏が、カリフォルニア大学バークレー校の博士課程在学中に始めたプロジェクトで、同校のBAIR（Berkeley Artificial Intelligence Research）を中心に開発が行われ、2017年にリリース。

主な特徴としては、

- ・ 画像認識などの処理が得意
- ・ 開発コミュニティが活発
- ・ 動作が速い

🔗公式サイト：<https://caffe.berkeleyvision.org/>

## 5. Microsoft Cognitive Toolkit

- ・元々CNTKという名前で2016年にリリースされたが、同年10月に現在の名前に変更された。
- ・Skypeなどにも使われてる。

主な特徴としては、

- ・リソース効率が良い
- ・ONNXフォーマット（ディープラーニングモデルを異なるフレームワーク間で交換するためのフォーマット）を初めてサポート
- ・巨大データセット処理時のパフォーマンス低下を最小化するためのアルゴリズムが組み込まれているので、複雑マシンで巨大データセットを扱える
- ・コミュニティはそこまで発展していない

🔗公式サイト：[https://www.microsoft.com/en-us/research/product/cognitive-toolkit/?lang=fr\\_ca](https://www.microsoft.com/en-us/research/product/cognitive-toolkit/?lang=fr_ca)

## 6. MxNet

- ・MxNetは、CMU（カーネギーメロン大学）、NYU（ニューヨーク大学）、NUS（シンガポール国立大学）、MIT（マサチューセッツ工科大学）など、様々な大学からの研究者が協力して開発され、2016年にAWS（Amazon Web Service）によるサポートが発表された。

主な特徴としては、

- ・高いスケーラビリティ
- ・モバイルデバイスも対応可
- ・Python, R, Scala, JavaScript, C++など多くの言語に対応
- ・研究者の間での人気はあまり高くなく、コミュニティもそこまで発達していない
- ・命令的プログラムと宣言的プログラムの併用が可能

🔗公式サイト：<https://mxnet.apache.org/versions/1.7/>

## 7. Chainer

- ・日本のベンチャー企業である、PFN（Preferred Networks）によって開発。
- ・2019年12月にその開発を終了し、PyTorchへ移行されることが発表されたが、define-by-runというアプローチを他のフレームワークに先駆けて提唱するなど、後続くフレームワークにも少なからず影響を与えた。

主な特徴としては、

- ・国内での人気が高い
- ・TensorFlowのように計算グラフを定義してから計算を実行するのではなく、計算と同時に計算グラフを定義するので、モデルの再構築が楽（これがPyTorchなどの動的計算グラフにつながる）

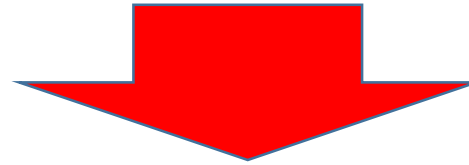
🔗公式サイト：<https://chainer.org/>

## ◆機械学習の定義

機械学習

=

データからパターンやルールを  
機械自身に見つけさせる仕組み



### ▶ パターン、ルール抽出のイメージ 図表10-1

売上レポート



店舗ごとの  
仕入れ・売上  
日ごとの仕入・売上



曜日ごとの  
販売数に傾向が  
ありそう

機械にパターンやルールを発見させる

## ◆機械学習の主な用途

機械学習が最も多く  
活用される場面



識別



予測

判別済みのデータ（「犬」というラベルがついた画像）と未判別のデータ（ラベルがついていない画像）があるという状況で、未判別のデータを正しく識別するというタスク

将来の販売数など、過去のデータにもとづいて将来どれくらい売れそうかなどの「観測されていない未来の値を予測するタスク」



# ◆ルールベースと機械学習の違い

ルールベース



機械学習ではなく、人が決めたルールにもとづく作業効率化の仕組み

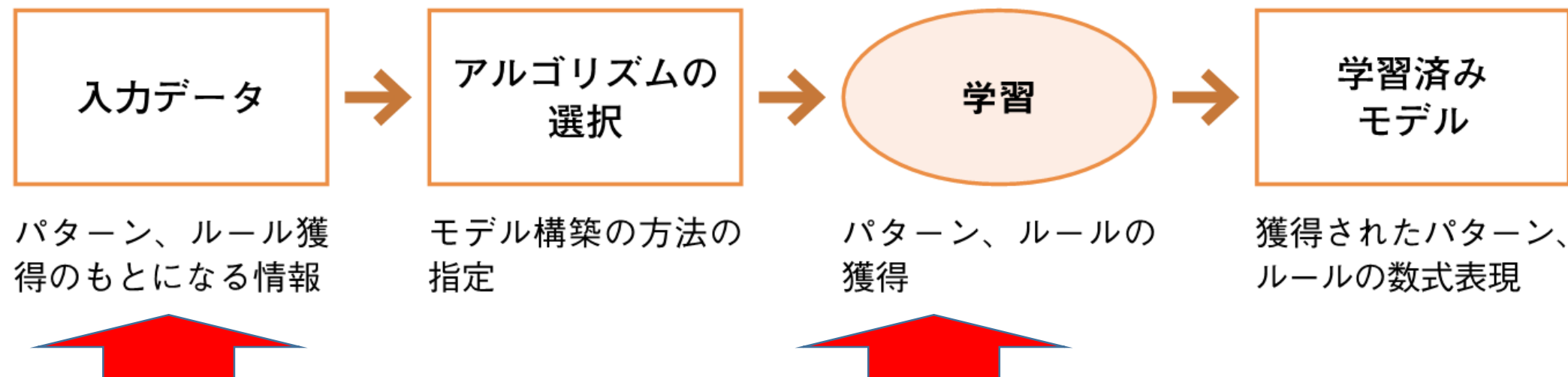
ルールベースは非常に柔軟でわかりやすく、精度をコントロールしやすい一方で、ルールが多くなると管理が難しくなるうえに、人が気づかないパターンはルールに取り込めないというデメリットもある。

## ▶ 機械学習とルールベースの比較 図表11-1

	機械学習	ルールベース
概要	データからルールを見つける	人がルールを決める
メリット	データの特徴に従い、アルゴリズムによって何らかの数学的根拠をもとにルールを見つける	<ul style="list-style-type: none"><li>・ 自明なルールを決めることが非常に簡単</li><li>・ ルールを複数組み合わせることで複雑なルールを作ることも可能</li></ul>
デメリット	<ul style="list-style-type: none"><li>・ 一定以上のデータがないと識別や予測に効果的なルールを見つけられない</li><li>・ 課題やデータに対して適切なモデルやアルゴリズムを選択する必要がある</li></ul>	<ul style="list-style-type: none"><li>・ ルールのメンテナンスが大変</li><li>・ 人が認識している以上の複雑で精緻なルールは作れない</li></ul>

# ◆機械学習によるパターンとルール抽出

## ▶ 機械学習によるモデル構築のイメージ 図表11-2



パターンやルールを発見するもとになる「データ」を用意

「機械学習アルゴリズム」を用いて「モデル」を構築する。  
「モデル」：  
パターンやルールを数式で表現したもの  
「機械学習アルゴリズム」：  
「モデル」を構築するための一連の数学的な処理。

注) 機械学習とルールベースは組み合わせてもよい

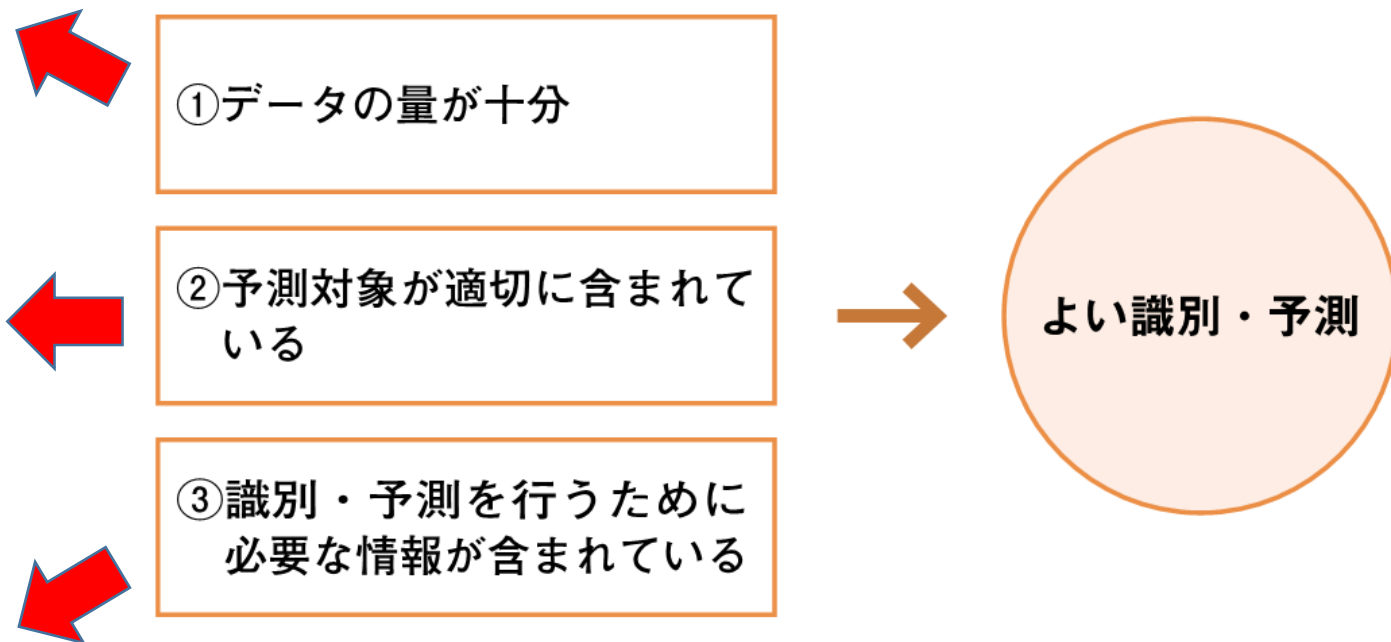
# ◆機械学習が強みを発揮するにはデータが重要

データ量が少ないと信頼できる結果は得られない。たとえば、1カ月後の売上予測を行いたいのに、先週と今週のデータを見て増加傾向にあるからといって、1カ月後もこの調子で増加すると考えるのは早計

気温に応じて売れやすくなる商品があったときに、気温の情報がデータとして取得されていなければ、気温に関連するルールを見つけられない。

猫と犬の識別を行いたいのに猫のデータがなければ適切な識別は行えない。同様に、ある商品を購入しそうな人を予測したいのに、購入者のデータしか取得できない場合も正しい予測は行えない。

## ▶ 機械学習を行うために必要なデータの要件 図表11-3



# ◆機械学習に不向きなこと

## ▶ 予測が困難な課題の例 図表11-4

- ・ 偶発的に起こるもの  
サイコロの出目、ルーレットの出目、コイントスの裏表、etc
- ・ 現象が起こるメカニズムが複雑または現象を説明するデータが十分に取得できないもの  
地震、メガヒット商品の発生、超高額購入顧客の発生、etc
- ・ 過去のデータがないもの  
新しい施策の効果、新商品の販売数、新しい施設の売上、etc

# ◆機械学習の目的 ①自動化によるコスト削減

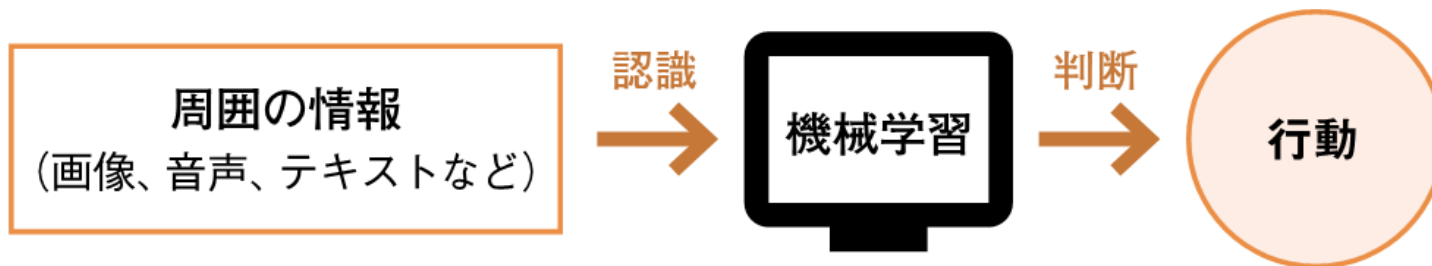
機械学習は「これまで人が判断してきた作業を自動化できる可能性がある」という点が非常に優れている。

## ▶ 人が行う認識、判断のプロセスと機械が行う認識、判断のプロセス 図表12-1

### ・ 人の場合

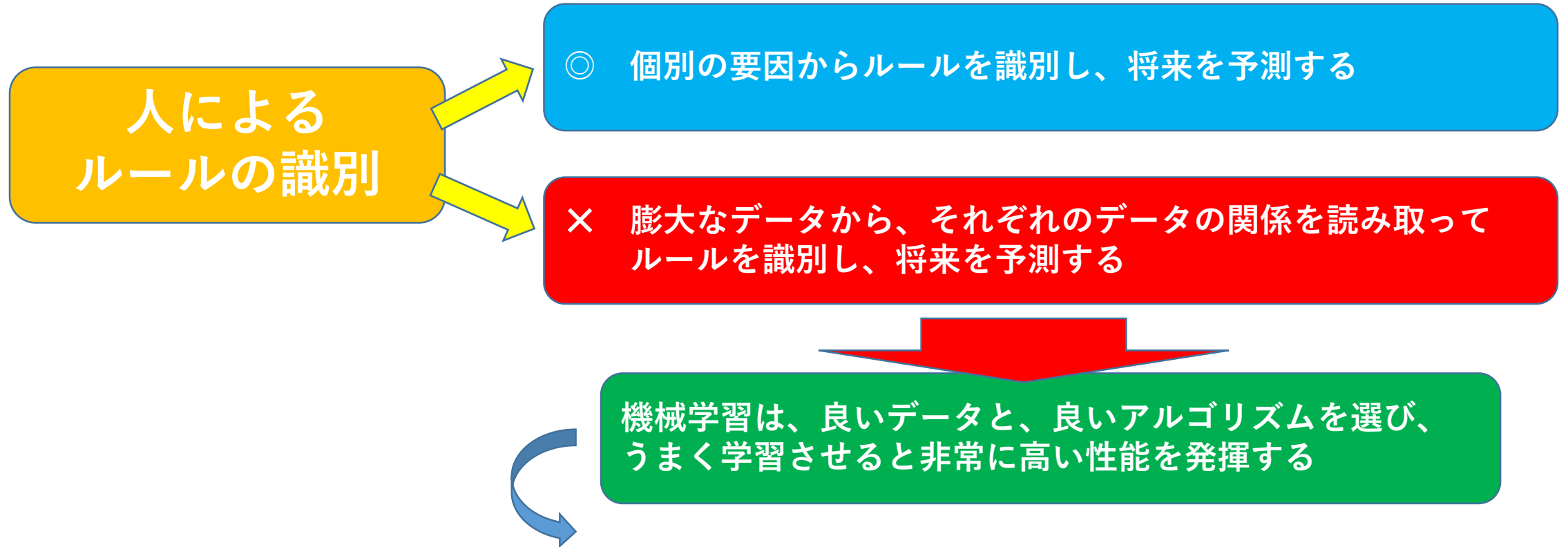


### ・ 機械の場合



機械学習自体には身体的なインターフェースはないため、別途ロボットなどが必要になる

## ◆機械学習の目的 ②高精度の識別、予測による効率化



DeepMind社の機械学習技術による、「データセンター内の冷却電力の効率化」。  
データセンター内外のさまざまな場所に取りつけたセンサーのデータを学習することで、設備の稼働状況や周辺の変化に応じた冷却設備の稼働コントロールを実現。  
結果、従来と比べて約40%の電力の削減が図られたと発表されている。

# ◆機械学習の3つの手法

## ①「教師あり学習」

入力に対してあらかじめ正解がわかっている場合に、正解を導くパターンやルールを学習する手法。

ここでの「教師」とは、正解データのこと。

## ②「教師なし学習」

正解のないデータから類似グループをまとめたり、重要な特徴を抽出したりする学習手法。

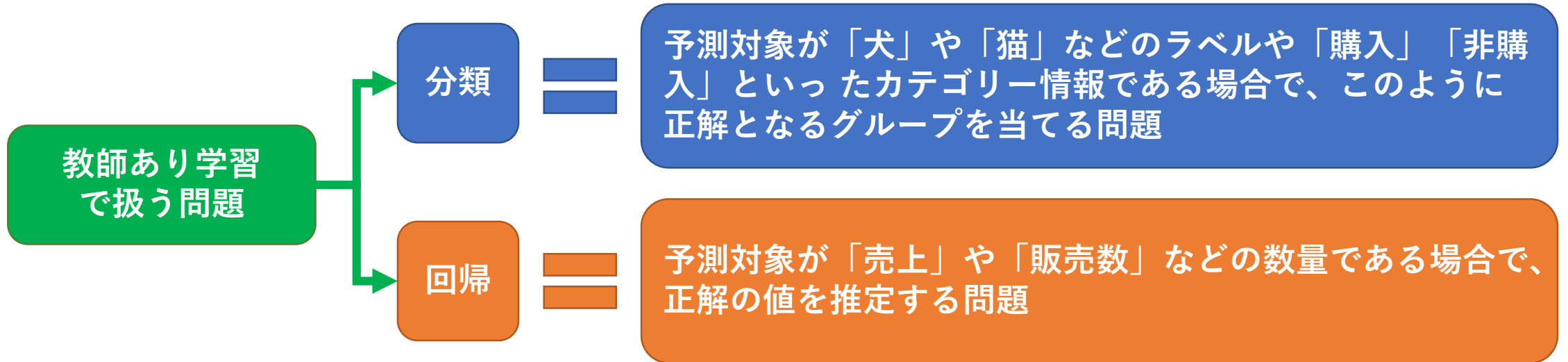
## ③「強化学習」

コンピューターが自ら試行錯誤しながら最適な戦略を学習する手法。

### ▶ 機械学習の3つの手法 図表13-1

手法	機能
教師あり学習	正解のラベルや数値がわかっているデータをもとに学習モデルを構築し、ラベルや数値が未知のデータに対して予測や識別を行う手法。正解データがない場合は人がデータに正解をつける必要がある。たとえば大量のペットの画像データがあったとして、それぞれに「犬」や「猫」といったラベルをつけて学習させ、学習モデルを作る。以後は、ラベルがついていないデータを学習モデルに与えると、犬の画像なら犬、猫の画像なら猫に自動的にラベルづけされる
教師なし学習	正解のないデータから、共通する特徴を持つグループを見つけたり、データを特徴づける情報を抽出したりする学習手法。たとえば、購買データから購買行動が似ているユーザーをグルーピングしたり、アンケートデータからユーザーの嗜好（ブランド好き、アウトドア好きなど）を抽出したりすることができる
強化学習	ゲームやギャンブルなど、結果がでるまでに時間がかかったり、多数の繰り返しが必要になったりするタスクに関して、実際に行動しながら最適な戦略を学習する手法。たとえばゲームなどで、コンピューター自身がどのように行動したら高得点が得られるかを試行錯誤しながら学んでいく

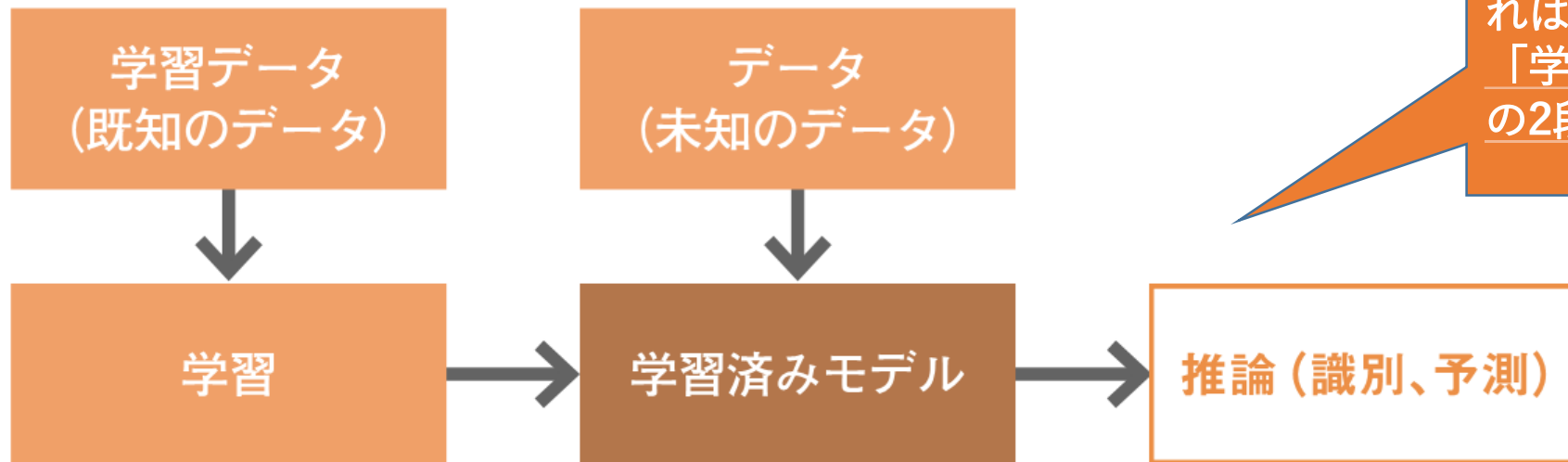
## ◆教師あり学習でできること





# ◆機械学習（教師あり学習）の基本的な流れ

## ▶ 機械学習(教師あり学習)の基本的な流れ 図表13-2



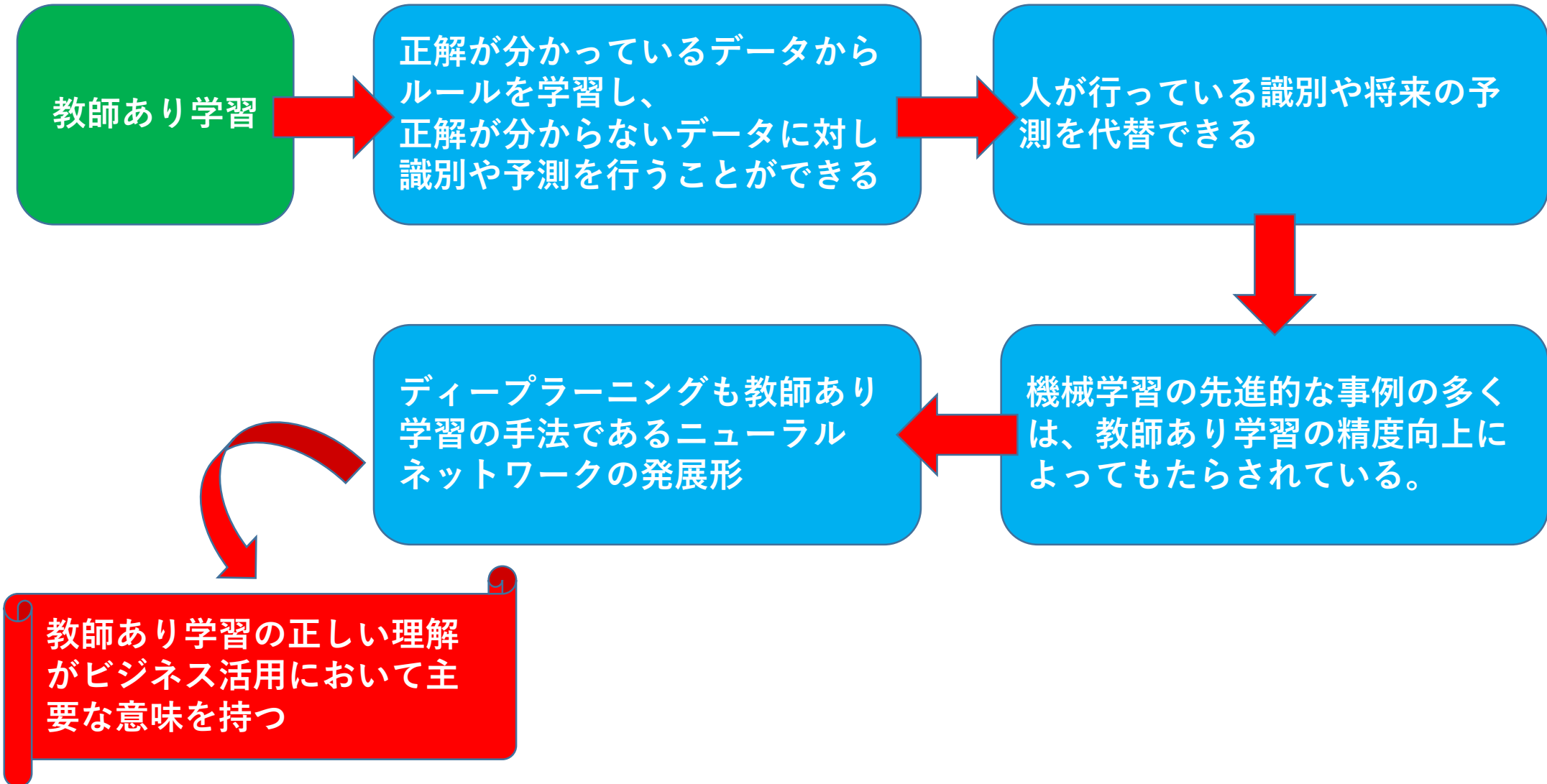
教師あり学習の一般的な流れは、大きく分けると「学習」と「推論」の2段階で構成されている。

「データから正解を導くパターンやルールを見つけ、入力に対して正解を出力するモデルを構築する」処理。

「構築したモデルを学習データではないデータに適用して、識別や予測を行う」処理。

適切な学習により構築されたモデルは、推論においてもよい精度を持つようになる。  
そのため、よい推論が行えるようなモデルを構築することが機械学習における主な目標になる。

# ◆教師あり学習の特徴



# ◆教師なし学習とは？

## 教師なし学習

識別や予測の対象となる正解（教師）がない学習の手法

類似データの  
グルーピング  
↓  
クラスタリング

クラスタリングとはデータ間の距離に基づいて類似データをグルーピングすること

「乗用車ブランドの好意度を尋ねたアンケート」があるとして、車に対するイメージ（高級車好き、性能重視、安価志向、安全性重視など）が近いユーザーをまとめるといったことができる。

各グループの嗜好に合った内容のダイレクトメールを送る

次元削減

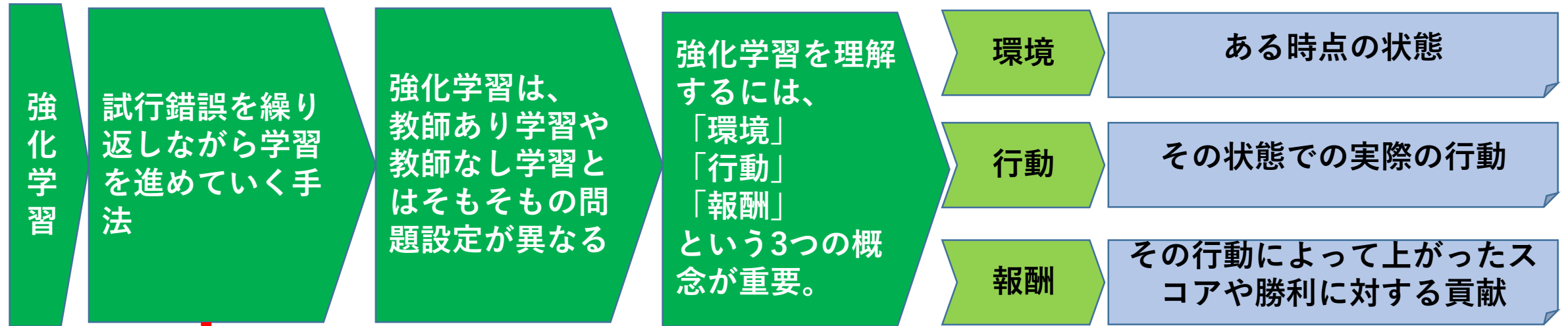
次元削減とはデータの次元（特徴量の数）を削減する手法

次元削減を行うことで、より本質的なデータを抽出することが可能となる

「主成分分析」と呼ばれる手法で、少ない情報で全体のばらつきを表現する主成分という特徴を抽出することができる

商品イメージに関するアンケートから、複数の評価項目をまとめて「安全・安心」「価格」「機能性」「話題性」に関するイメージのような少数の要素を抽出するといったことが可能

# ◆強化学習とは



主にゲームやギャンブルなど、結果が出るまでに時間がかかるタスクや多数の繰り返しが必要になるタスクに関して、コンピューターが実際に行動しながら最適な戦略を学習していく。

AlphaGoはまさに強化学習が力を発揮した事例。  
AlphaGo同士がコンピューター上で対局を繰り返しながら強くなっていき、人間を超える強さを獲得している。

# ◆機械学習のモデルとアルゴリズム

教師あり学習  
における  
モデルと  
アルゴリズム

教師あり学習は  
推定や予測  
をする問題に  
適用できる

推定や  
予測  
👉  
回帰

機械学習では  
識別や予測を  
行うために、  
モデルを学習  
アルゴリズム  
にもとづいて  
構築する。

モデル  
とは学習アルゴ  
リズムによって  
獲得された  
パターンやルー  
ルを数式的に表  
現したもの

ある商品の  
売上を  
気温と湿度  
から予測する  
場合の  
線形回帰モデ  
ルの例

## ▶ 線形回帰の例 図表14-1

$$\text{予測売上} = a \times \text{気温} + b \times \text{湿度} + \text{定数}$$

目的変数

重みパラメーター

説明変数

気温と湿度  
が0の場合の  
売上に相当

説明変数

識別や予測を  
行う対象とな  
る変数

目的変数

識別や予測を  
行うために必  
要となる変数

売上が気温と湿度に応じて変動する場合の線形回帰による売上予測の場合。重みパラメーターと定数部分を「最小二乗法」というアルゴリズムによって求める

# ◆モデル構築のフロー

顧客データ売上  
データ  
などの表形式の  
データ

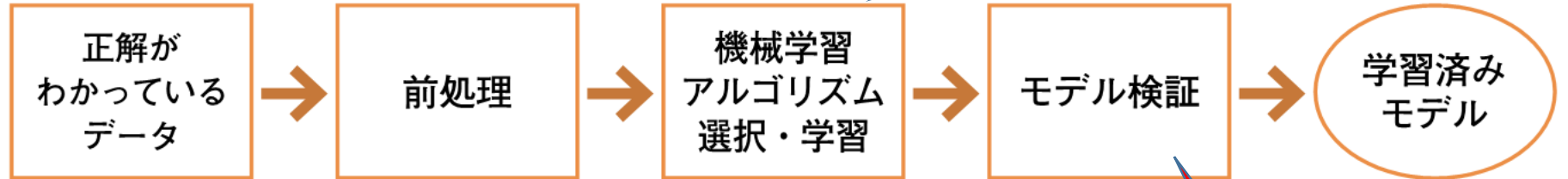
画像  
JPEG、GIF  
など

音声  
MP3、MP4

自然言語  
テキスト

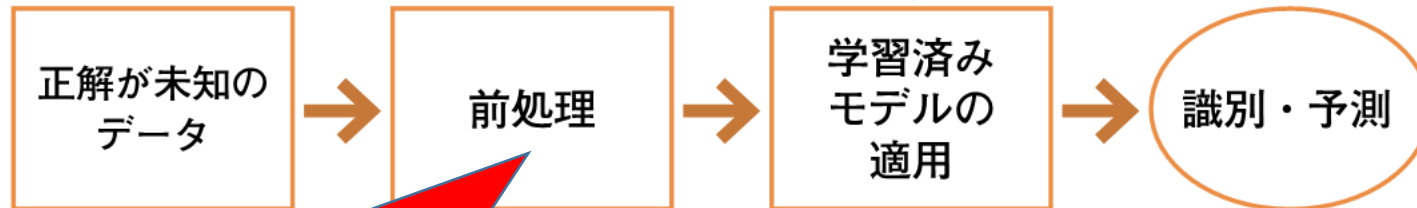
## ▶ 機械学習(教師あり学習)の全体像 図表14-2

### ・学習ステップのフロー



学習アルゴリズムによって  
数学的な手順にもとづいてパ  
ターンやルールを学習する

### ・推論ステップのフロー



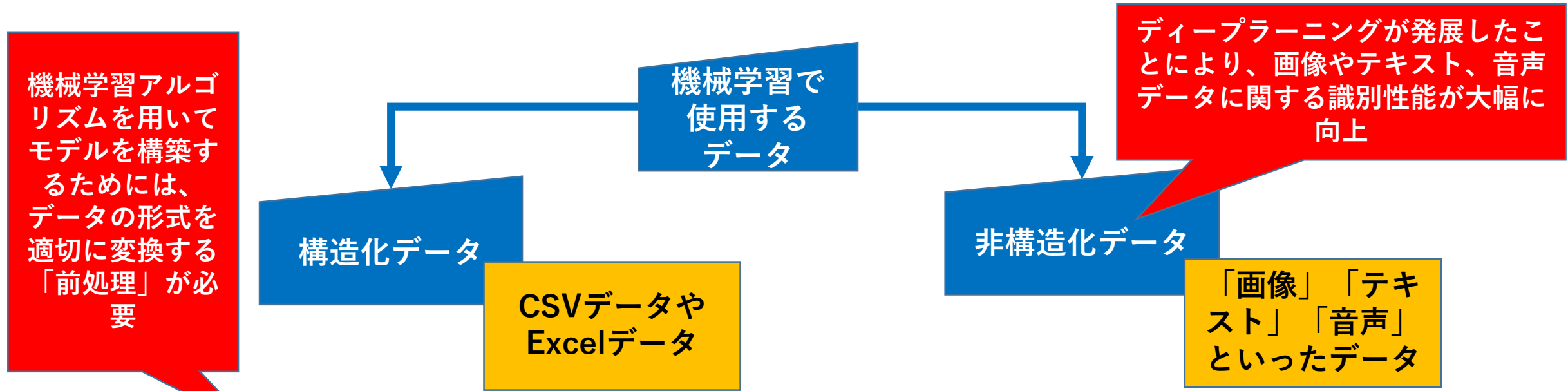
正解が未知デー  
タに対してどの  
程度の予測性能  
が期待できるか  
を確認する

機械学習アルゴリズムを適用可能な形式にデー  
タを加工する



各種データを数値形式のデータに変換する  
データに欠損や異常がある場合は補完する

# ◆データの種類と前処理・構造化データと非構造化データ



## ▶ 構造化データ・非構造化データの例 図表15-1

### 構造化データ

- ・ 各種業務システム内のデータ（受注、発注、在庫、人事、POS など）
- ・ 政府や調査会社の統計データ

### 非構造化データ

- ・ 画像データ（商品画像、SNS の投稿画像など）
- ・ 動画データ（監視カメラ映像、テレビ番組など）
- ・ テキストデータ（議事録、SNS 投稿テキストなど）
- ・ 音声データ（コールセンター会話録、会議録音データなど）

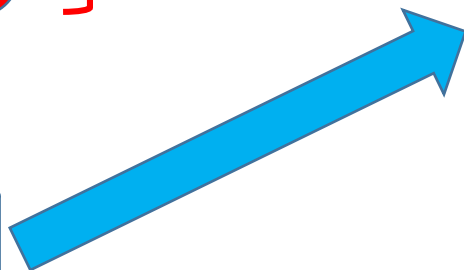
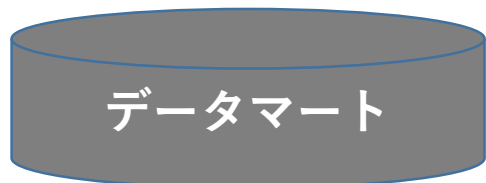
# ◆構造化データに対する機械学習

顧客ごとのある商品の  
購入確率を予測したい



顧客ごとの商品購入経験の  
有無と、商品購入に影響を  
与えると考えられる情報を  
行列形式に整理する

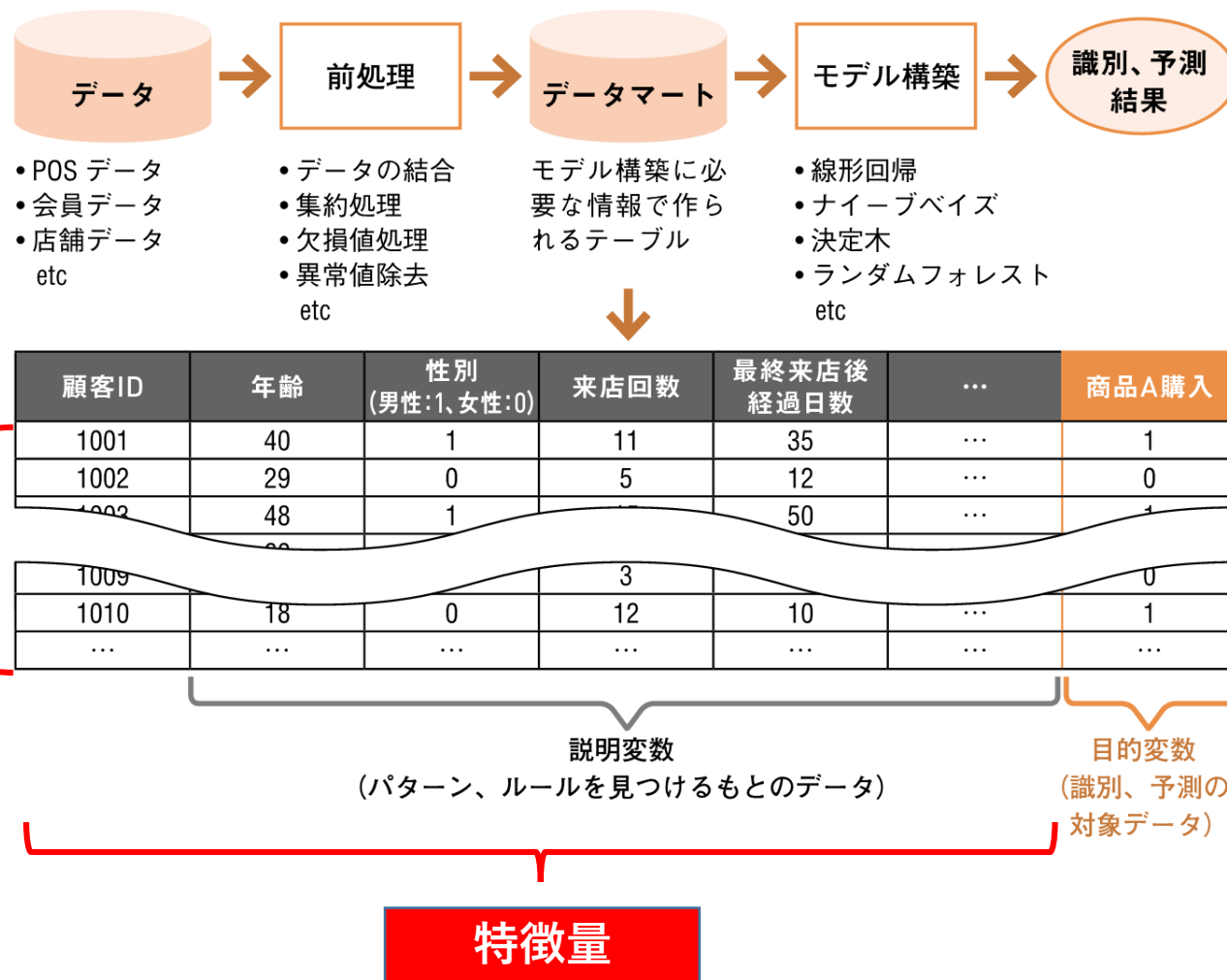
データが欠落していたり、  
異常な値のデータが含まれ  
ている場合、適切にデータ  
を補完したり、異常データ  
を除去する



前  
処  
理

事  
例

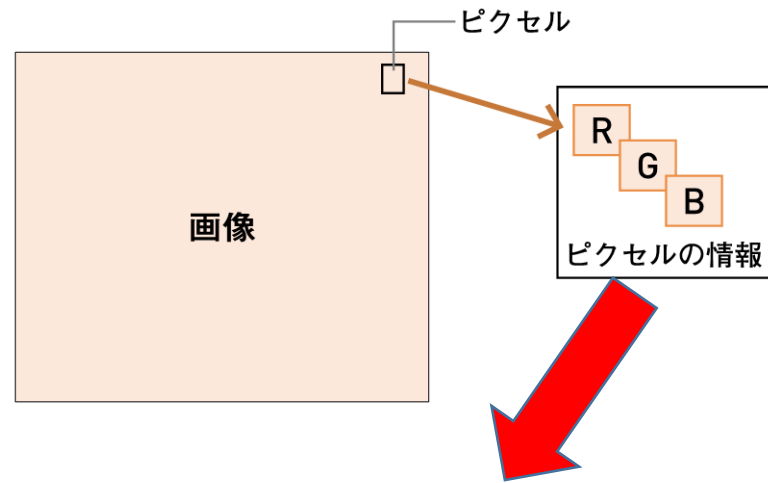
▶ 構造化データにおけるモデル構築のイメージ 図表15-2



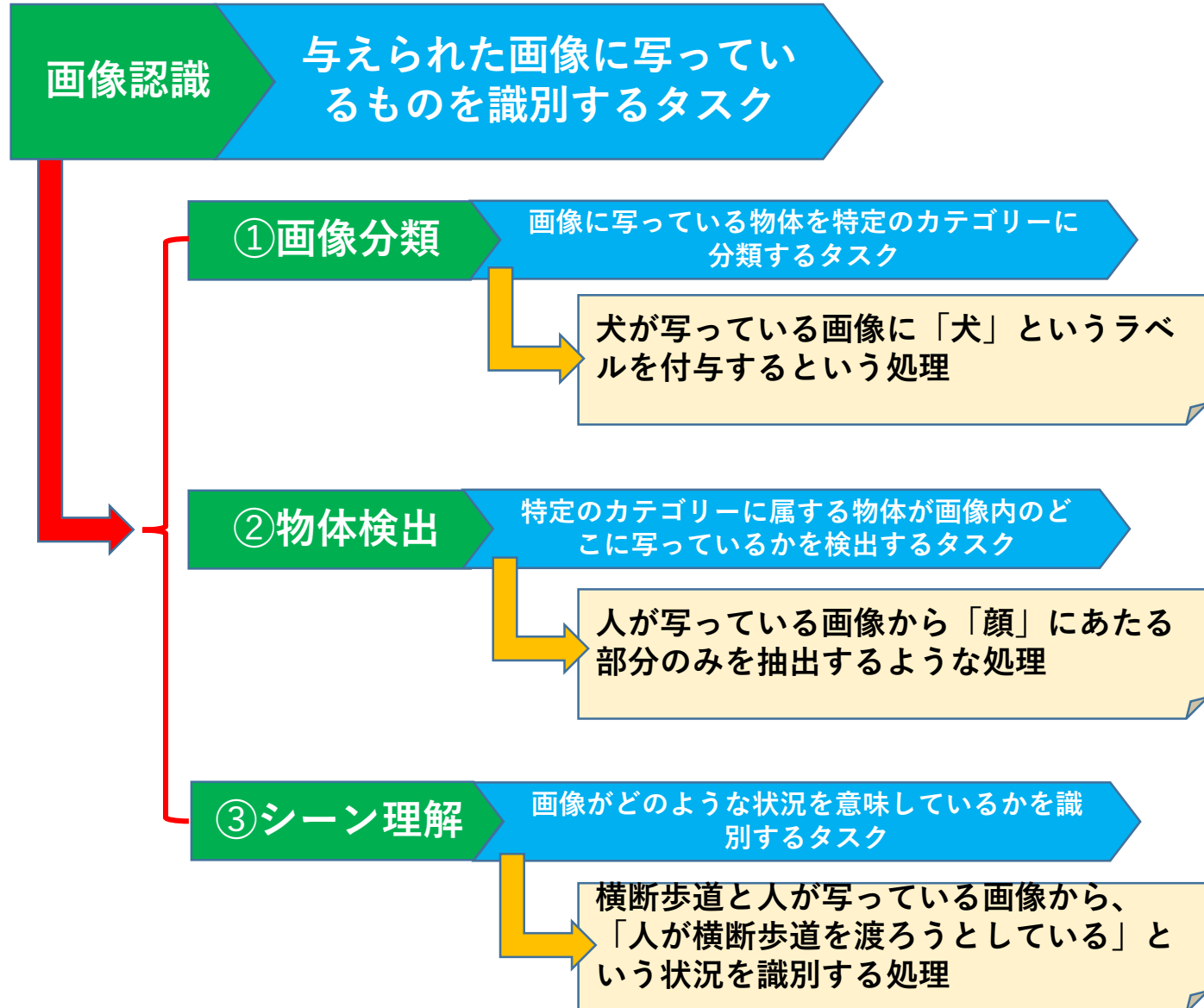


# ◆画像データに対する機械学習と画像認識のタスクの種別

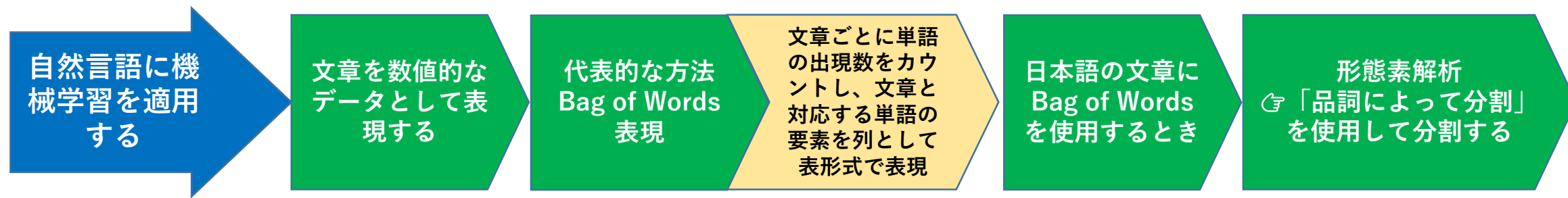
▶ 画像データの構造 図表15-3



各ピクセルの情報を入力データとして、機械学習のアルゴリズムに適用



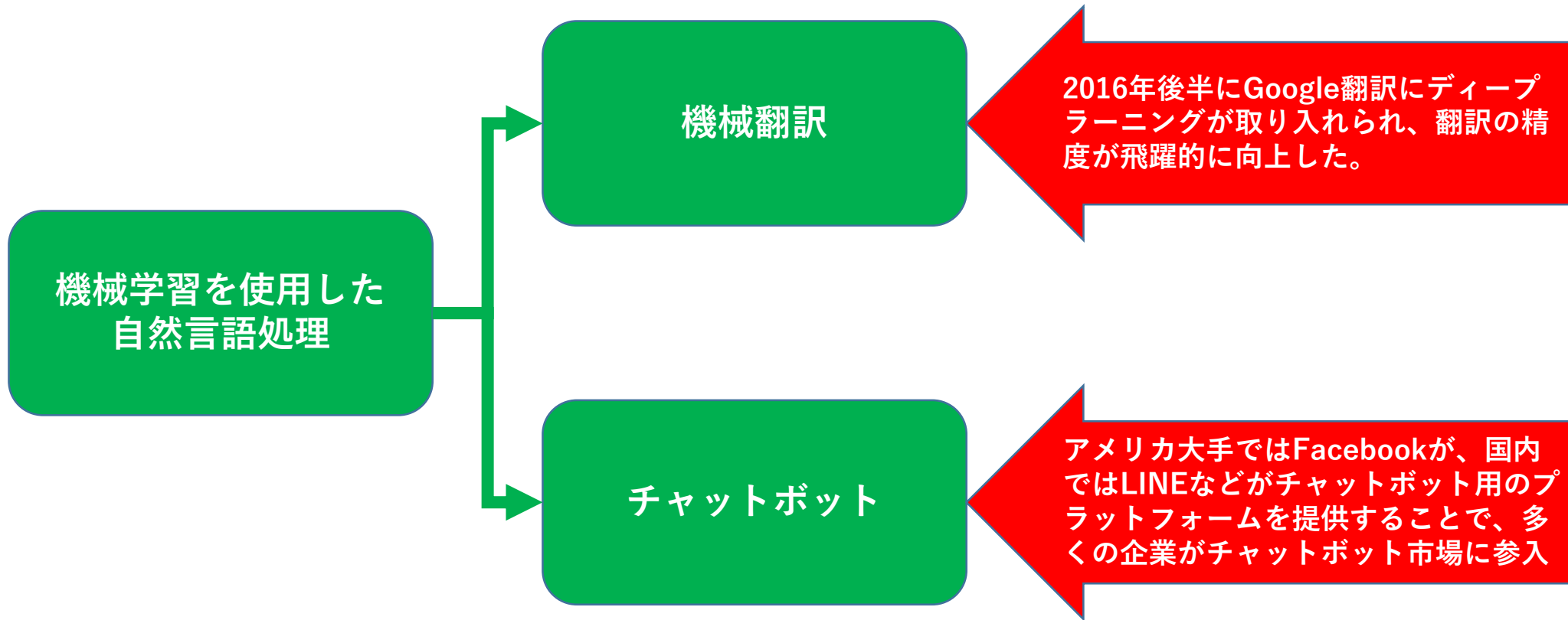
# ◆自然言語データに対する機械学習



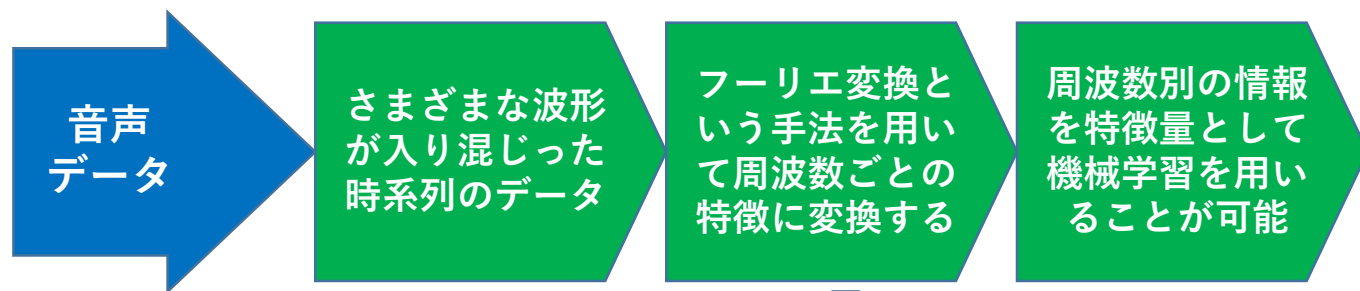
## ▶ Bag of Words表現 図表15-4



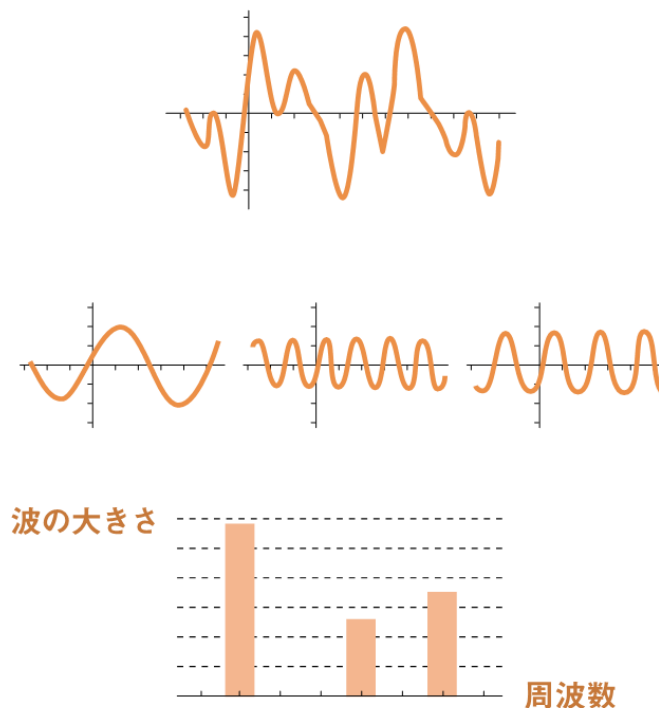
## ◆自然言語処理の活用事例



# ◆音声データに対する機械学習



▶ フーリエ変換のイメージ 図表15-5



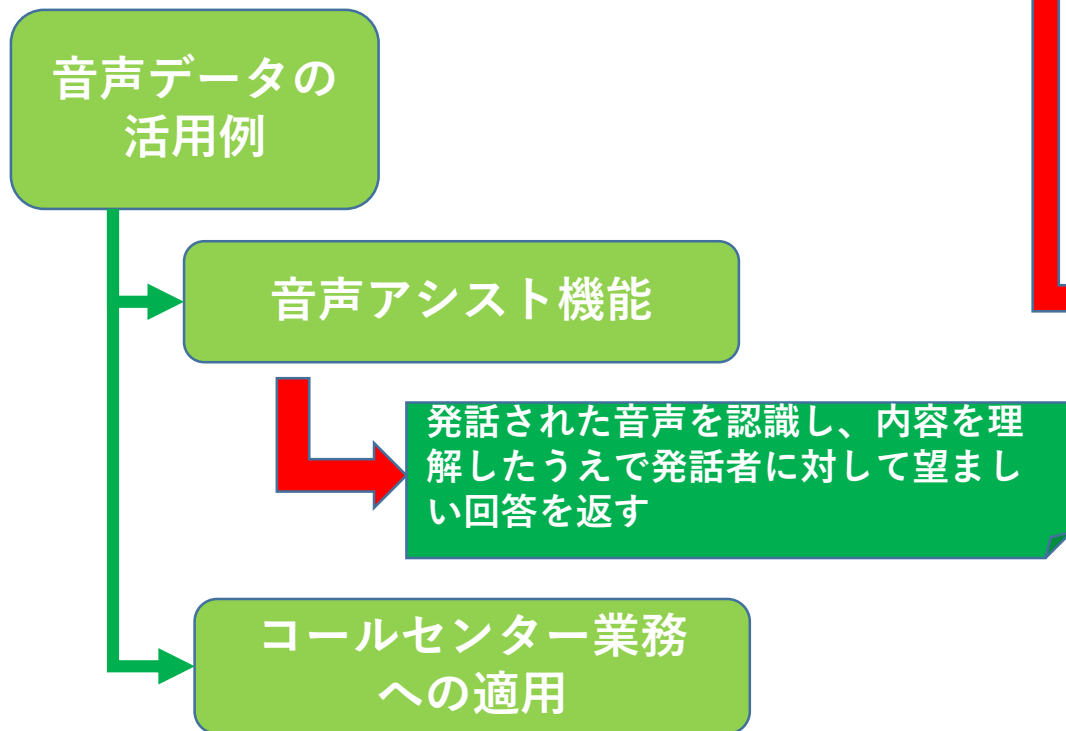
複雑な波形の元データ



波の大きさと周波数に応じて分解する



波の大きさと周波数で表す



# ◆学習アルゴリズムの選択

これまでにさまざまな学習アルゴリズムが考案されている

線形回帰

ランダムフォレスト

ナイーブベイズ

ニューラルネットワーク

決定木

サポートベクターマシン

それぞれ適用可能な状況やパフォーマンスを発揮するための前提条件が異なる

目的や状況に応じて適切な手法を選ばなくてはならない

## ▶ 教師あり学習の主なアルゴリズム 図表16-1

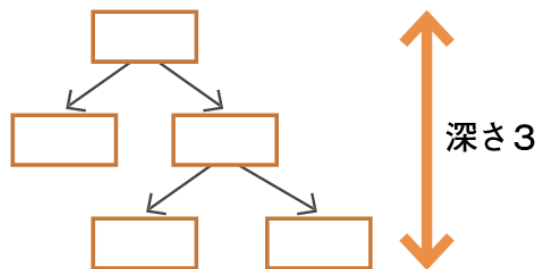
手法	説明
線形回帰 (Linear Regression)	目的変数と説明変数に線形関係を仮定してモデルを構築する手法。予測値に対する各特徴量の寄与度がわかるため、構造の解釈がしやすい
ナイーブベイズ (Naive Bayes)	特徴量間に独立性を仮定して、各カテゴリに所属する確率を推定する手法。主に文書分類などに利用される
決定木 (Decision Tree)	木のように段階的にルールを適用しながらデータを分類していく手法。ルールを辿っていくプロセスが可視化されるため、構造の解釈がしやすい
ランダムフォレスト (Random Forest)	事例と特徴量をランダムにサンプリングしたデータをもとに決定木を多数作り、多数決的に回帰・分類を行う手法
ニューラルネットワーク (Neural Network)	人間の脳の活動を模した学習の手法。ニューロンと呼ばれる層をつなぎ合わせてネットワークを構築する
サポートベクターマシン (Support Vector Machine)	カーネル関数といわれる高次元空間への写像関数を用いて分類を行う手法。線形分離不可能なデータも分類できる一方で、計算コストが高い

# ◆モデルの複雑さを調整するハイパーパラメータ

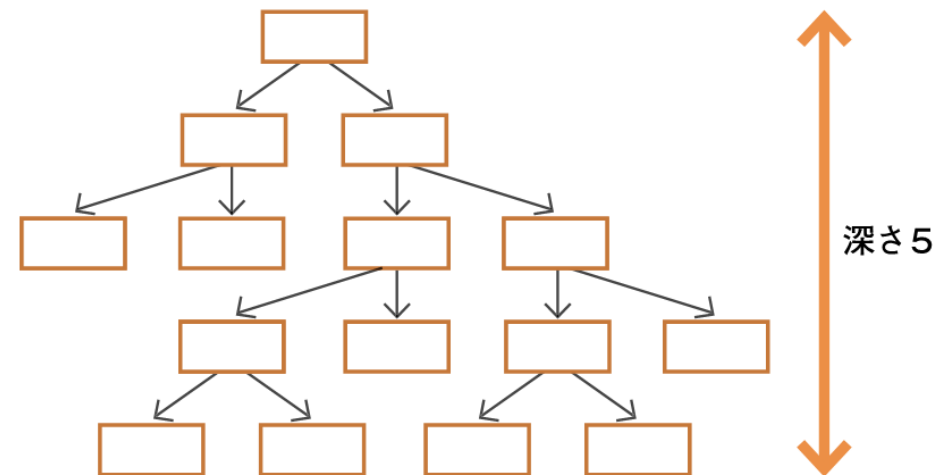


## ▶ 決定木におけるハイパーパラメーターの例 図表16-2

単純なモデル



複雑なモデル



木の深さがハイパーパラメーターに該当し、深いほうが複雑なモデルを意味する

# ◆ニューラルネットワークとディープラーニング

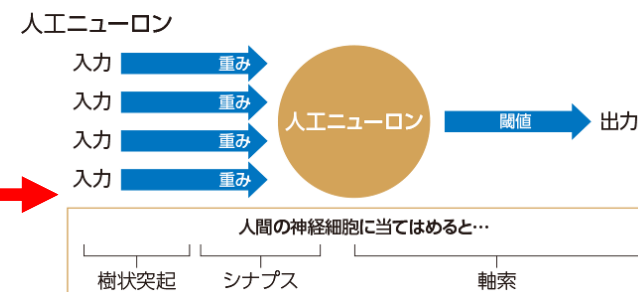
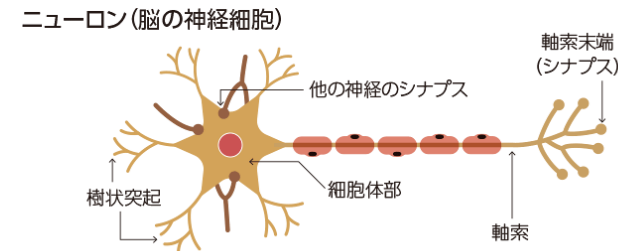
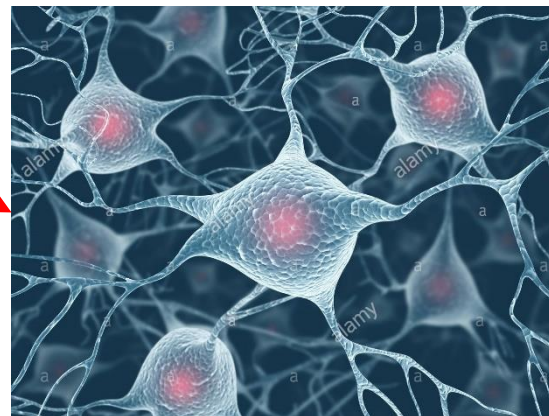
人の脳はニューロンと呼ばれる多数の神経細胞で構成されており、ニューロンが複雑に結合することによって情報処理が行われていると考えられている

この情報処理ネットワークを模したモデルがニューラルネットワークである

ニューラルネットワークはニューロンに該当する多数のユニットから構成される

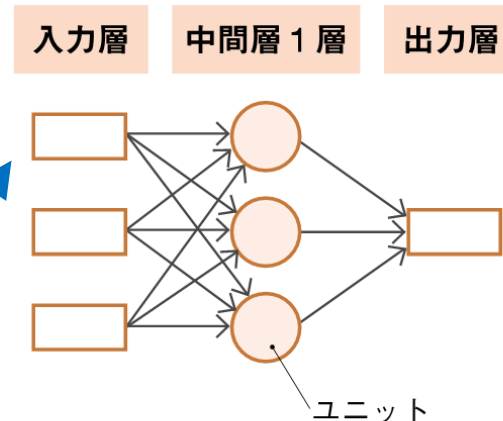
各ユニットは入力に対して、ユニット間の重みを加味した非線形な演算を行った結果を出力するという構造を持つ

このときニューラルネットワークの層を多数重ねたものを「ディープニューラルネットワーク」と呼び、ディープニューラルネットワークを用いた機械学習の手法を「ディープラーニング」と呼びます

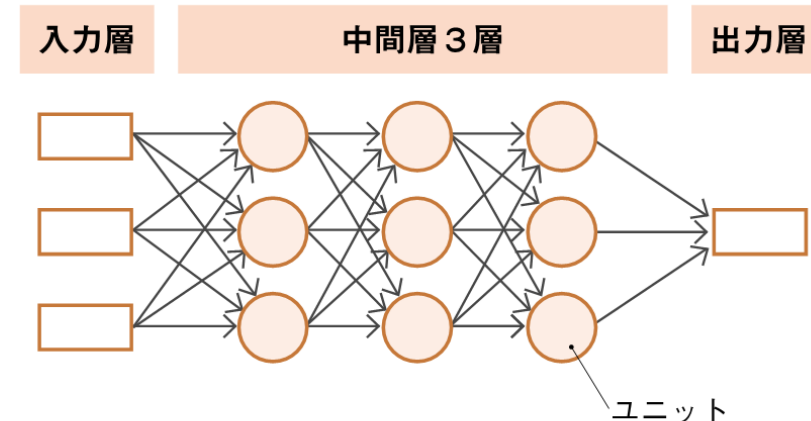


## ▶ ニューラルネットワークとディープラーニングのイメージ 図表17-1

ニューラルネットワーク



ディープラーニング (=多層ニューラルネットワーク)



ニューラルネットワークの中間層を多層構造にしたモデルをディープラーニングと呼ぶ



# ◆ディープラーニングの学習について

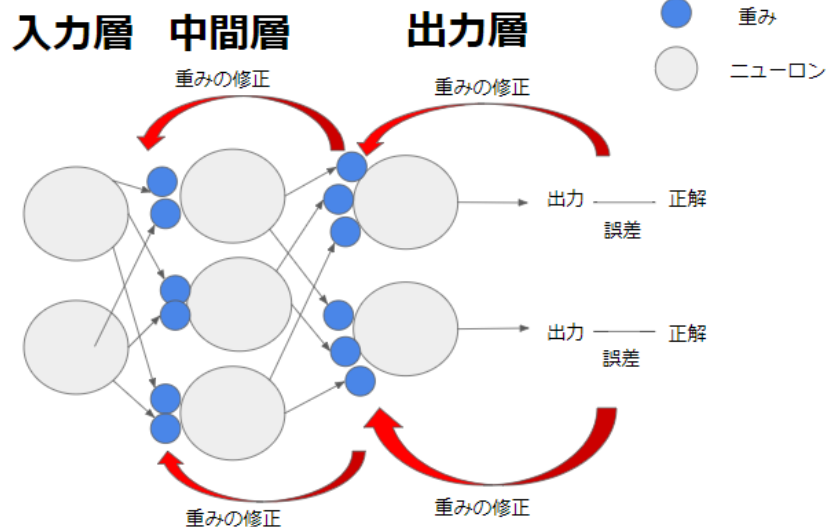
ディープ  
ラーニング  
の学習

「誤差逆伝搬法」  
(バックプロパ  
ゲーション)  
アルゴリズム  
を用いる

正解との誤差が段々小さ  
くなるようにユニット間  
の重みを調整していくこ  
とで実現される

層が増えると誤差逆伝  
搬法が機能しなくなる  
問題が起きるため、適  
切に学習できないとい  
うことがあった

アルゴリズムの改善と計算資  
源の向上により、近年では多  
くの問題に対してほかのモデ  
ルが追従できない精度を達成  
するようになってきた



## ▶ さまざまなディープラーニング関連のモデル 図表17-2

手法	説明
ディープニューラルネットワーク (Deep Neural Network : DNN)	ニューラルネットワークの層を多層にしたモデル。一般的には中間層が2層以上の深い構造を持つ
オートエンコーダ (AutoEncoder)	ニューラルネットワークの出力が入力と一致するように構成したモデルであり、主に次元削減のために利用される
畳み込みニューラルネットワーク (Convolutional Neural Network : CNN)	主に画像解析に用いられるモデルで、局所的な領域の情報を畳み込む処理を行う畳み込み層と特徴を集約するプーリング層を多層に組み合わせて構成される
リカレントニューラルネットワーク (Recurrent Neural Network : RNN)	主にテキスト解析や時系列解析に用いられるモデルで、隠れ層の値を再び次の層の入力として用いるという再帰的な構造を持つ