

現場で使える機械学習・データ分析  
基礎講座 (DAY1)  
SkillUP AI

# 本講座の方針

---

- ・ 本講座は、動画講義、対面講義、通し課題で構成されます
- ・ 動画講義では、機械学習の理論に関する部分を説明します
- ・ 対面講義では、動画講義の復習、Notebook演習、グループワークを行います
- ・ 通し課題では、DAY1-DAY4までを通して、1つの課題に取り組んで頂きます
- ・ 対面講義は、動画講義の復習から行います
- ・ カリキュラム自体は対面講義と動画講義は同じです
- ・ 内容が非常に高度ですので、**毎回必ず動画講義で予習してきてください**

## 本講座で必要になるPC環境

---

- ・ 本講座の対面講義では、 Jupyter Notebook(言語はPython)を用います
- ・ Jupyter Notebookは、 Anacondaというソフトをインストールすることで利用できるようになります
- ・ 対面講義では、 Anaconda3-5.1.0がインストールされていることを前提に講義を進めますので、 事前にインストールをしておいてください
  - ・ 最新のAnacondaの場合、 TensorFlowがPython3.7に公式対応していないため、 インストールがうまくいかない場合があります
- ・ その他必要なライブラリは、 0.preparation.ipynbを参考に事前にインストールしておいてください

## 参考図書（全日程分、およそ日程の進行順に記載しています）

---

- ① 『ITエンジニアのための機械学習理論入門』(中井悦司、技術評論社)
- ② 『はじめてのパターン認識』(平井有三、森北出版)
- ③ 『わかりやすいパターン認識』  
(石井健一郎、前田英作、上田修功、オーム社)
- ④ 『機械学習 -データを読み解くアルゴリズムの技法-』  
(Peter Flach 著、竹村 彰通ら訳、朝倉書店)
- ⑤ 『深層学習 Deep Learning』(人工知能学会、近代科学社)
- ⑥ 『これならわかる深層学習入門』(瀧 雅人、講談社)
- ⑦ 『深層学習』  
(Ian Goodfellow、Yoshua Bengioら著、KADOKAWA)
- ⑧ 『続・わかりやすいパターン認識—教師なし学習入門—』  
(石井健一郎、上田修功、オーム社)

# 本講座の構成

DAY1	DAY2
<ul style="list-style-type: none"><li>• 機械学習概論<ul style="list-style-type: none"><li>• 人工知能とは</li><li>• 機械学習とは</li><li>• 機械学習アルゴリズムの実装とワークフロー</li><li>• 機械学習アルゴリズム概観</li></ul></li><li>• 教師あり学習の基礎<ul style="list-style-type: none"><li>• 線形回帰</li><li>• ロジスティック回帰</li><li>• 多変量モデルへの拡張</li></ul></li><li>• モデルの評価指標<ul style="list-style-type: none"><li>• 回帰問題 (MAE/MSE/RMSE)</li><li>• 分類問題 (精度/適合率/再現率/F1-score)</li></ul></li></ul>	<ul style="list-style-type: none"><li>• モデルの検証・正則化<ul style="list-style-type: none"><li>• 訓練誤差と汎化誤差</li><li>• 過学習</li><li>• 正則化 (L2/L1)</li><li>• ホールドアウト法・交差検証法</li></ul></li><li>• 前処理<ul style="list-style-type: none"><li>• 正規化 / 標準化</li><li>• 無相関化 / 白色化</li></ul></li><li>• 教師あり学習の発展的トピック<ul style="list-style-type: none"><li>• サポートベクターマシン</li></ul></li></ul>

# 本講座の構成

DAY3	DAY4
<ul style="list-style-type: none"><li>前処理<ul style="list-style-type: none"><li>特徴選択</li></ul></li><li>教師あり学習の発展的トピック<ul style="list-style-type: none"><li>木モデル (決定木・ランダムフォレスト)</li><li>ニューラルネットワーク</li></ul></li></ul>	<ul style="list-style-type: none"><li>教師あり学習の発展的トピック<ul style="list-style-type: none"><li>深層学習</li><li>k-最近傍法</li></ul></li><li>教師なし学習<ul style="list-style-type: none"><li>クラスタリング</li><li>特徴抽出・次元削減</li></ul></li><li>モデルの改善<ul style="list-style-type: none"><li>ハイパーパラメータ最適化</li></ul></li></ul>

## DAY1の目標とDAY4までの流れ

---

- DAY1では、講義内容をもとに、精度は気にせず、とにかく教師あり機械学習モデルを構築できるところまでを目指します
- DAY2以降で、精度を向上させる方法や様々なアルゴリズムを紹介し、教師あり機械学習モデル構築時の手札を増やしていきます
- DAY4で通し課題の最終発表の時間を設けているので、最初のモデルから何を変えてどの程度精度が向上したかを発表して頂きます
- DAY4では、教師なし学習や深層学習などについても触れる予定です

# DAY1の目次

---

- 機械学習概論（15分）
  - 人工知能とは（2分）
  - 機械学習とは（3分）
  - 機械学習アルゴリズムの実装とワークフロー（5分）
  - 機械学習アルゴリズム概観（5分）
- 教師あり学習の基礎（2時間50分）
  - 線形回帰（60分）&グループワーク1（30分）
  - ロジスティック回帰（60分）
  - 多変量モデルへの拡張（20分）
- モデルの評価指標（10分）
  - 回帰問題（5分）  
(MAE/MSE/RMSE)
  - 分類問題（5分）  
(精度/適合率/再現率/F1-score)
- グループワーク2（30分）
- 通し課題に関する説明と質疑（15分）

# 機械学習概論

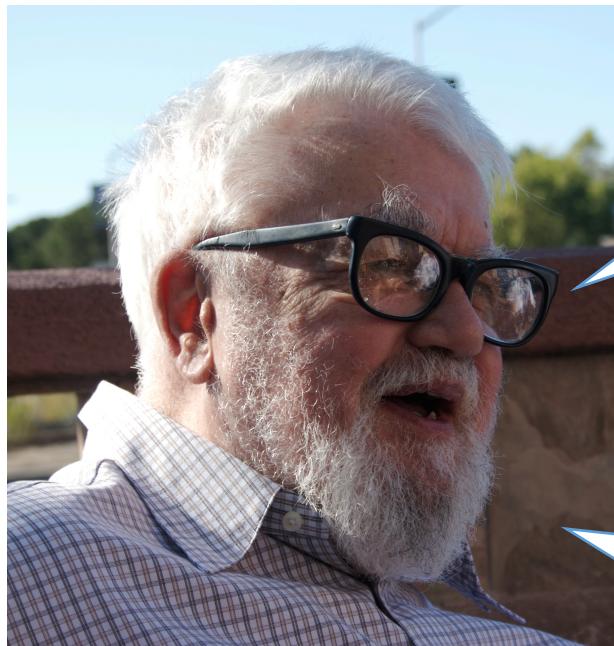
1. 人工知能とは
2. 機械学習とは
3. 機械学習アルゴリズムの実装とワークフロー
4. 機械学習アルゴリズム概観

# 人工知能とは？ 機械学習とは？ 深層学習とは？

---

- ・ 人工知能とは？SFの世界？
- ・ 機械学習とは？機械が学習するってどういうこと？
- ・ 深層学習とは？何が深いの？

# 人工知能（Artificial Intelligence）の定義



Q. 人工知能とは何でしょうか？

知的な機械、特に、  
**知的なコンピュータプログラムを作る科学と技術です<sup>[1]</sup>**

Q. 知能とは何でしょうか？

**実際の目標を達成する能力の計算的な部分です**  
人間、動物、そして機械には、種類や水準が  
さまざまな知能があります<sup>[1]</sup>

ジョン・マッカーシー\*

(John McCarthy, 1927-2011)

\*：人工知能研究の第一人者。「Artificial Intelligence」という用語を初めて公の場で使用した人物であり、プログラミング言語LISPの開発者でもある

[1] 人工知能学会：人工知能のFAQ、<https://www.ai-gakkai.or.jp/whatsai/AIfaq.html> (2018.09.12アクセス)

「人工知能 = 機械学習 = 深層学習」ではないので要注意！！

- ・機械学習→人工知能を実現するための技術領域のひとつ
- ・深層学習（ディープラーニング）→機械学習の方法論のひとつ

## 人工知能 (AI)

機械学習

深層学習

SVM

線形回帰

クラスタ分析

決定木

etc…

エキスパートシステム

探索アルゴリズム

etc…

# 機械学習（Machine Learning）の定義



明示的にプログラムしなくても学習する能力を  
コンピュータに与える研究分野<sup>[2]</sup>のこと

具体的には…

コンピュータプログラムが、あるタスクにおいて、  
用意されたデータを使い、性能の評価値を向上させること

アーサー・リー・サミュエル\*

(Arthur Lee Samuel, 1901-1990)

\* : 世界初の学習型プログラム「Samuel Checkers-playing Program」を開発した  
計算機科学者。ハッシュテーブルの考案者でもある。

[2] Wikipedia : 機械学習、<https://ja.wikipedia.org/wiki/機械学習> (2018.09.12アクセス)

# タスクとは

---

- ・機械学習分野では、機械学習モデルが対象とする問題設定のことをタスクと呼ぶ
  - ・例、株価予測、画像認識、画像物体検出、文章分類、機械翻訳

# 機械学習の全体像

## 学習フェーズ



過去のデータ



データ・タスクに  
対応した判断ルール

## 運用フェーズ



新たなデータ



データ・タスクに  
対応した判断ルール

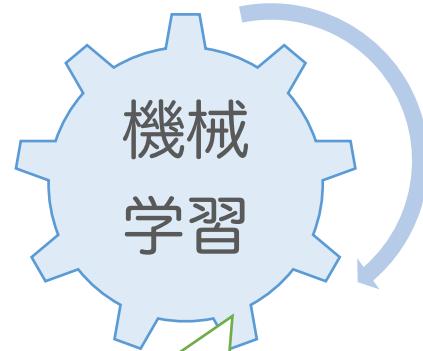


予測結果  
 = 犬

# 機械学習の全体像

## 学習フェーズ

過去のデータ



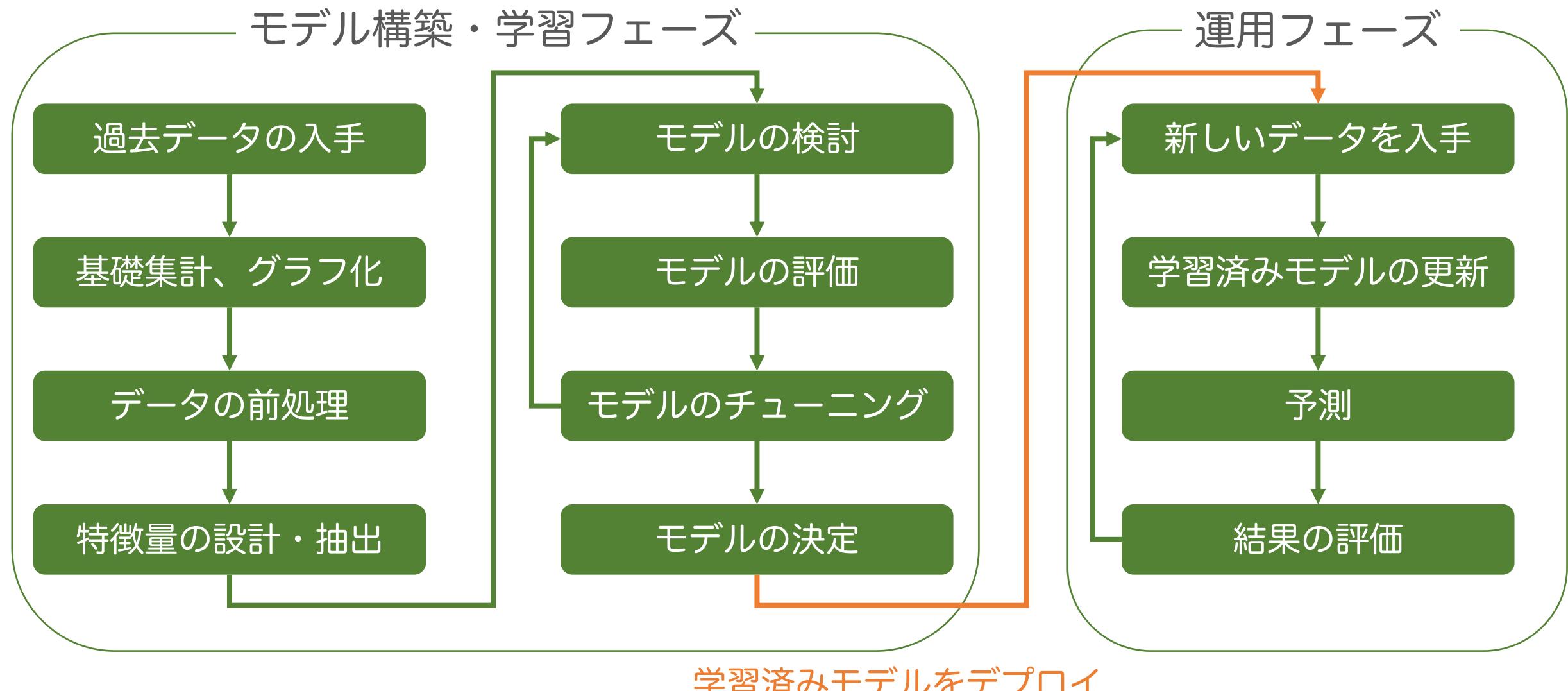
データ・タスクに  
対応した判断ルール

データの理解がないと  
意味のあるデータが選別不可

性能向上には  
理論の理解が不可欠

ビジネスの理解がないと  
判断ルールを活かせない

# 機械学習システムの開発・運用フロー



# 機械学習システムの開発・運用フロー | ビジネス観点で重要なこと

## モデル構築・学習フェーズ

過去データの入手

基礎集計、グラフ化

データの前処理

特徴量の設計・抽出

### 課題定義

- 機械学習で何を解くか？
- どれほどのビジネスバリューをもたらすか？
- その課題は本当に機械学習で解く必要がある課題か？

# 機械学習システムの開発・運用フロー | ビジネス観点で重要なこと

## モデル構築・学習フェーズ

過去データの入手

基礎集計、グラフ化

データの前処理

特徴量の設計・抽出

モデルの検討

モデルの評価

モデルのチューニング

モデルの決定

ROI観点でどこまで  
チューニングすべきか

- わずかの改善に多大な工数を割いてはならない
- 精度1%の改善でも50%→51%より90%→91%のほうがより工数がかかる

# 機械学習システムを実装する際に必要な知識とスキル

## 機械学習アルゴリズム

教師なし学習のアルゴリズム  
教師あり学習のアルゴリズム  
強化学習のアルゴリズム

## 基礎知識

数学記号  
微分積分  
線形代数  
確率  
統計学  
最適化手法

## 基本的なITスキル

Python  
Linux  
Bash  
SSH  
MySQL  
Json  
encoding  
Git

# RとPython、どっちがいいの？

	R	Python
習得の容易さ	◎ データ分析に特化しているため、記述が簡潔	○ 汎用言語のため、プログラミングの基本を理解する必要がある
応用発展性	○	◎ データ分析以外のライブラリも充実
開発環境	○ Rstudioは使いやすい	◎ Jupyter Notebookはもっと使いやすい
情報量	◎ 日本語の書籍やサンプルコードが多い	○ Stack OverFlowを見れば、大概解決する。(英語)
統計関連のライブラリ	○ 書籍やサンプルコードが多い	○ ライブラリはあるが、使用例や情報が少ない
機械学習関連のライブラリ	○	○
日本語対応	○ 不便を感じることはない	△ グラフ化など、たまに不便を感じる
レポーティング	○ knitrでレポートは作成できるが、使いづらい。NotebookでRカーネルをインストールして使う方がいいかも。	○ Notebook自体がレポートになる。 柔軟にいろいろできる

# 主要ツールやライブラリ概観

---

プログラミング言語



機械学習ライブラリ



ベクトル・行列計算に  
特化したライブラリ



計算可視化ライブラリ



データ処理ライブラリ



科学計算ライブラリ



より使いこなせるようになる方は『機械学習・ディープラーニングのためのPython基礎講座』へ！

# 機械学習アルゴリズムの分類

## 教師あり学習 (Supervised Learning)

- 入力とそれに対する正しい出力 (=教師信号) が学習用データとしてモデルに与えられる
- モデルは出力を教師信号に近づけるように学習

## 教師なし学習 (Unsupervised Learning)

- 入力のみが学習用データとして、モデルに与えられる
- モデルは入力の関係性や構造をうまく表現するように学習

## 強化学習 (Reinforcement Learning)

- 入力のみが学習用データとして与えられる
- 出力は与えられないが、出力がどれだけ良いかはわかる
- モデルは出力の良さを最大化するように学習

※モデル：データを説明する判断ルール（関数や確率分布）のこと

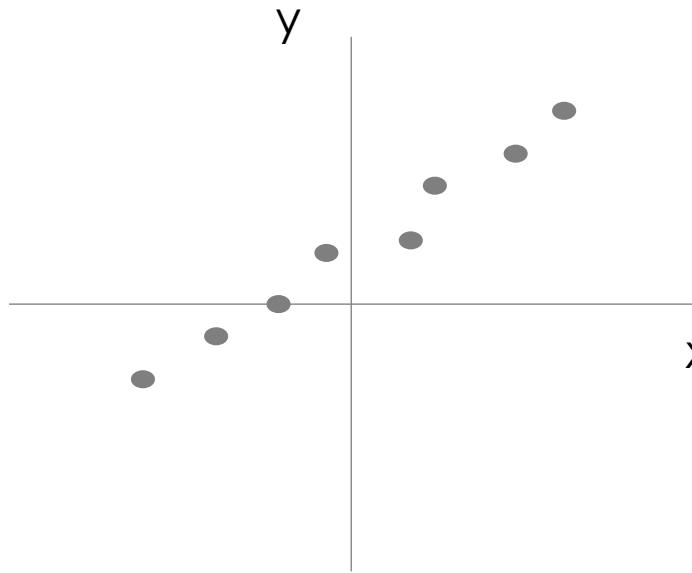
# 教師あり学習の代表的なアルゴリズム

↓ 詳しくはDAY2にて！

アルゴリズム	目的	線形/非線形	正規化、標準化
線形回帰	回帰	線形	必要
ロジスティック回帰	回帰、分類	線形	必要
サポートベクターマシン	回帰、分類	非線形	必要
k近傍法	回帰、分類	非線形	必要
決定木	回帰、分類	非線形	不要
ランダムフォレスト	回帰、分類	非線形	不要
ブースティング	回帰、分類	非線形	不要
ナイーブベイズ	分類	非線形	不要
状態空間モデル、階層ベイズモデル	回帰、分類	非線形	必要
ニューラルネットワーク(深層学習)	回帰、分類	非線形	必要

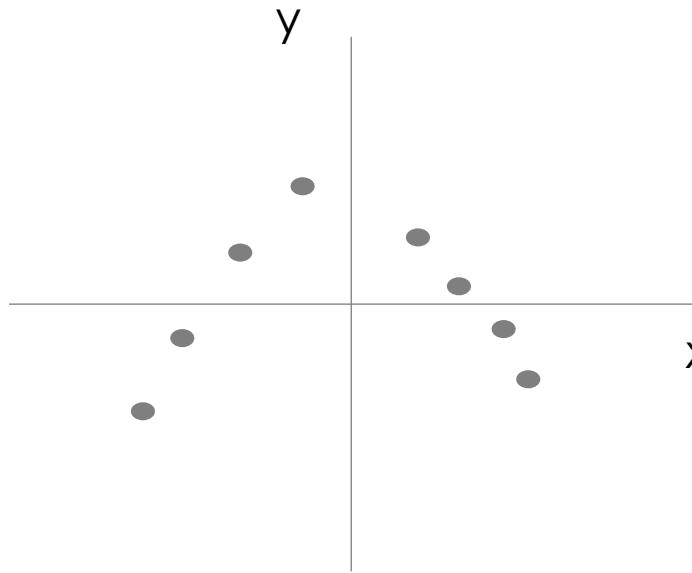
これは線形?非線形?

---



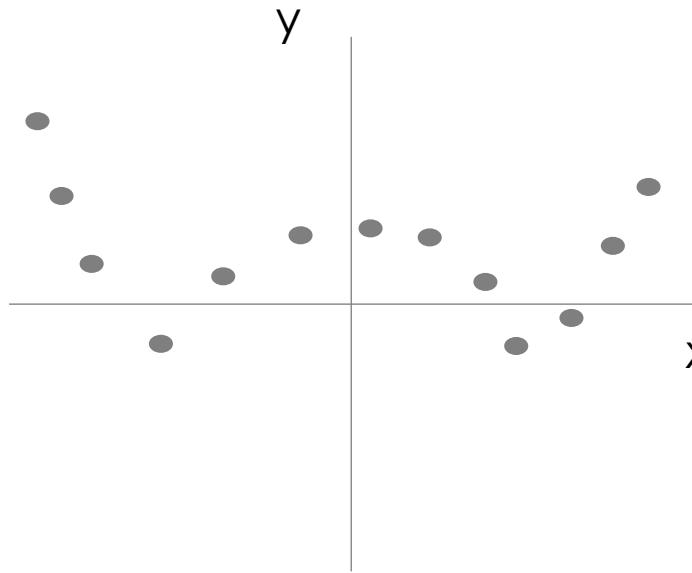
これは線形?非線形?

---



これは線形?非線形?

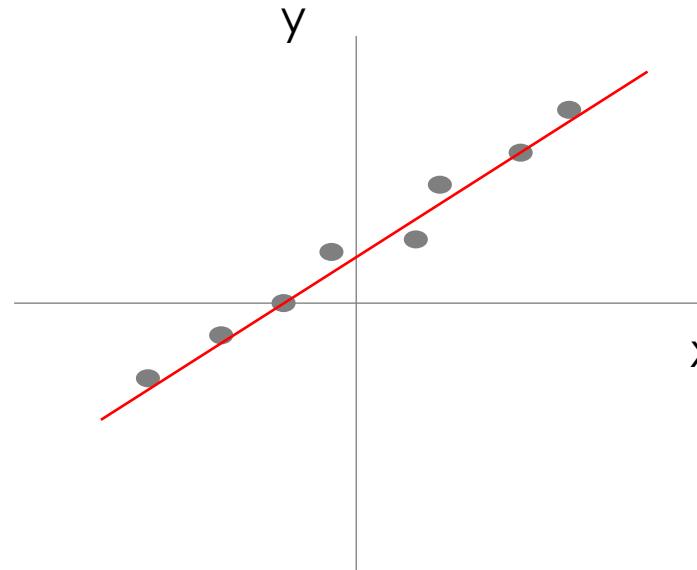
---



「線形」「非線形」という言葉は、  
データ分析でどのように使われるでしょうか？

# 線形と非線形

---

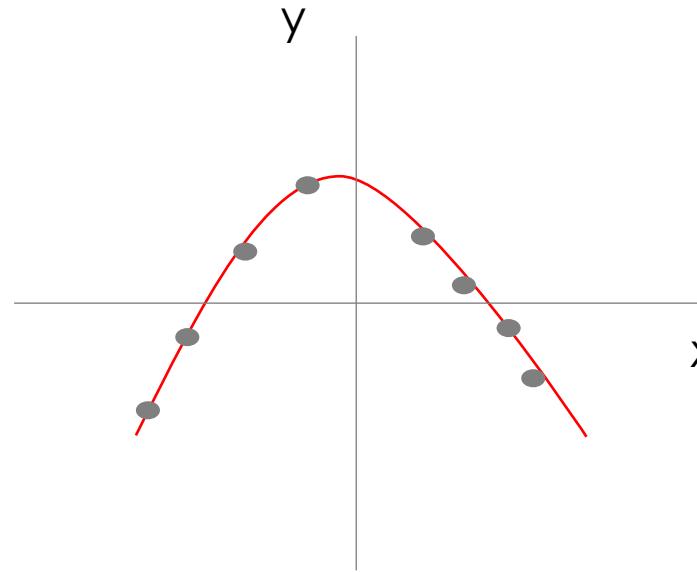


$$y = ax + b$$

データでは線形の関係が見られる。  
線形の手法を使って回帰モデルをつくる。

# 線形と非線形

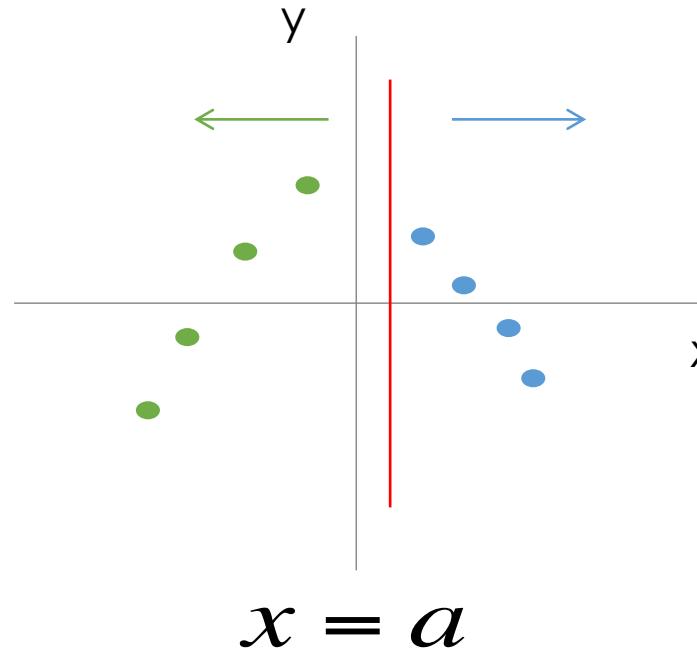
---



$$y = ax^2 + bx + c$$

データでは非線形の関係が見られる。  
線形の手法(3つの項の線形和)を使って回帰モデルをつくる。

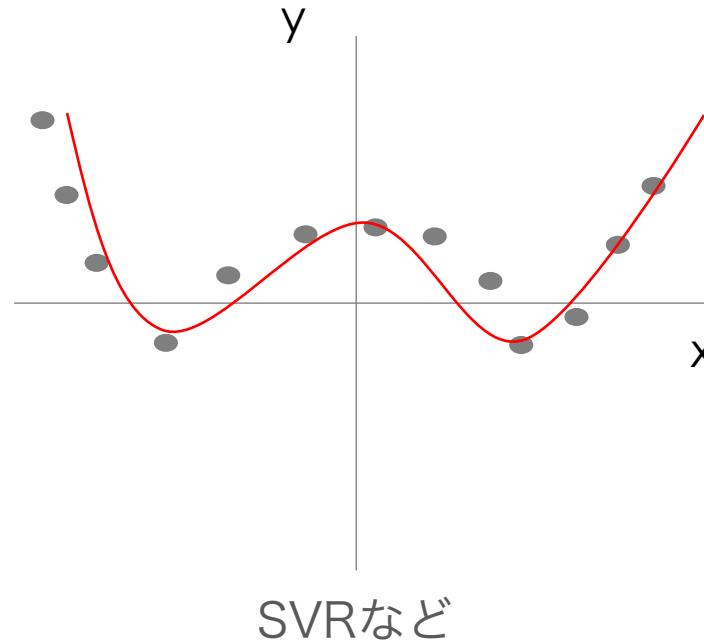
# 線形と非線形



データでは非線形の関係が見られる。  
線形の識別境界をつくる。

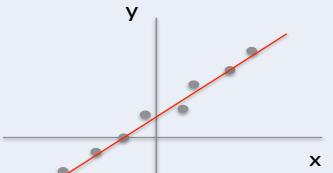
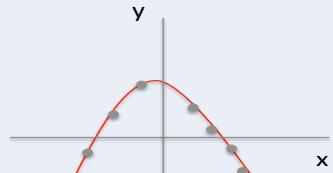
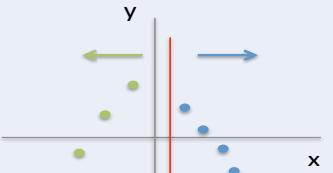
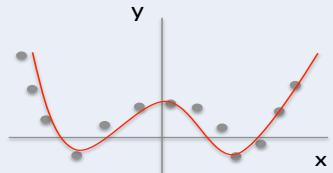
# 線形と非線形

---



データでは非線形の関係が見られる。  
非線形の回帰モデルをつくる。

# 線形と非線形

例	 $y = ax + b$	 $y = ax^2 + bx + c$	 $x = a$	 SVRなど
データ間の関係性	データでは線形の関係が見られる。	データでは非線形の関係が見られる。	データでは非線形の関係が見られる。	データでは非線形の関係が見られる。
手法	線形の手法を使って回帰モデルをつくる。	線形の手法を使って回帰モデルをつくる。	線形の識別境界を作る。	非線形の手法を使って回帰モデルをつくる

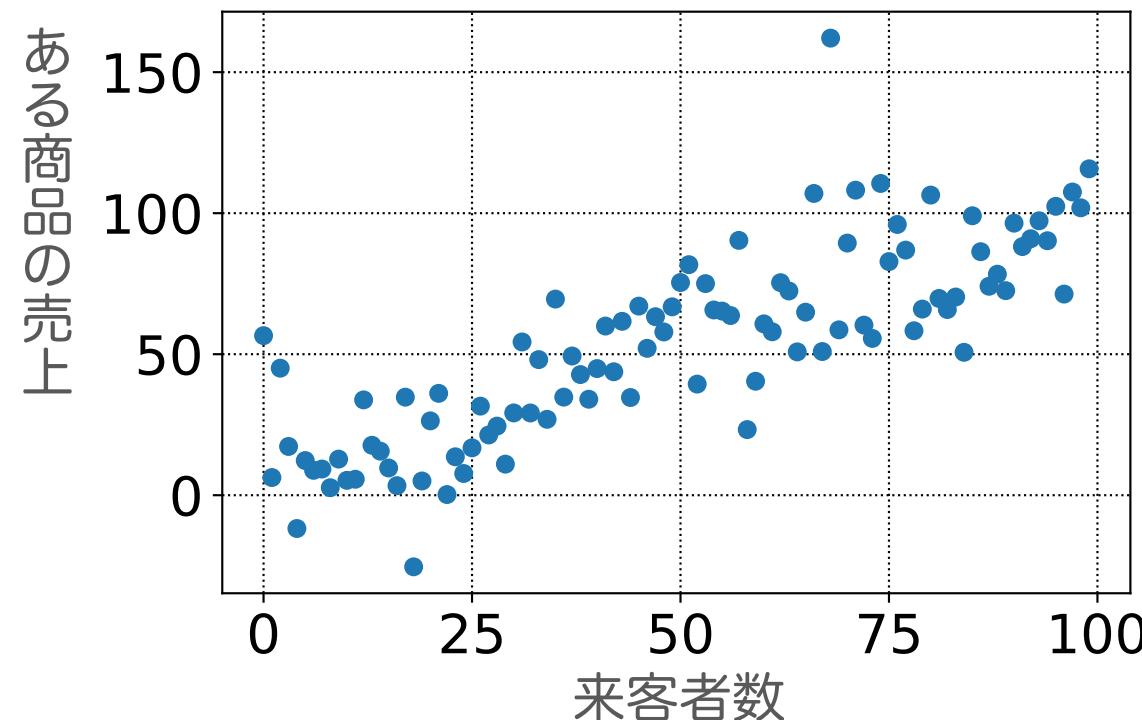
# データの関係が線形か非線形かは意識しましょう

---

1. データが線形であれば、線形で解いた方が楽
2. データが非線形でも、線形で解ける場合がある
  - 例)  $y = ax^2 + bx + c$  のように多項式の線形和の関係に落としめる場合
3. アルゴリズムを選択する際には、対象とするデータが線形か非線形かを常に意識して選ぶことが重要

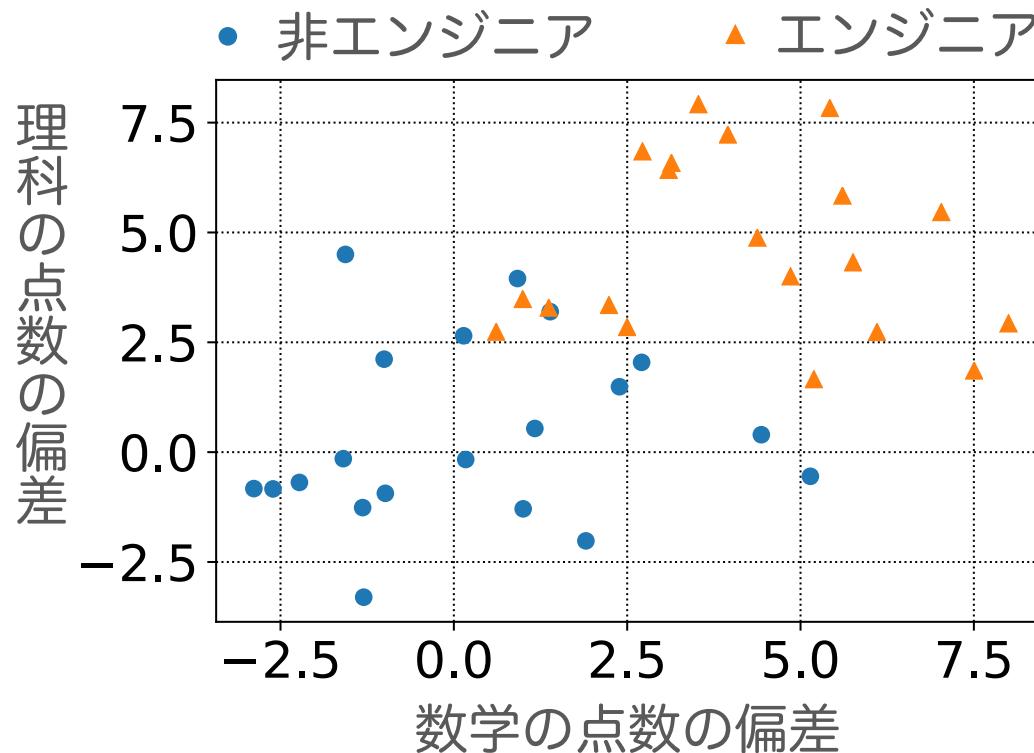
# 回帰 (Regression) とは

- ・回帰とは、数値を予測すること
  - ・教師信号が連続値のケースは回帰問題に属する
  - ・例) 来客数からある商品の売上を予測



# 分類 (Classification) とは

- ・分類とは、カテゴリを予測すること
  - ・教師信号がカテゴリを示す離散値のケースは分類問題に属する
  - ・例) 理科と数学の偏差から将来エンジニアになるか予測 (2カテゴリ、二值分類)



# 回帰と分類の違い

- 回帰と分類は、目的変数が異なる

	回帰	分類
説明変数 (特徴量、特徴ベクトル、素性、素性ベクトルとも呼ばれる)	$\mathbf{x} = (x_1, x_2, \dots, x_p)^T \in \mathbb{R}^p$ ( $x$ は $p$ 個の実数からなるベクトルという意味の数式)	
目的変数	$y \in \mathbb{R}$ ( $y$ は実数という意味の数式)	$y \in \{0, 1, \dots, C - 1\}$ ( $y$ は $C$ 個のカテゴリを示す離散値という意味の数式)

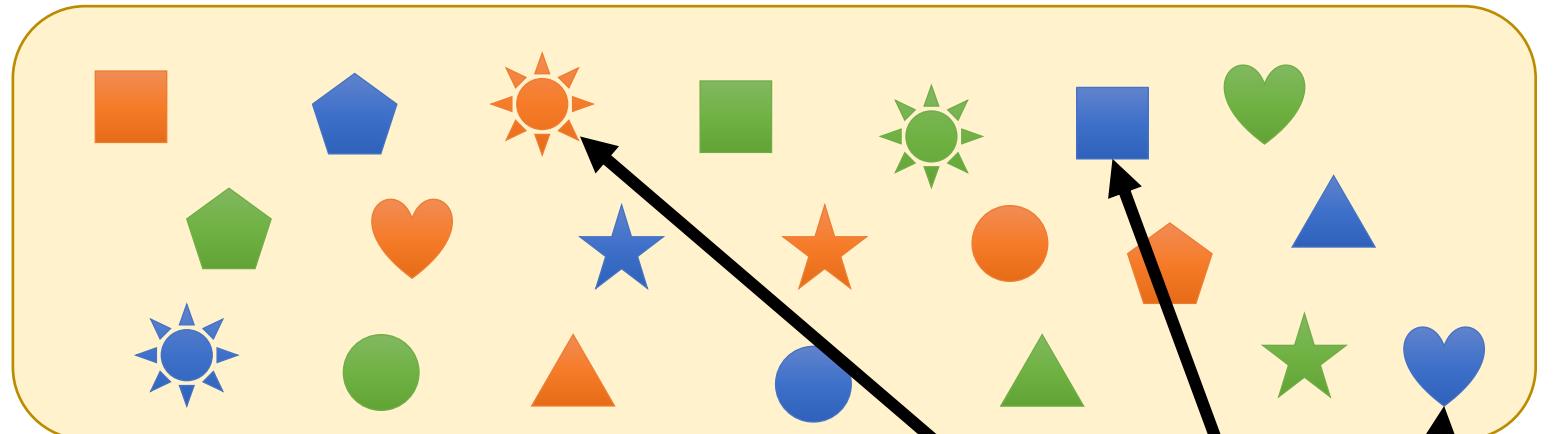
数学記号は、Appendixを参照

# 教師なし学習(クラスタリング)の代表的なアルゴリズム

日本語名	英語名	クラスタ数の決め方	分散	手法の概要
K-平均法	K-means	分析者が決める	等分散なデータでないとうまく分類できない	各データとクラスタ中心の距離の総和が最小になるように分類する方法
X-means	X-means	BICが決める		BICでクラスタ数を決めるようにしたK-means
混合ガウスモデル	Gaussian Mixture Model	分析者が決める		ガウス分布の重ね合わせによる確率分布を求め、属するクラスタを求める方法
ディリクレ過程 混合ガウスモデル	Dirichlet Process Gaussian Mixture Model	クラスタリングと同時に推定	等分散でないデータもうまく分類できる	クラスタに分けられるデータ要素の確率だけなく、クラスタ数も推定する方法  小さなクラスタとして判定され得るような微妙なデータがあると、それをクラスタリングしてしまい、真のクラスタ数よりも多いクラスタ数を推定してしまうことがある。

# 機械学習アルゴリズムのポイント

モデル（関数・確率分布）の集合

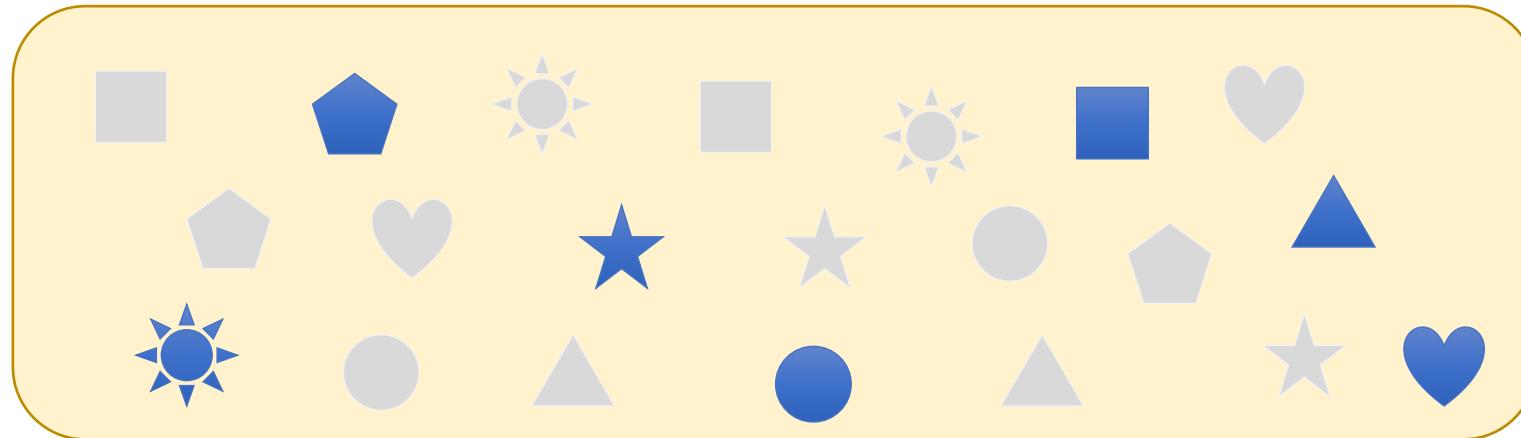


機械学習アルゴリズムとは  
学習データをうまく  
説明できるモデルを探すこと

しかし候補となるモデルが多すぎて、最適なものを見つけるのは困難…

# 機械学習アルゴリズムのポイント

モデル（関数・確率分布）の集合

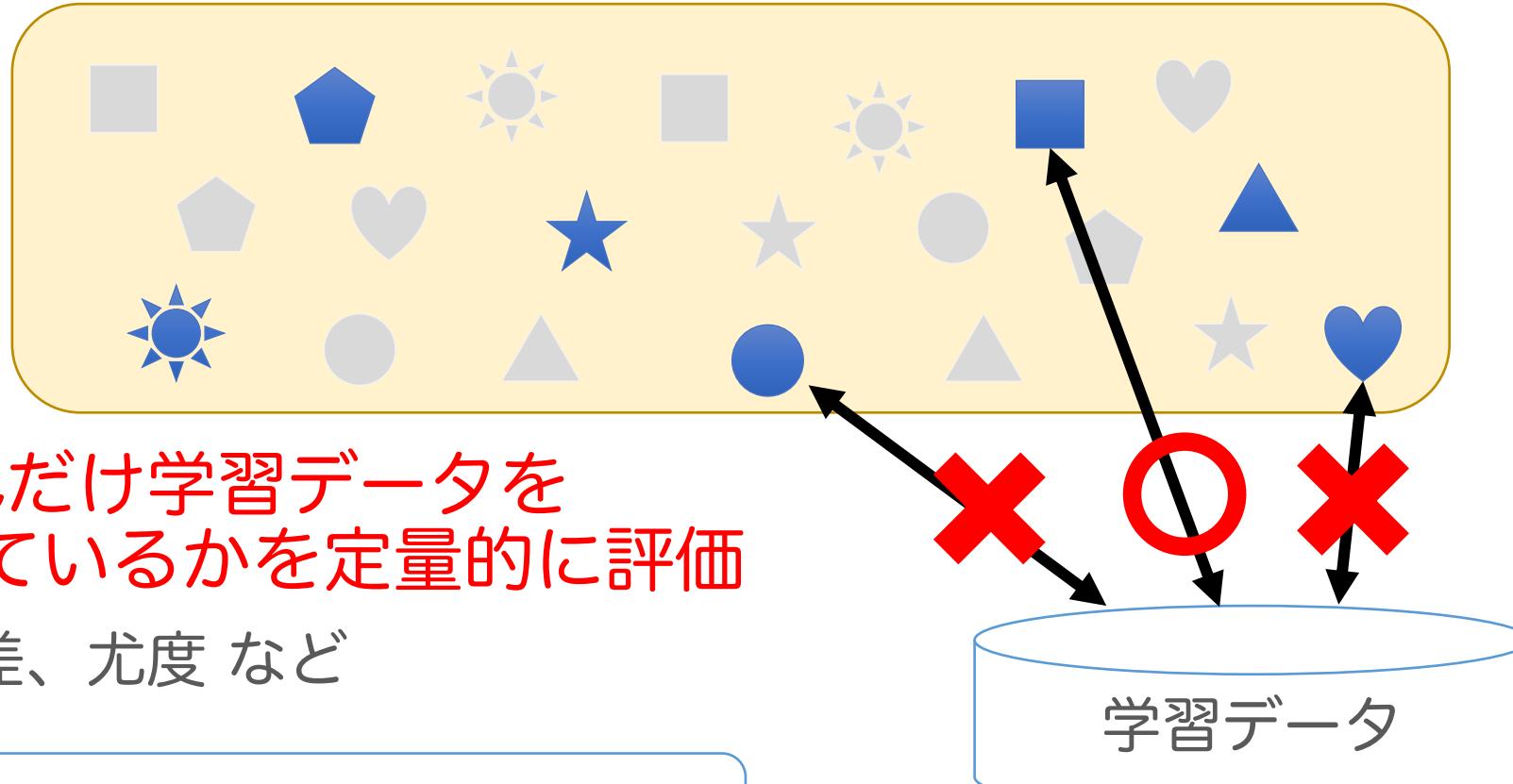


一般的にはデータの性質を考慮して、モデルの集合を制限  
例）線形関数、正規分布 など

1. モデルを決める

# 機械学習アルゴリズムのポイント

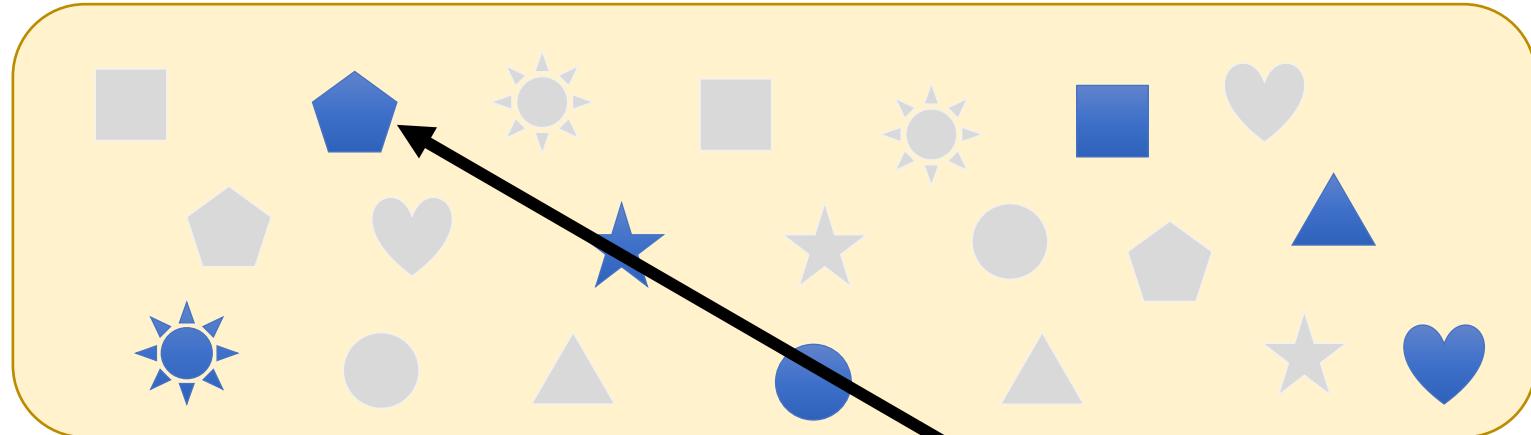
モデル（関数・確率分布）の集合



2. モデルの評価基準を決める

# 機械学習アルゴリズムのポイント

## モデル（関数・確率分布）の集合



良い評価値を出すモデルを探索  
(どのように最適化するか)

例) 最小二乗法、確率的勾配降下法など

学習データ

3. 良い評価値のモデルを探す（最適化）

# 機械学習アルゴリズムのポイント

---

1. データの性質に合わせてモデルを決める

2. モデルの評価基準を決める

3. 良い評価値を出すモデルを探す（最適化）

問題やデータに対する事前知識を用いて  
1~3を設定することで機械学習アルゴリズムが完成する

# 機械学習アルゴリズムのポイント

- 問題やデータに対する事前知識を用いて、「モデル」「評価基準」「最適化（の方法）」を設定することで機械学習アルゴリズムが完成する

## モデル

- データの性質に合わせてモデルを決める

## 評価基準

- モデルの評価基準を決める

## 最適化

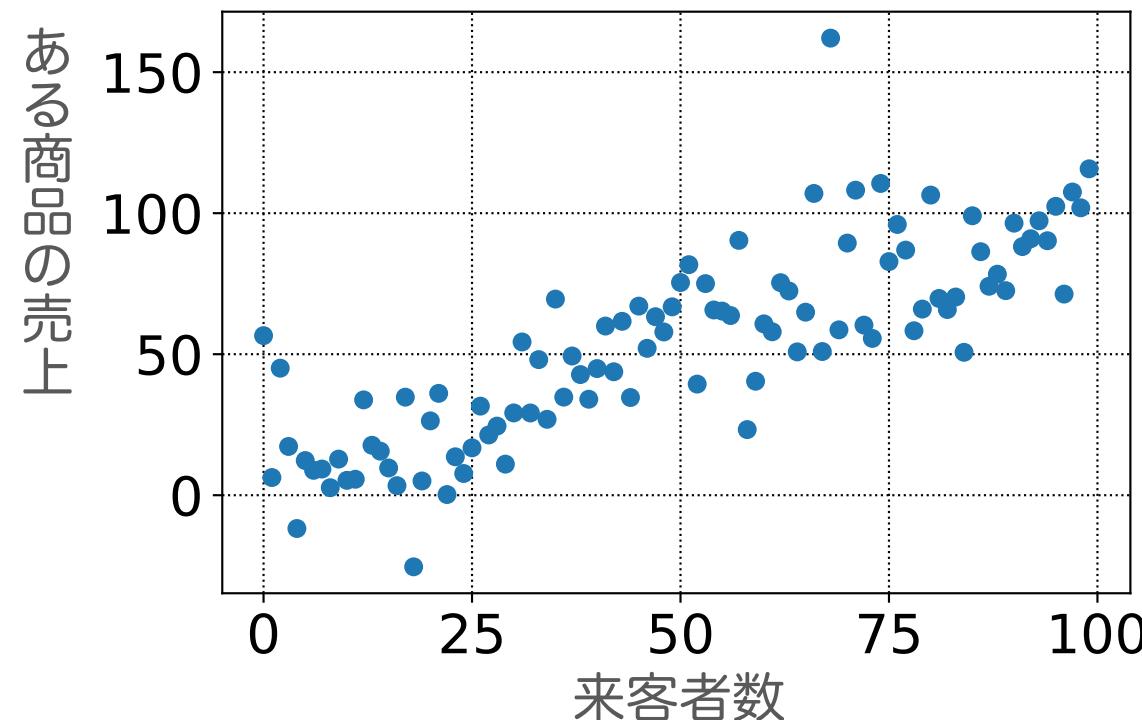
- 良い評価値を出すモデルを探す

# 線形回帰

1. 回帰とは
2. 線形モデル（多項式モデル）
3. モデルの評価基準（二乗誤差）
4. モデルパラメータの最適化（最小二乗法）
5. 線形回帰モデルの解釈
6. 線形回帰まとめ（モデル、評価基準、最適化の観点から）

# 回帰 (Regression) とは

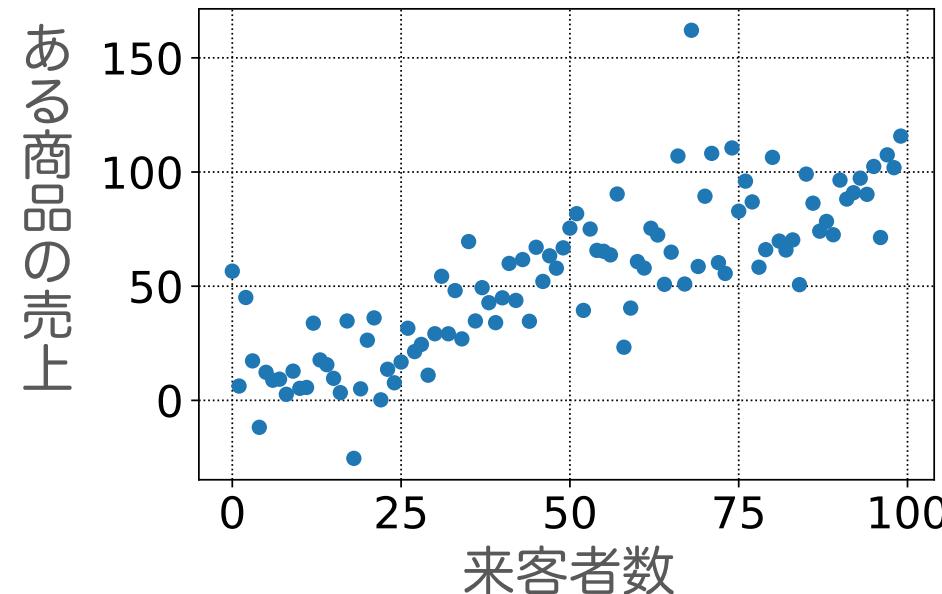
- ・回帰とは、数値を予測すること
  - ・教師信号が連続値のケースは回帰問題に属する
  - ・例) 来客数からある商品の売上を予測



## 問題

---

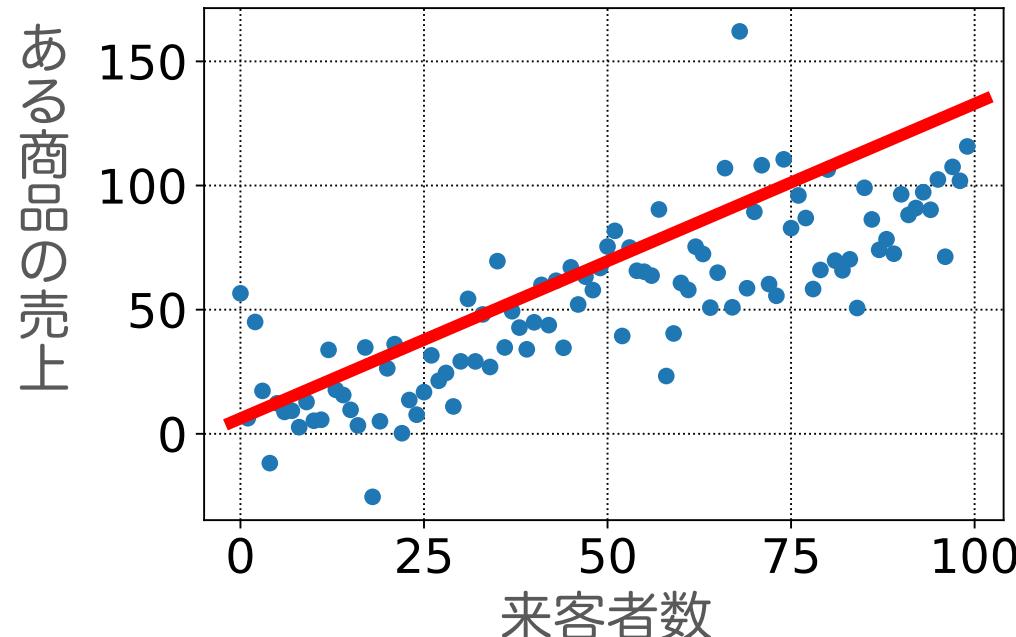
- ・小売店にて、ある商品の売り上げと来店客数の関係を調査したところ  
下図の結果がえられた
- ・このデータから、ある来店客数のときにおける商品の売上を予測するモデルを  
つくりたい
- ・どのようなモデルをつくればいいか？



## 問題に応えるための1つのアイデア

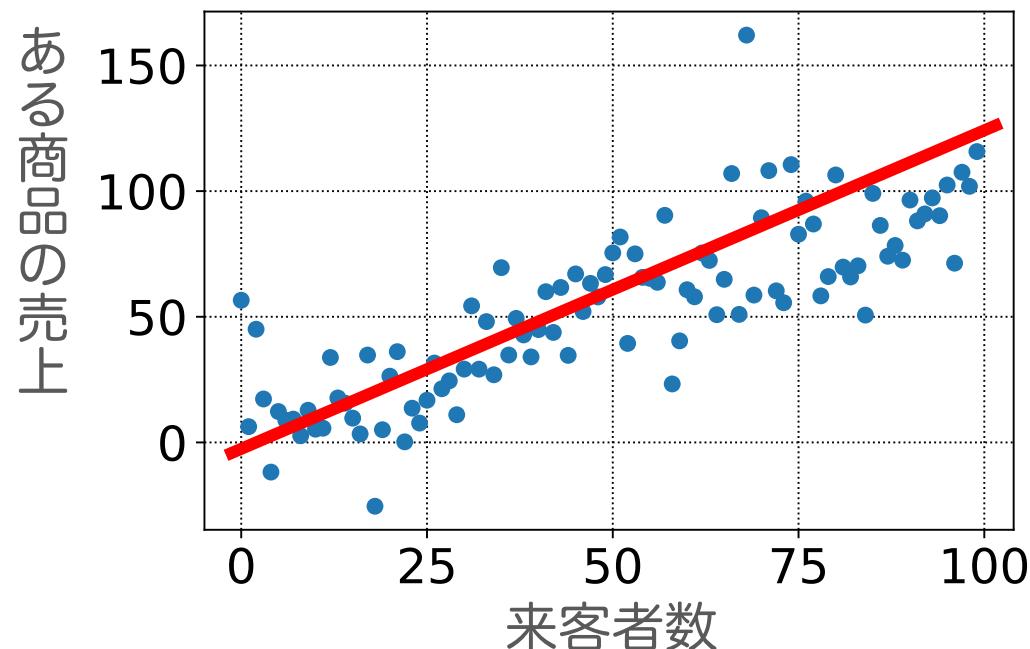
---

- ・直線を引くのがよさそう
- ・直線といえば、 $y = ax + b$ を思いつくが、傾き $a$ と切片 $b$ はどうやって決めればいいのか？



# 線形回帰モデル (Linear Regression Model) | 説明変数が1次元のとき

- ・ (2次元の) 直線: 切片と傾きが学習可能なモデルパラメータ
- ・ 切片と傾きのベクトル表現を  $w = (w_0, w_1)^T$ 、モデルへの入力 (説明変数) を  $x = (1, x_1)^T$ としたとき、モデルの出力  $\hat{Y}(x; w)$  は次式で定義する



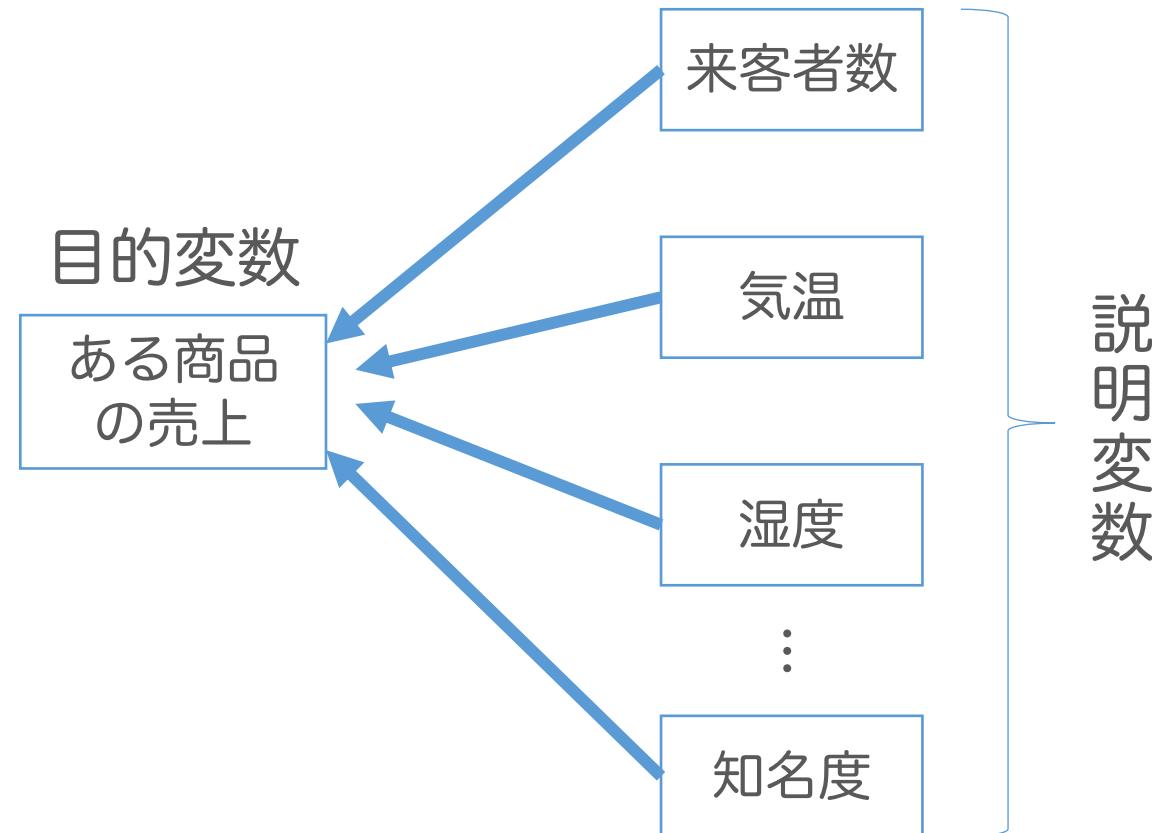
$$\hat{Y}(x; w) = w^T x = w_0 + w_1 x_1$$

ある重み  $w$  のときに、  
入力(説明変数) $x$ を入  
れた時の出力。  
推定値なので、 $Y$ の  
上に $\hat{\cdot}$ がついている。

ベクトル  $w$ を  
転置したもの  
と入力(説  
明  
変  
数  
のベ  
ク  
ト  
ル  
 $x$ の積。

# 売上と来客者数の関係さえわかれば十分？

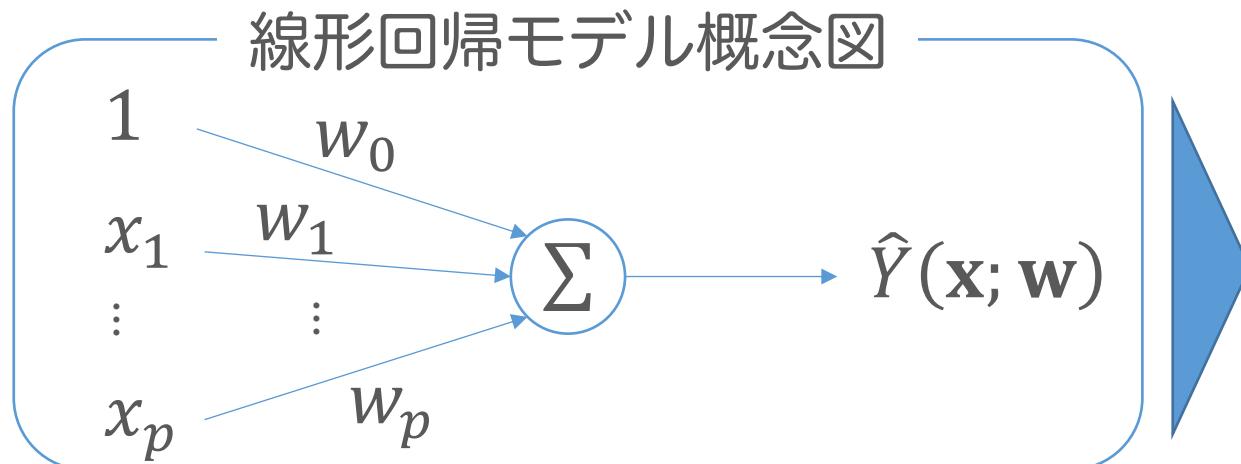
- ・ 売上に関係がありそうな要因は、来客者数以外にも色々ありそう
- ・ 線形回帰モデルの考え方で、他の要因も考慮して予測できないだろうか？



# 線形回帰モデル (Linear Regression Model) | 説明変数が多次元の場合

- 学習可能なパラメータを  $\mathbf{w} = (w_0, w_1, \dots, w_p)^T \in \mathbb{R}^{p+1}$ 、モデルへの入力  $\mathbf{x} = (1, x_1, \dots, x_p)^T$ としたとき、線形回帰モデルの出力  $\hat{Y}(\mathbf{x}; \mathbf{w})$  を次式で定義する

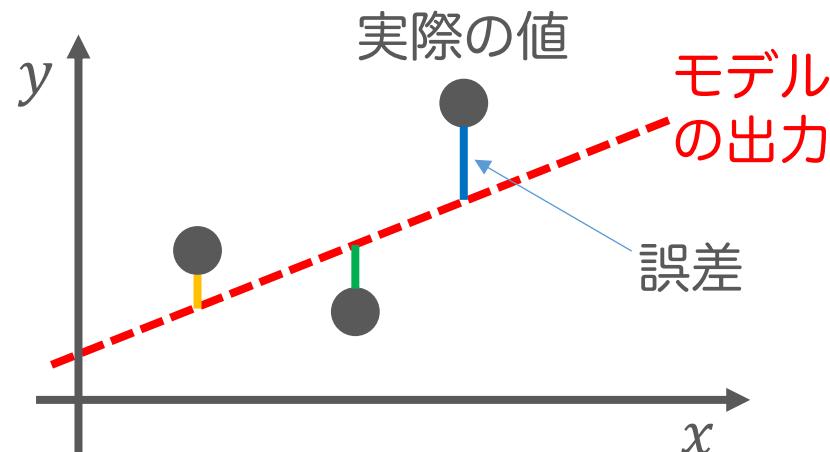
$$\hat{Y}(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} = \sum_{i=0}^p w_i x_i = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_p x_p$$



各説明変数に重みづけして  
足し合わせるシンプルなモデル  
学習データをよく説明できる  
ように重み  $w$  を学習

# モデルの評価基準（二乗誤差）

- モデルの良さは出力と教師データの二乗誤差で評価



## 二乗誤差の考え方

1辺の長さが誤差の正方形の面積の総和

$$E_D = \frac{1}{2} (\text{Yellow Square} + \text{Green Square} + \text{Blue Square})$$

を最小化すれば、誤差も小さくなる！

最適化の観点で良い評価指標（詳しくは最適化講座にて！）

## モデルの評価基準（二乗誤差）

---

- $N$ 個の学習データのもとで二乗誤差の総和 $E_D$ は次式で定義される

$$E_D = \frac{1}{2} \sum_{n=1}^N \left( \underbrace{\hat{Y}(\mathbf{x}^{(n)}; \mathbf{w})}_{y\text{の推定値}} - \underbrace{y^{(n)}}_{y\text{の正解値}} \right)^2$$

$\mathbf{x}^{(n)}$  :  $n$ 番目のデータの説明変数

$y^{(n)}$  :  $n$ 番目のデータの目的変数  
=教師データ

# モデルパラメータの最適化（最小二乗法）

- 二乗誤差を最小化する方法 = **最小二乗法**
- 線形回帰モデルは、二乗誤差を最小化する重みは解析的に（=数式変形によって）求めることが可能
- 二乗誤差が最小になる条件式  $\frac{\partial E_D}{\partial w} = 0$  を  $w$ について解くと…

最適な重み

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

導出の詳細は

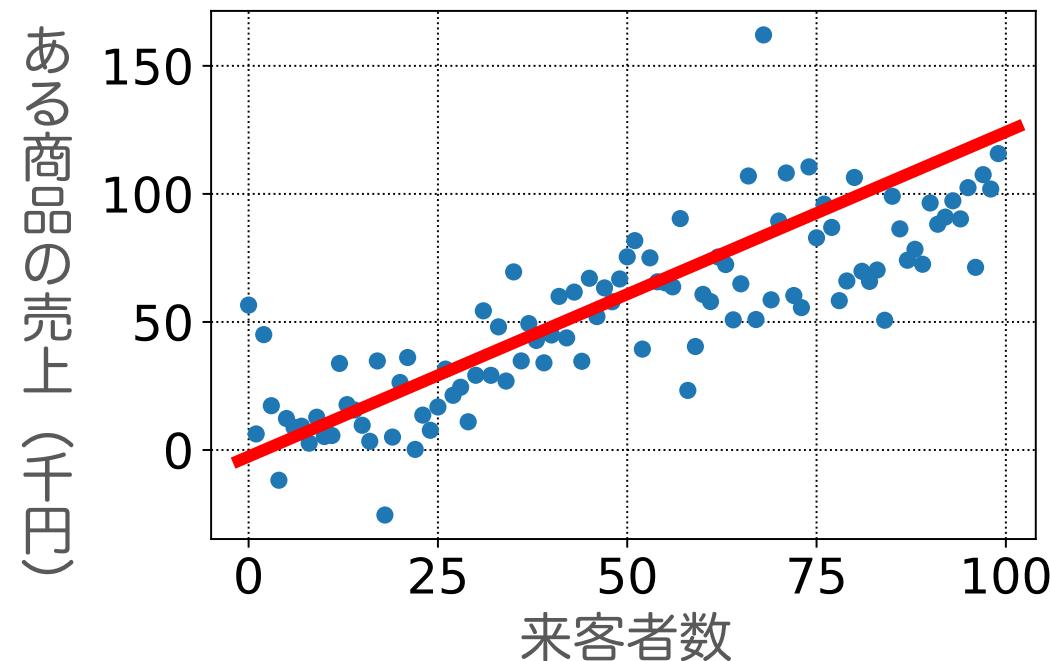
『ITエンジニアのための機械学習理論入門』  
p.64~p.66を参照

$$\Phi = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \vdots \\ \mathbf{x}^{(N)T} \end{pmatrix} = \begin{pmatrix} 1 & x_1^{(1)} & \cdots & x_p^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \cdots & x_p^{(N)} \end{pmatrix}$$

$$y = (y^{(1)}, y^{(2)}, \dots, y^{(N)})^T$$

# 線形回帰モデルの解釈 | 直線の意味とは？係数の意味とは？

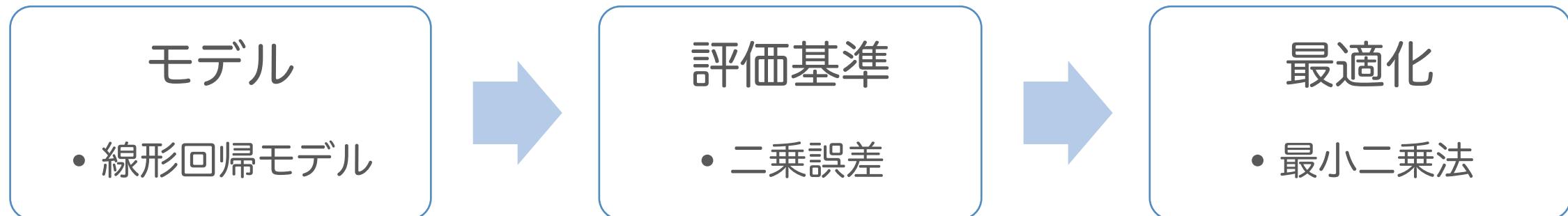
- この直線は来客者数50人なら売上は約6万円になるという関係を示唆
- モデルを得たら直観に反していないか必ずチェックしよう
  - 可視化してデータに重ねてみる、係数をチェックしてみるなど



- ・ 線形回帰モデルの係数は、偏回帰係数とも呼ばれる
- ・ 偏回帰係数は、他の変数を固定してその変数だけを動かしたときに得られる目的変数Yの変化量を意味する
- ・ 係数プロットという一覧化の手法で検討するのが一般的
  - ・ DAY3の「特徴選択」で詳細を説明します

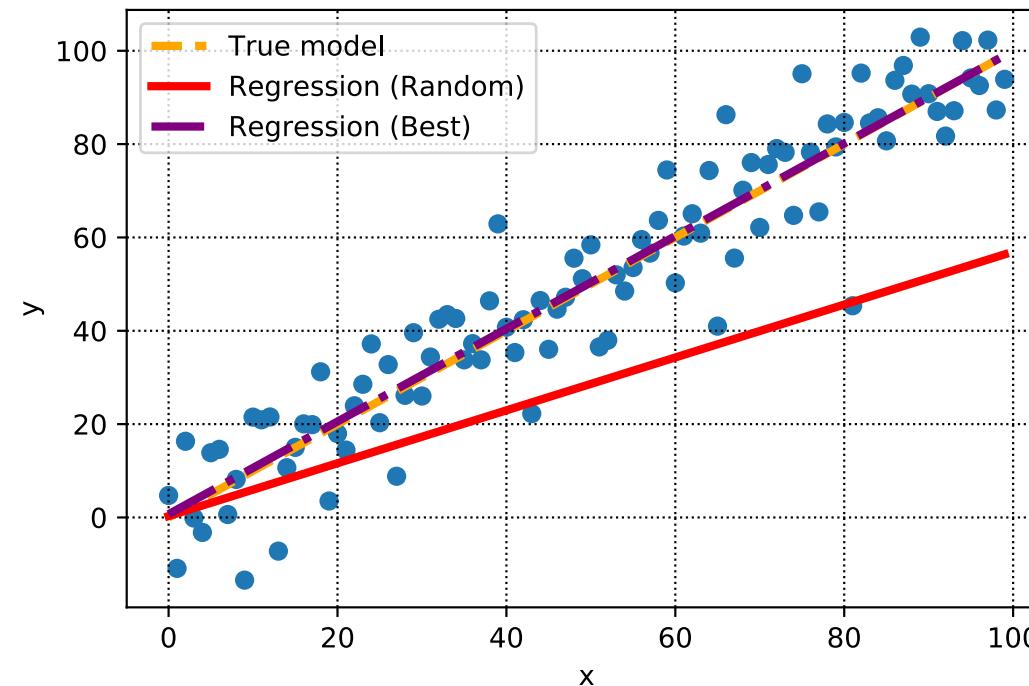
# 線形回帰まとめ（モデル、評価基準、最適化の観点から）

- ・ 線形回帰（回帰分析）：ある変数を他の変数で説明するための直線を求める  
こと
  - ・ 前者の変数のことと目的変数、後者の変数を説明変数と呼ぶ
- ・ 例）1日の売上を目的変数、気温を説明変数にし、その関係性を表す直線を求める
  - ・ この直線を求まれば、明日の天気予報を用いて売上を予測できる



## [演習] 線形回帰 (1\_linear\_regression\_pseudo\_data.ipynb)

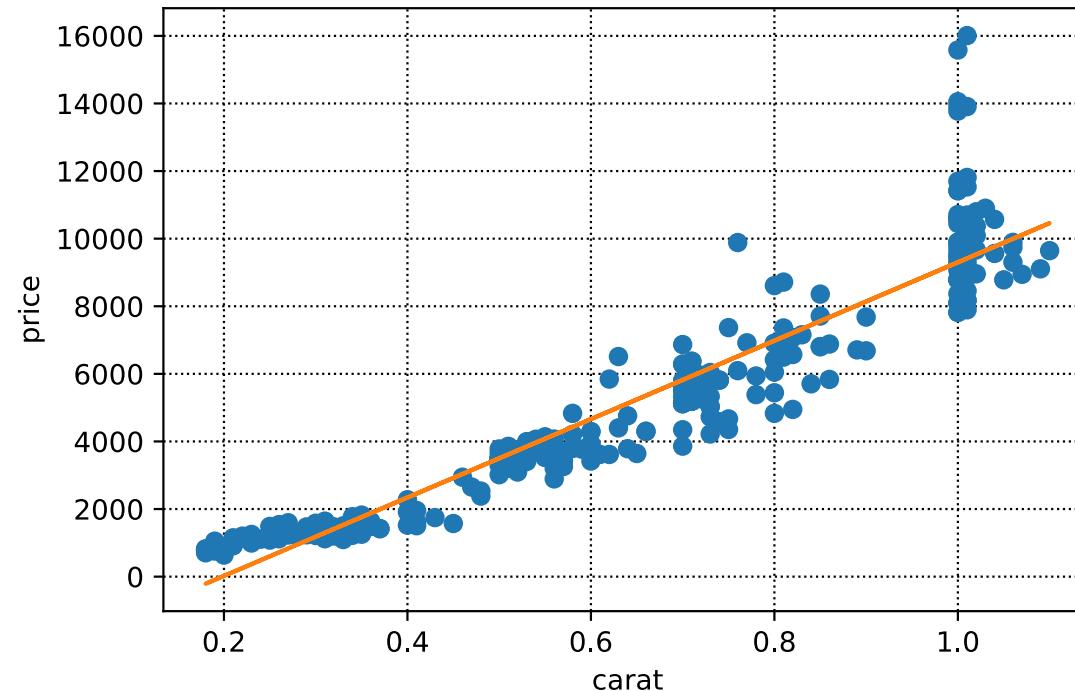
- まずは簡単なデータで線形回帰を試してみましょう
- ランダムに直線のパラメータを決めたものよりも、最小二乗法によって決めたパラメータの方が二乗誤差が小さくなることを確認しましょう



## [演習] 線形回帰 (2\_linear\_regression\_real\_data\_trainee.ipynb)

---

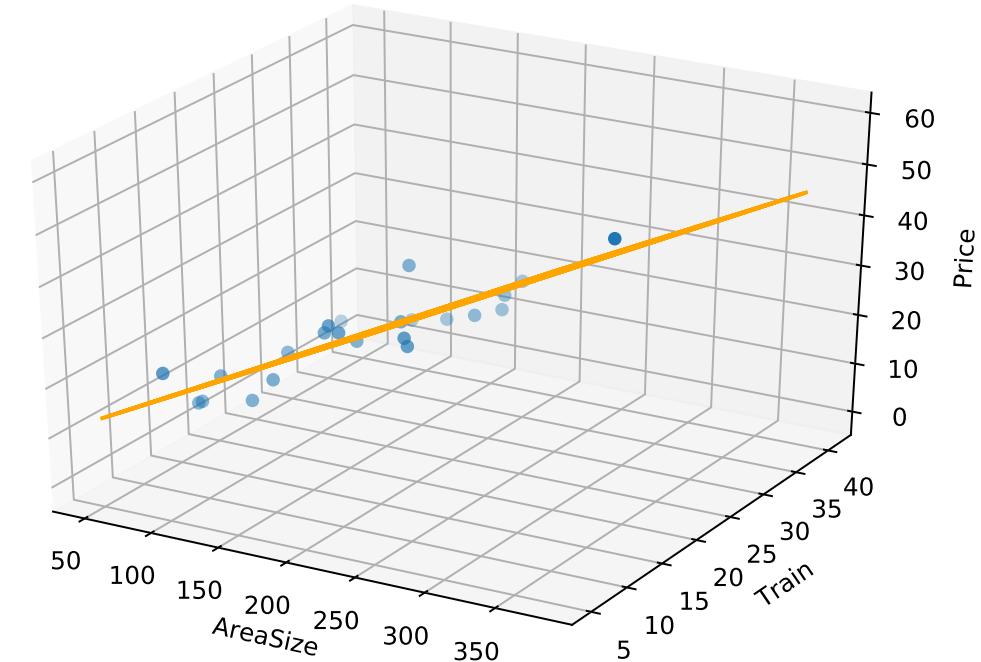
- ・ 実際のデータで線形回帰を実行してみましょう
- ・ ダイヤモンドのカラット数からその価格を予測するモデルを作りましょう



## [演習] 線形回帰 (3\_linear\_regression\_multi\_psedo\_data.ipynb)

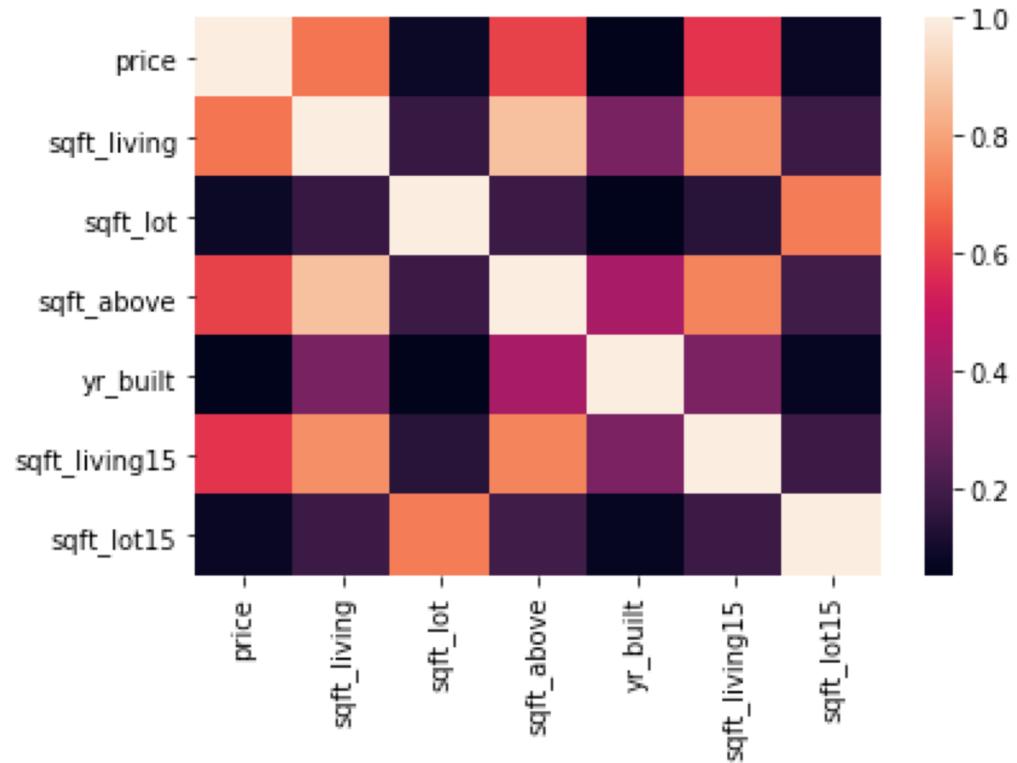
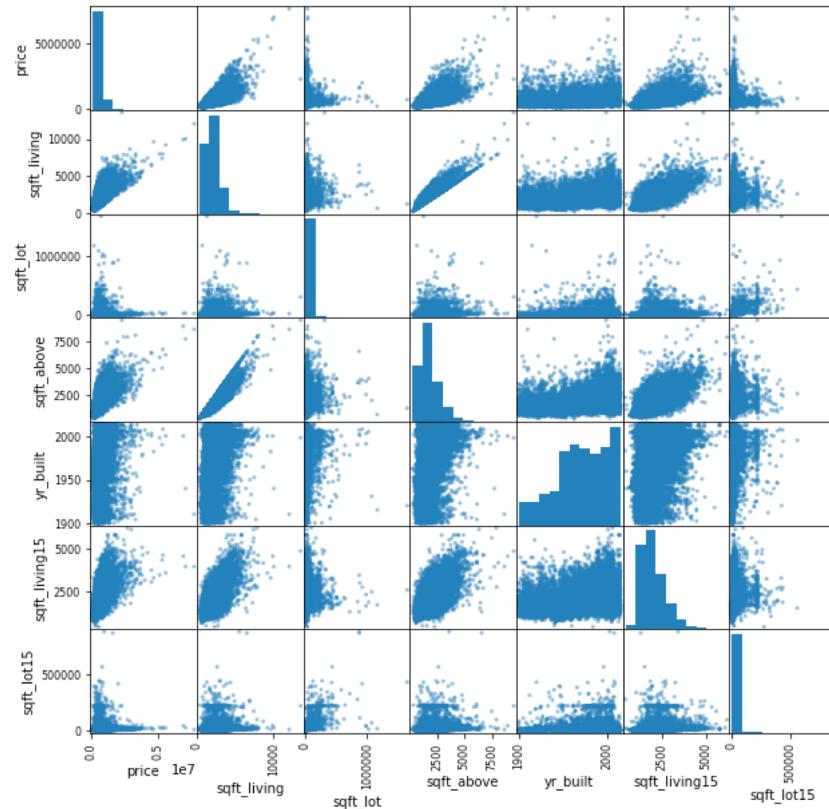
- 説明変数が多い住宅価格データセットで線形回帰をしてみましょう
- さらに説明変数ごとに基本的な統計量を確認してみましょう

	AreaSize	HouseSize	PassedYear	Price	Train	Walk
count	23.000000	23.000000	23.000000	23.000000	23.000000	23.000000
mean	144.139130	81.356522	9.834783	19.395652	23.260870	6.217391
std	70.086095	26.436955	4.023071	11.605151	10.247915	5.062846
min	50.000000	48.700000	3.100000	5.500000	5.000000	1.000000
25%	99.000000	66.300000	5.700000	11.600000	16.000000	2.500000
50%	139.600000	77.900000	10.500000	17.600000	19.000000	5.000000
75%	173.300000	86.750000	13.150000	25.900000	34.500000	7.500000
max	379.800000	163.700000	14.700000	59.500000	41.000000	20.000000



## [演習] 線形回帰 (4\_linear\_regression\_multi\_real\_data\_trainee.ipynb)

- データ数が大きな住宅価格データセットで線形回帰をしてみましょう
- 統計量や散布図行列、相関係数なども確認してみましょう



## グループワーク（線形回帰）

---

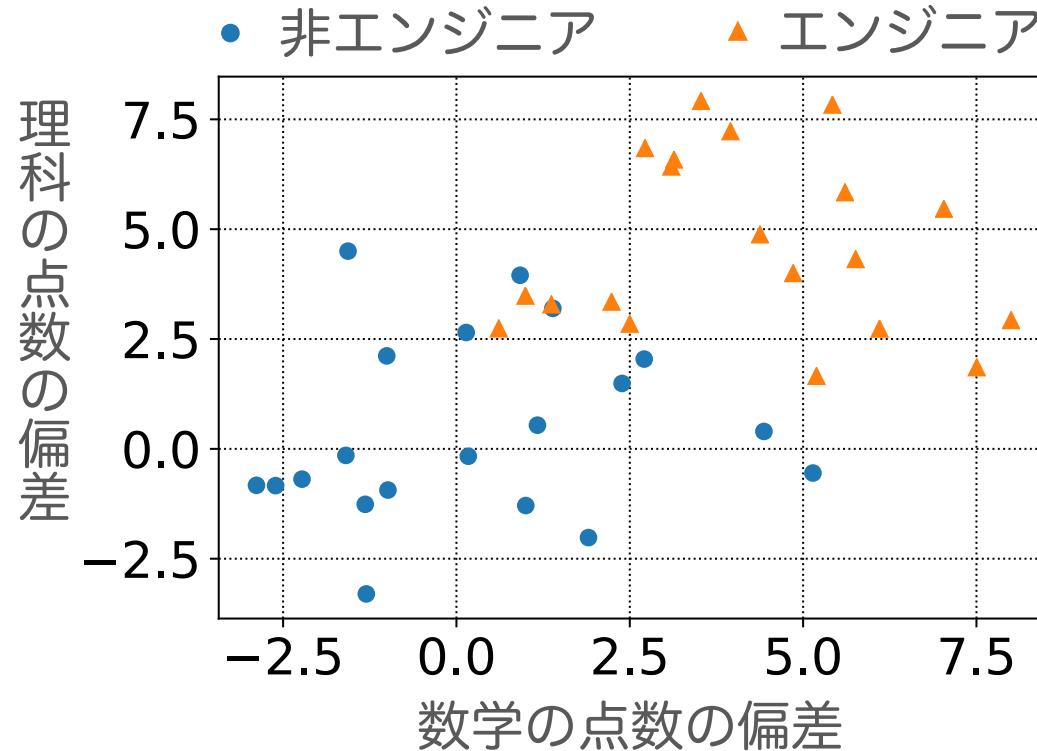
- ある洋菓子チェーン店はすべての店舗が駅前に立地している
  - あなたは売上予測モデルを構築し、最も売上が高くなりそうな場所に新しい店舗を出店することを思いついた
  - 3人程度のグループを作り、以下のディスカッションを進めましょう
- 
- 説明変数として用いるべき要因をできるだけ多く挙げてください(5分)
  - 挙げた要因から重要そうなものを5つ程度選択してください  
その要因が、なぜ重要なのかを説明できるようにしてください(5分)
  - その要因に関するデータはどのように取得するのか、  
そして、どの程度データが集まる見込みがあるか検討しましょう(5分)
  - 自分のグループで最終的に使うことを決めた要因、その理由、  
データの取得方法、データ量の見込みを全体に向けて発表しましょう(各2分)

# ロジスティック回帰

1. 分類とは
2. データ生成過程の導入
3. シグモイド関数
4. 最尤推定
5. 確率的勾配降下法
6. ロジスティック回帰まとめ（モデル、評価基準、最適化の観点から）

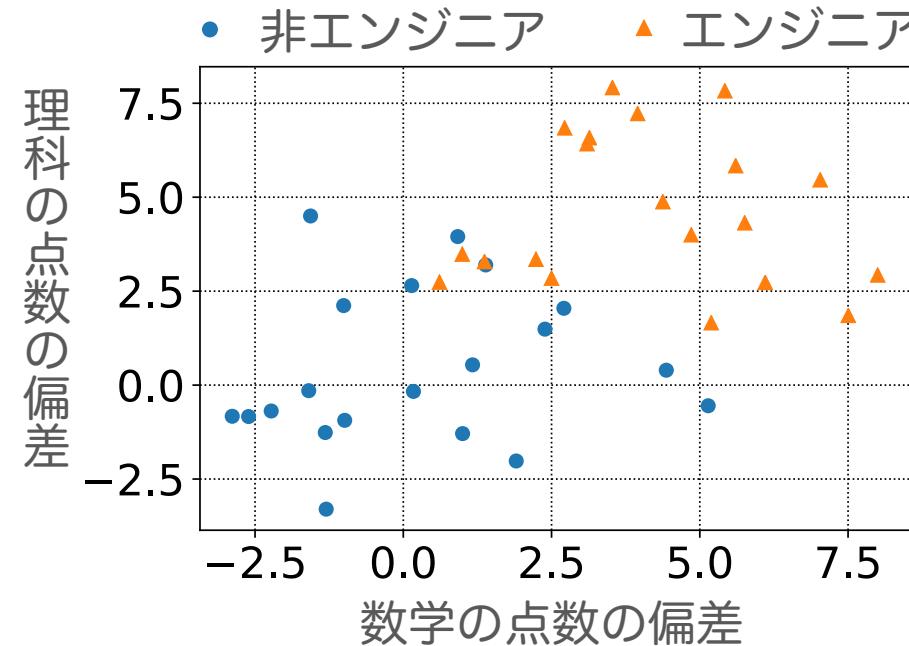
# 分類 (Classification) とは

- 分類とは、カテゴリを予測すること
  - 教師信号が離散値（カテゴリ変数）のケースは分類問題に属する
  - 例）理科と数学の偏差から将来エンジニアになるか予測（2カテゴリ、二値分類）



## 問題

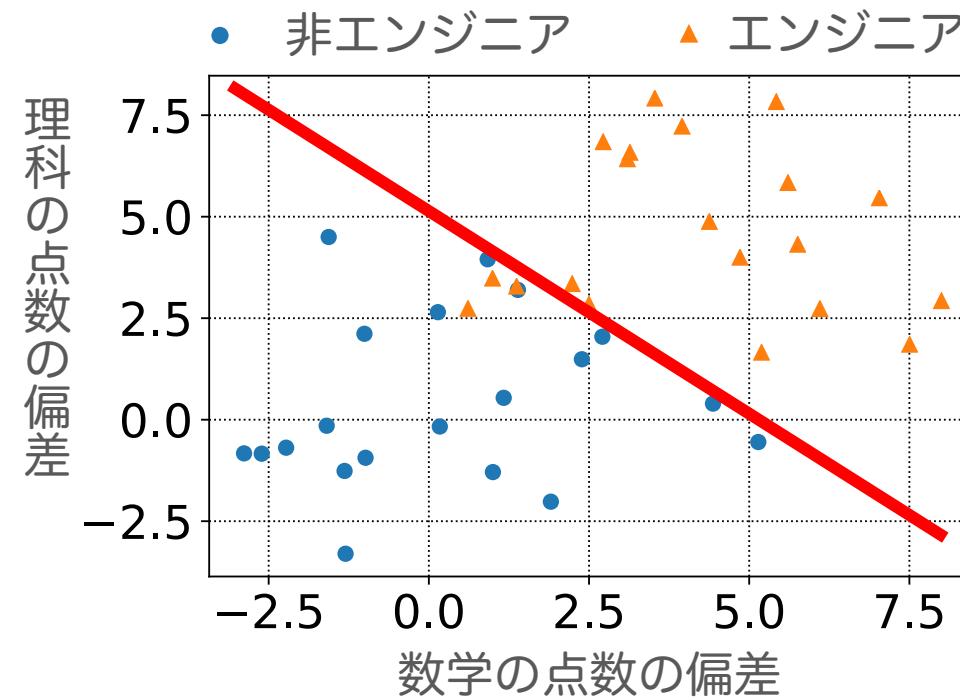
- あるエンジニアと非エンジニアを対象として、学生時代の理科および数学の偏差を調査したところ下図の結果が得られた
- このデータから、ある学生が将来エンジニアになるかを予測するモデルをつくりたい
- どのようなモデルをつくればいいか？



# 問題に応えるための1つのアイデア

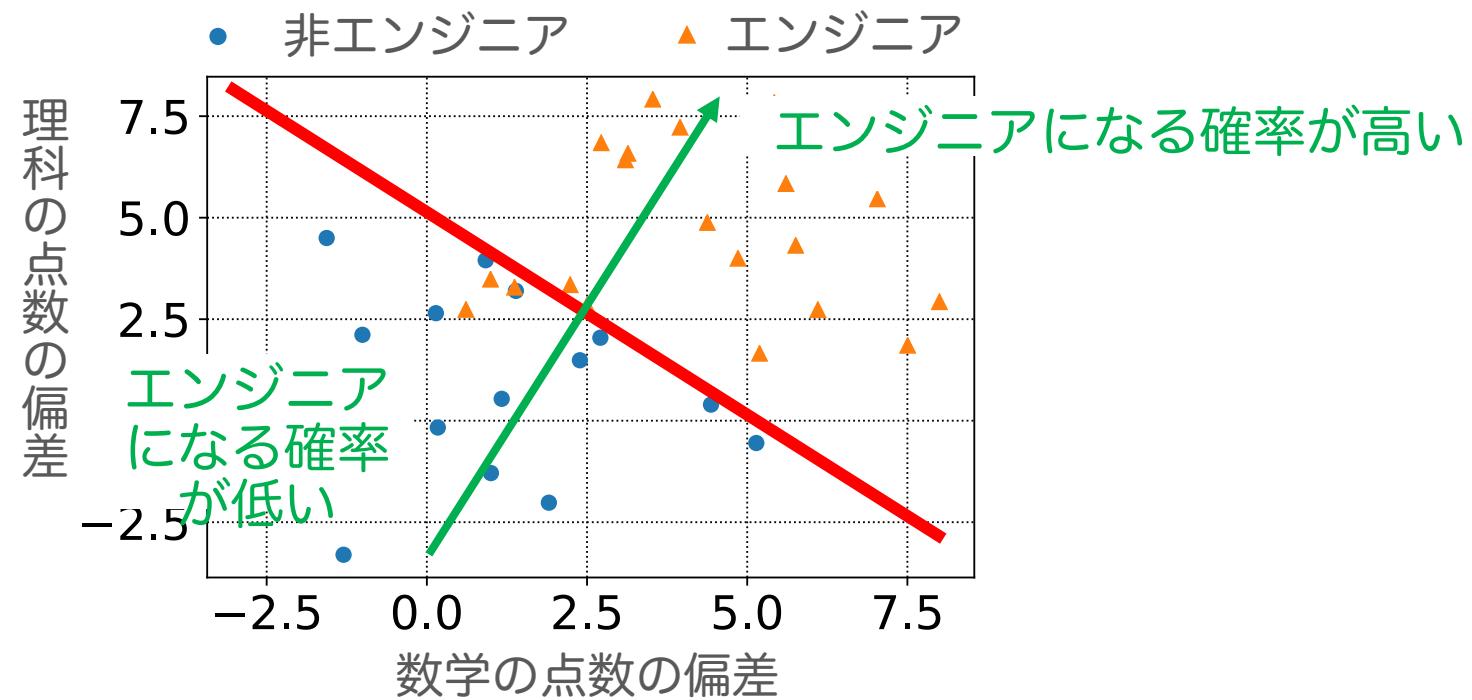
---

- 直線の境界線を引くのがよさそう
  - これには、線形回帰モデルの考え方方が使えるか？



## 問題に応えるための1つのアイデア

- さらに、エンジニアになる確率（＝確信度）を表現できるとうれしい



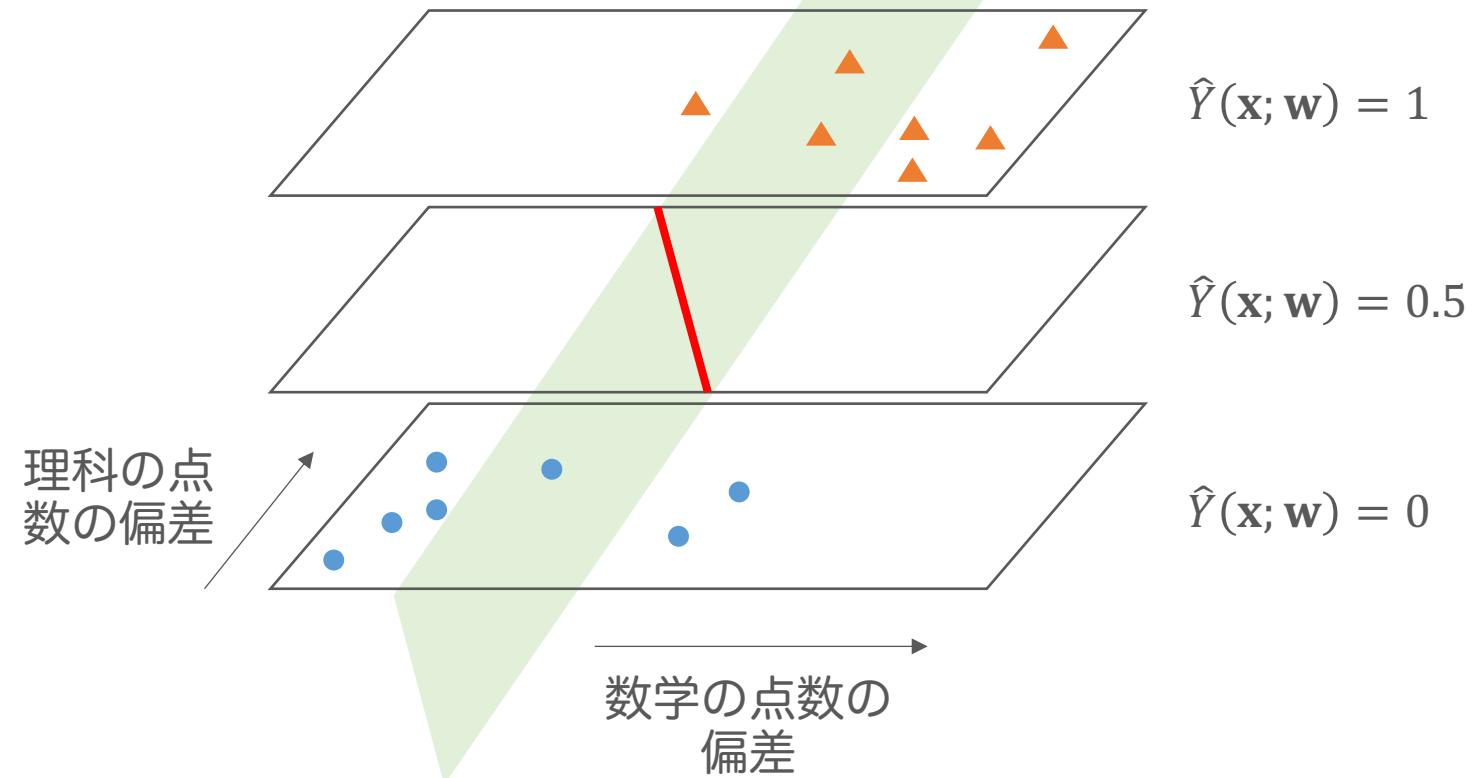
# 線形回帰モデルをどうやって使うのか？

エンジニアを1、非エンジニアを0として、線形回帰モデルを当てはめてみる

↓ 線形回帰  
モデルによる予測結果

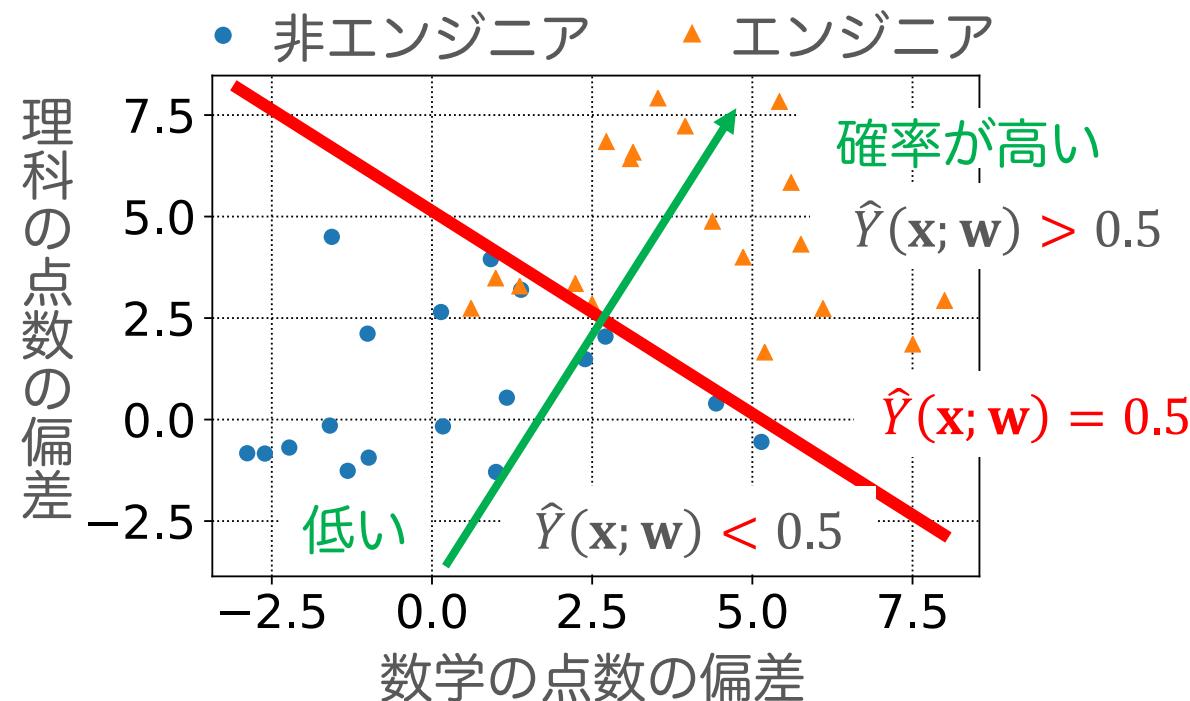
線形回帰モデルの出力が以下を満たす  
直線を**境界線**とする

$$\hat{Y}(\mathbf{x}; \mathbf{w}) = 0.5$$



緑平面の中の $\hat{Y}(\mathbf{x}; \mathbf{w}) = 0.5$ となる直線を境界線として採用する。  
得られた赤線を境界線として採用すれば、数学の点数の偏差と理科の点数の偏差で2値を分類できる

# データ発生確率の導入



線形回帰モデルの出力が以下を満たす直線を**境界線**とする

$$\hat{Y}(x; w) = 0.5$$

出力 $\hat{Y}(x; w)$ の値が大きい  
= エンジニアである確率が高い  
(出力が小さければ確率は低い)

出力 $\hat{Y}(x; w)$ を**確率値**に変換

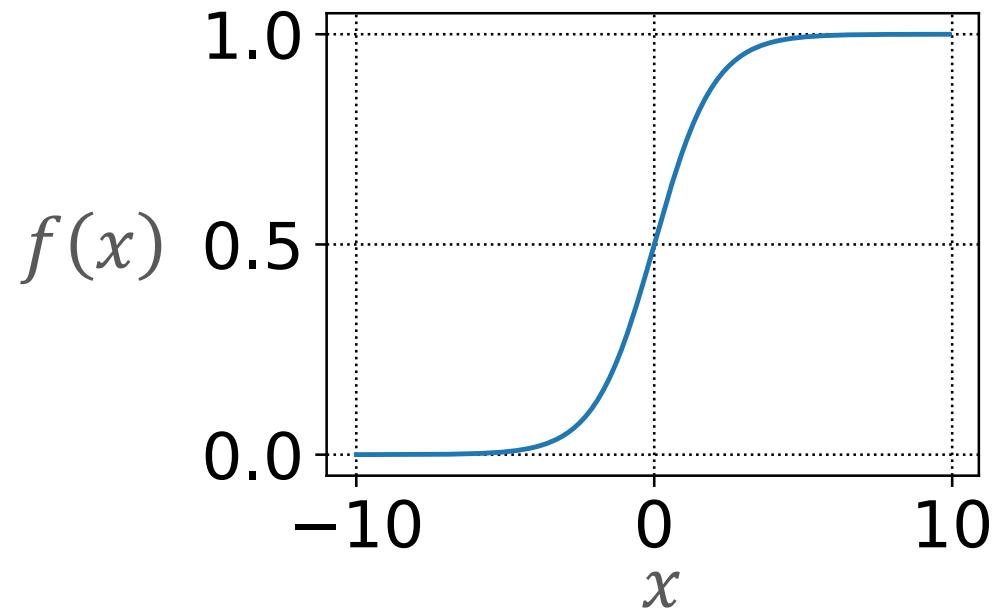
# シグモイド関数

- 出力値をシグモイド関数によって確率値に変換
  - シグモイド関数：入力が $-\infty$ から $\infty$ に大きくなるにつれて、出力が0から1までなめらかに変化する関数

シグモイド関数

定義： $f(x) = \frac{1}{1 + \exp(-x)}$

性質： $\frac{d}{dx} f(x) = f(x)(1 - f(x))$

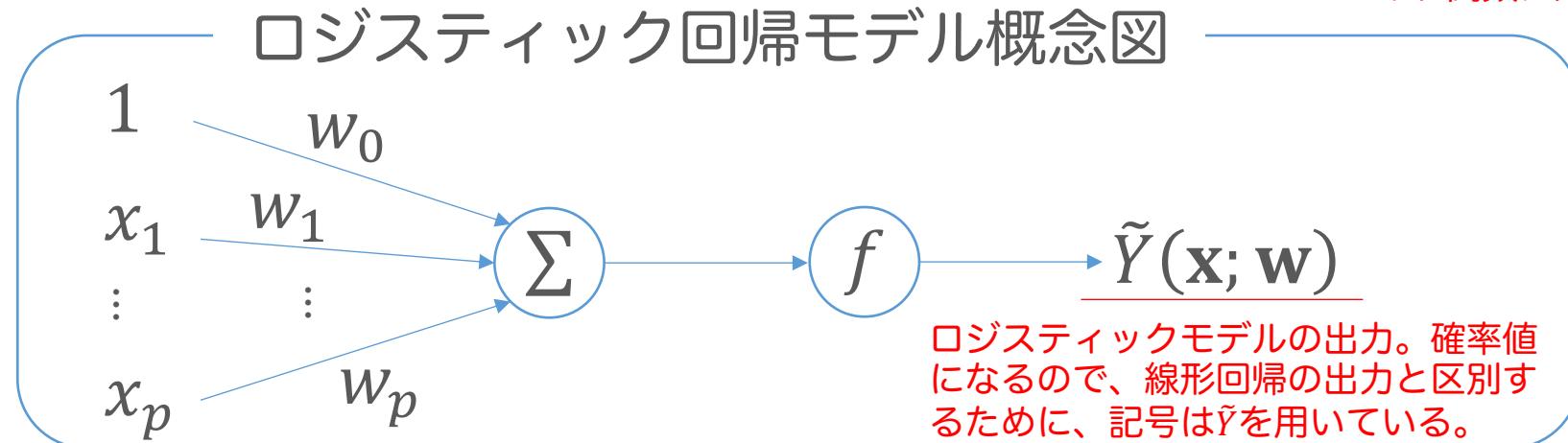


# シグモイド関数 | ロジスティック回帰モデル (Logistic Regression Model)

- 学習可能なパラメータを  $\mathbf{w} = (w_0, w_1, \dots, w_p)^T \in \mathbb{R}^{p+1}$ 、モデルへの入力を  $\mathbf{x} = (1, x_1, \dots, x_p)^T$ 、シグモイド関数を  $f(x)$ としたとき、ロジスティック回帰モデルの出力  $\tilde{Y}(\mathbf{x}; \mathbf{w})$  を次式で定義する

$$\tilde{Y}(\mathbf{x}; \mathbf{w}) = f(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\sum_{i=0}^p w_i x_i)}$$

線形回帰の出力をシグモイド関数に入力している



## 最尤推定 (Maximum Likelihood Estimation)

---

- ・ ロジスティック回帰モデルの出力 $\tilde{Y}(x; w)$ は、データ $x$ がエンジニアである確率を意味する
- ・ どのデータに対してもエンジニアである/ではないという確率 (=確信度)ができるだけ大きくなるように境界線を引きたい
- ・ できるだけ尤 (もっと) もらしい境界線のパラメータ $w$ を学習することを**最尤 (さいゆう) 推定**という

# 最尤推定 (Maximum Likelihood Estimation)

- モデルの評価指標として、全データに対する確率の対数値  
**(対数尤度；たいすうゆうど)  $\ln P$  を使用**
- 対数尤度を最大化するようにパラメータ  $w$  を学習**

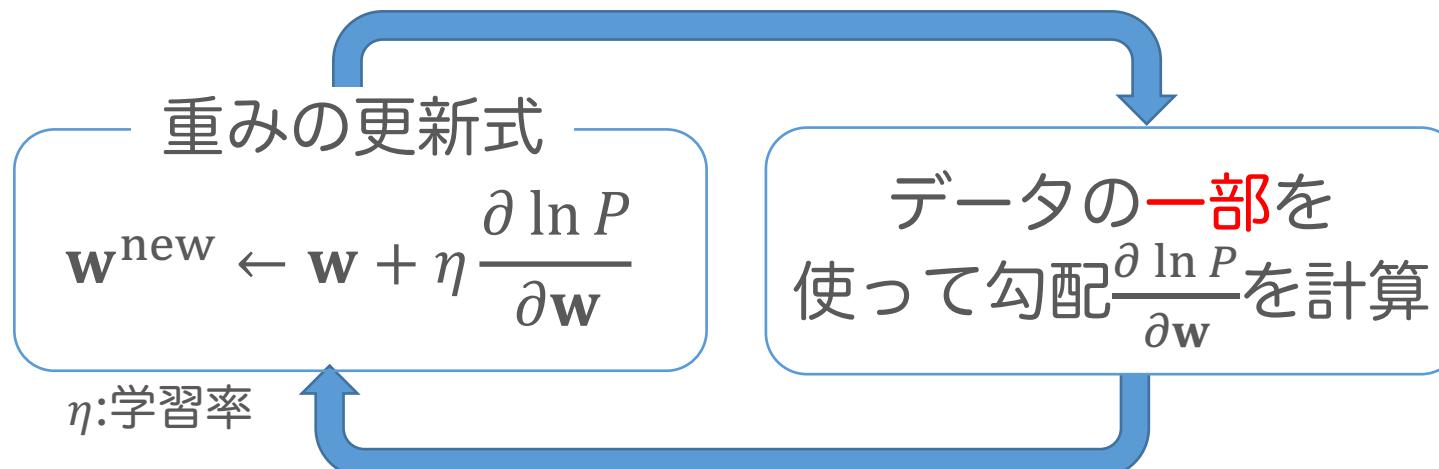
$$\ln P = \sum_{n=1}^N \left\{ y^{(n)} \ln \tilde{Y}(\mathbf{x}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \ln (1 - \tilde{Y}(\mathbf{x}^{(n)}; \mathbf{w})) \right\}$$

正解ラベルがエンジニアであるデータは  
そのデータがエンジニアであるという確信度  
が高いほど評価値が良い

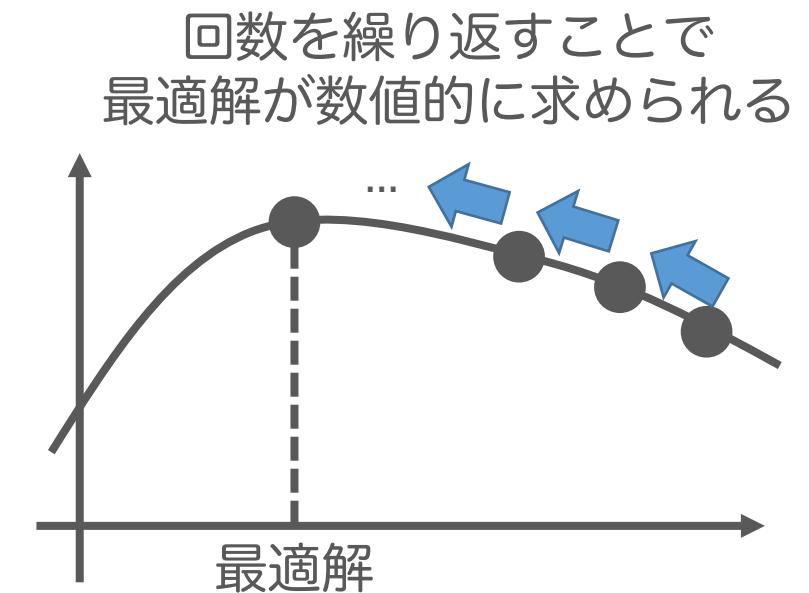
正解ラベルがエンジニアではないデータは  
そのデータがエンジニアではないという確信度が  
高いほど評価値が良い

# 確率的勾配法

- 対数尤度が最大となる条件式  $\frac{\partial \ln P}{\partial w} = 0$  を解いて、最適な重みを求めたいところだが解析解が求まらない…
- そこで繰り返し計算によって最適化な値を探索する
- 勾配  $\frac{\partial \ln P}{\partial w}$  は解析的に求まるので、**確率的勾配法**で最適化



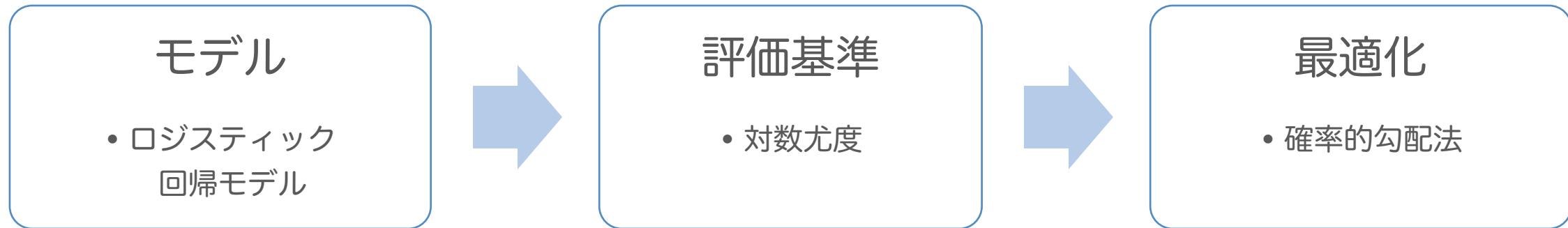
勾配の計算の詳細は  
『ITエンジニアのための機械学習理論入門』 p.164を参照



# ロジスティック回帰まとめ（モデル、評価基準、最適化の観点から）

---

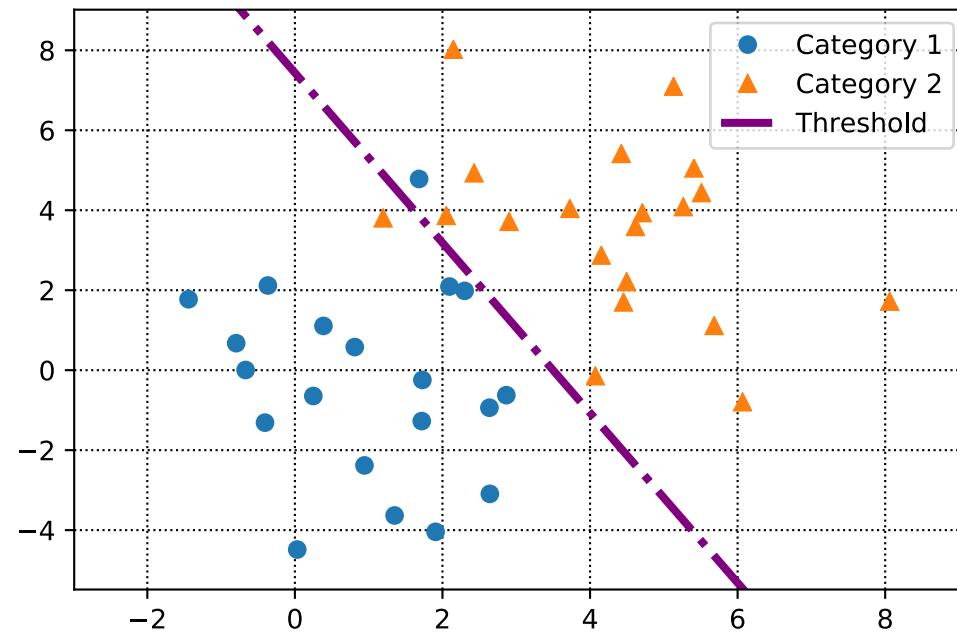
- ・ロジスティック回帰：回帰によってそのカテゴリである確率を予測する手法
- ・例）メールに含まれる単語を説明変数として、新たに着信したメールが迷惑メールかどうかを判別する



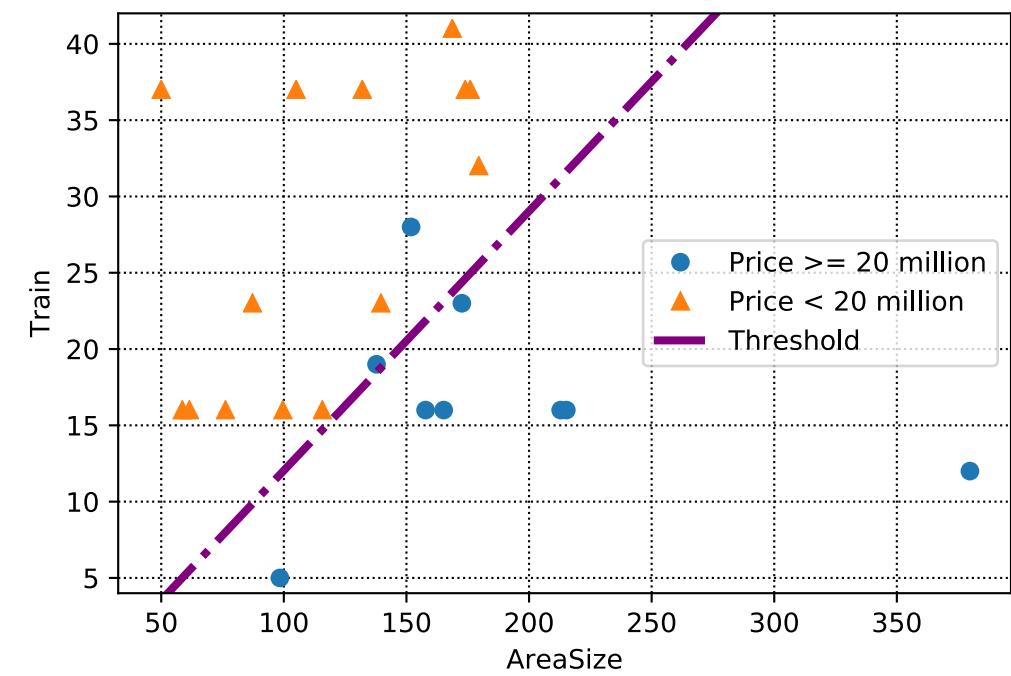
## [演習] ロジスティック回帰 (5\_logistic\_regression\_psedo\_data.ipynb)

- まずは簡単なデータでロジスティック回帰を試してみましょう
- エンジニア予測問題と住宅価格のカテゴリ予測問題でそれぞれモデルを作って確認してみましょう

エンジニア予測



住宅価格のカテゴリ予測



## [演習] ロジスティック回帰 (6\_logistic\_regression\_real\_data\_trainee.ipynb)

---

- ・説明変数が多いデータセットでロジスティック回帰を実装してみましょう
- ・住宅がリノベーションされているかどうかを予測するモデルを作ってみましょう

# 多変量モデルへの拡張

1. 多変量回帰
2. 多クラス分類（ソフトマックス関数と交差エントロピー）

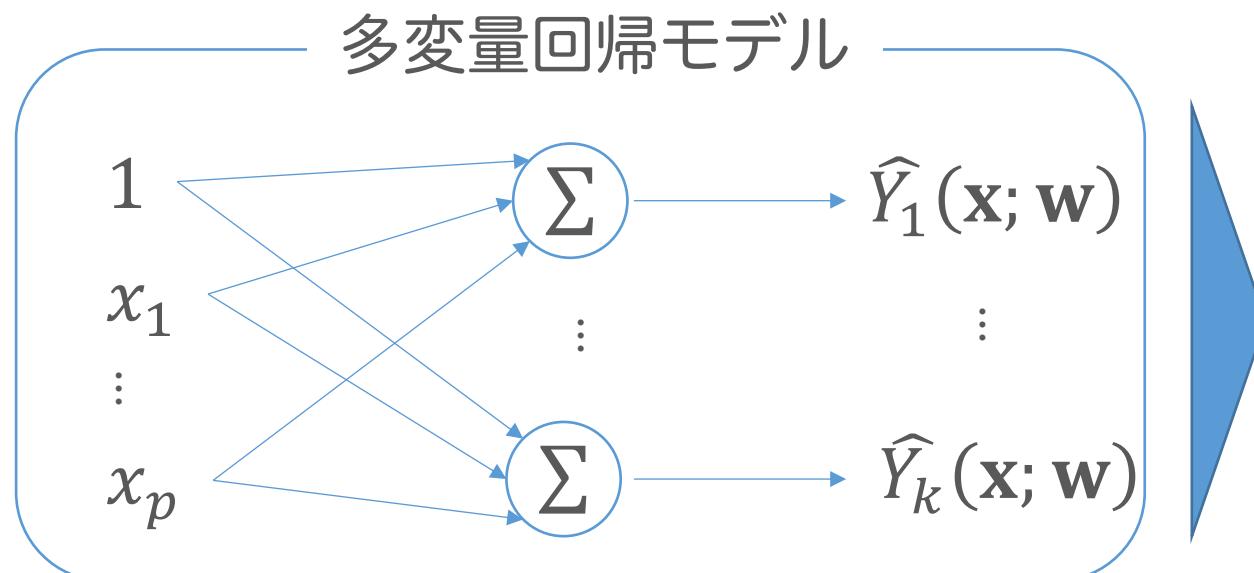
## 問題

---

- ・エンジニア、医者、弁護士、デザイナーを対象として、学生時代の理科および数学の偏差を調査した
- ・このデータから、ある学生が将来これら4つの職業になる確率をそれぞれ求めたい
- ・どのようなモデルをつくればいいか？

# 多変量回帰

- 今まででは目的変数が1次元のケースを取り扱ってきた
- しかし、実際には目的変数が多次元であることもありうる
- 線形回帰において出力が多次元のとき→**多変量回帰**

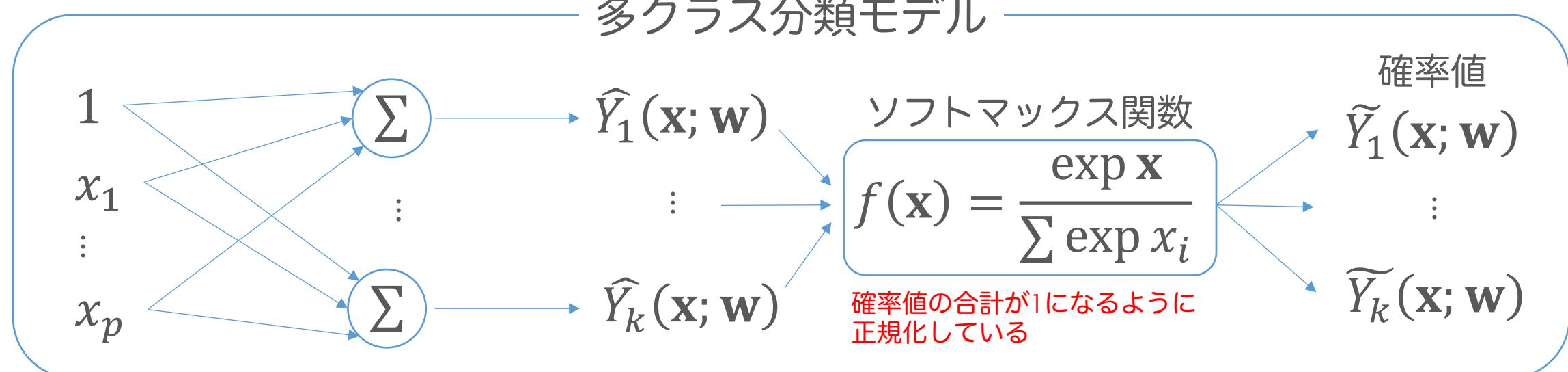


目的変数の各次元ごとに  
線形回帰モデルを置く

評価基準もそれぞれの次元に  
対する二乗誤差の和でOK

# 多クラス分類

- 分類問題のときは、各出力は「各カテゴリである確率」をそれぞれ意味する
  - 1つ目は犬である確率、2つ目は猫である確率、3つ目はパンダである確率…
- このような多クラス分類モデルのときは、シグモイド関数の代わりに  
**ソフトマックス関数**を用いて出力を確率値に変換する



## 多クラス分類

---

- 多クラス分類モデルにおける評価指標には**交差エントロピー**を用いることが多い
  - 深層学習でも一般的に用いられる
- 交差エントロピー=負の対数尤度

$$-\ln P = - \sum_{n=1}^N \sum_{i=1}^k \frac{y_i^{(n)} \ln \tilde{Y}_i(\mathbf{x}^{(n)}; \mathbf{w})}{y_i}$$

負の対数尤度  
=当てはまりの悪さを示す

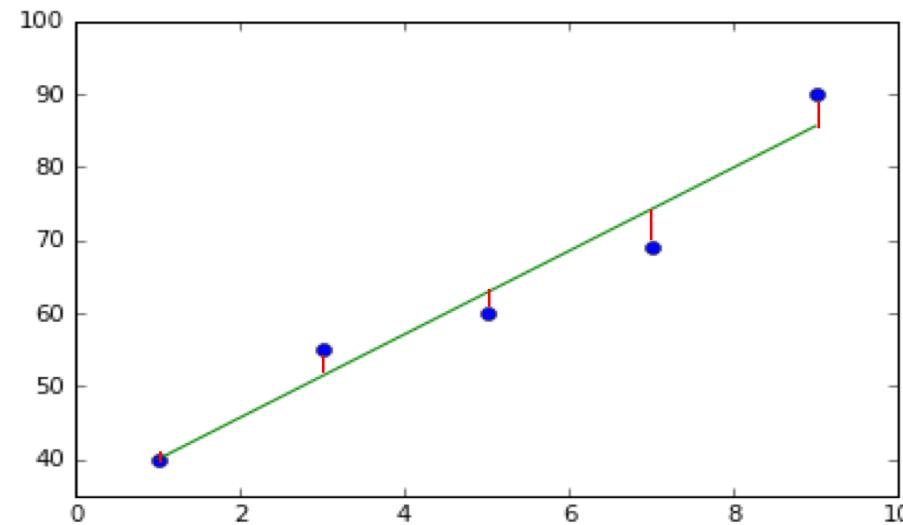
# モデルの評価指標

1. 回帰問題の評価指標
2. 分類問題の評価指標

# 回帰問題の評価指標 (MAE・RMSE・MSE)

- RMSE…予測と実際の差の二乗の平均の平方根
  - 平方根平均二乗和誤差とも呼ばれる
- MAE…予測と実際の差の絶対値の平均
  - 平均絶対値誤差とも呼ばれる
- MSE…予測と実際の差の二乗の平均
  - 平均二乗和誤差とも呼ばれる
- MAEは誤差の平均という点で解釈性がよい
  - 右図ならばMAEは「赤線の平均」

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (f_i - y_i)^2}$$
$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |f_i - y_i|$$
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (f_i - y_i)^2$$



# 分類問題の評価（混同行列/Accuracy/Recall/Precision/F1）

- 予測結果と真の結果でクロス集計をしたものを**混同行列**という
- Accuracy…正解率  $(TP+TN) / (TP+TN+FN+FP)$
- Recall…実際に正しいもののうち、  
正であると予測された割合  $(TP)/(TP+FN)$
- Precision…正と予測したもののうち、  
どれくらい正しかったか  $(TP)/(TP+FP)$

		真の結果	
		正 (Positive)	負 (Negative)
予測結果	正 (Positive)	TP個	FP個
	負 (Negative)	FN個	TN個

# 分類問題の評価（混同行列/Accuracy/Recall/Precision/F1）

---

- F1…RecallとPrecisionの調和平均  
両者のバランスと精度を評価した指標
- 目的によってどれで評価するかは変えるべき  
例) がんの診断→患者視点ならPrecision優先, 医師目線ならRecall優先

		真の結果	
		正 (Positive)	負 (Negative)
予測結果	正 (Positive)	TP個	FP個
	負 (Negative)	FN個	TN個

# [演習] モデルの評価指標 (7\_model\_evaluation.ipynb)

- 回帰問題、分類問題それぞれにおける評価指標の実装方法を確認してみましょう

```
# 値を予測  
y_pred = regr.predict(X)  
  
# MSEを計算  
mse = mean_squared_error(y, y_pred)  
print("MSE = %s"%round(mse,3))  
  
# MAEを計算  
mae = mean_absolute_error(y, y_pred)  
print("MAE = %s"%round(mae,3))  
  
# RMSEを計算  
rmse = np.sqrt(mse)  
print("RMSE = %s"%round(rmse, 3))
```

MSE = 14.885  
MAE = 3.057  
RMSE = 3.858

```
# ラベルを予測  
y_pred = clf.predict(X)  
  
# 正答率を計算  
accuracy = accuracy_score(y, y_pred)  
print('正答率 (Accuracy) = {:.3f}%'.format(100 * accuracy))  
  
# Precision, Recall, F1-scoreを計算  
precision, recall, f1_score, _ = precision_recall_fscore_support(y, y_pred)  
  
# カテゴリ「2000万以上」に関するPrecision, Recall, F1-scoreを表示  
print('適合率 (Precision) = {:.3f}%'.format(100 * precision[0]))  
print('再現率 (Recall) = {:.3f}%'.format(100 * recall[0]))  
print('F1値 (F1-score) = {:.3f}%'.format(100 * f1_score[0]))
```

正答率 (Accuracy) = 91.304%  
適合率 (Precision) = 87.500%  
再現率 (Recall) = 100.000%  
F1値 (F1-score) = 93.333%

## グループワーク（ロジスティック回帰）

---

- ある食品Aについて、ECサイトの登録情報をもとにその顧客が買うか買わないかを予測するロジスティック回帰モデルを作成した
- モデル作成に用いた顧客の登録情報：

年齢	性別	現住所の県名	類似商品 の購入回数	同じブランド商品 の購入回数
----	----	--------	---------------	-------------------

## グループワーク（ロジスティック回帰）

---

- 1. このモデルはどのように使えば利益を挙げられるでしょうか？
  - 利益が挙げられないという場合は、その問題点と改善方法を議論しましょう
- 2. 1の使い方に応じた、モデルの評価指標を議論してみましょう
- 3. 2で決めた指標について実製品 or サービスとして提供可能な基準(KPI)を設定してみましょう
- 4. 1~3の理由を説明できるように整理し、全体に向けて発表しましょう

# 『現場で使える機械学習・データ分析基礎講座』 を通しての課題

## 通し課題

---

- kaggleのデータセットを用いてモデルを構築し結果を公開することを本講座の通し課題とします
- 公開先は、kaggleのKernelsまたはGithubとします
  - Slackの所属チャンネルにipynbファイルを直接貼っても良いこととします
- 課題は、以下の2つから選んで下さい

課題	データセット名	問題設定
課題①	Kickstarter Projects	クラウドファンディングが成功するか(state)を予測
課題②	Car Fuel Consumption	100kmあたりのガソリン消費量(consume)を予測

<https://www.kaggle.com/>

The screenshot shows the main landing page for Kaggle Competitions. At the top, there's a navigation bar with links for 'Search kaggle', 'Competitions', 'Datasets', 'Kernels', 'Discussion', 'Jobs', and a 'Sign In' button. Below the navigation, a large banner reads 'Welcome to Kaggle Competitions' and 'Challenge yourself with real-world machine learning problems'. The page features three main sections with circular icons: 'New to Data Science?' (with a gear icon), 'Build a Model' (with a code editor icon), and 'Make a Submission' (with a trophy and podium icon). Each section has a brief description and a 'Learn more' or 'InClass' button. A small 'Dismiss' button is at the bottom right of the main content area.

General	InClass	Sort by	Grouped
18 Active Competitions			All Categories <input type="button" value="▼"/>
<input type="text" value="Search competitions"/> <input type="button" value="Search"/>			
	<b>Zillow Prize: Zillow's Home Value Prediction (Zestimate)</b> Can you improve the algorithm that changed the world of real estate? Featured · 3 days to go · housing, real estate	\$1,200,000 3,778 teams	
	<b>Mercari Price Suggestion Challenge</b> Can you automatically suggest product prices to online sellers? Featured · a month to go ·	\$100,000 1,411 teams	
	<b>Statoil/C-CORE Iceberg Classifier Challenge</b> Ship or iceberg, can you decide from space? Featured · 16 days to go · weather, shipping, image, binary classification	\$50,000 3,137 teams	
	<b>Toxic Comment Classification Challenge</b> Identify and classify toxic online comments Featured · a month to go · arguments, text	\$35,000 1,191 teams	
	<b>Corporación Favorita Grocery Sales Forecasting</b> Can you accurately predict sales for a large grocery chain? Featured · 8 days to go · food and drink, tabular, regression, forecasting	\$30,000 1,583 teams	
	<b>IEEE's Signal Processing Society - Camera Model Identification</b> Identify from which camera an image was taken Featured · a month to go · image	\$25,000 199 teams	
	<b>Recruit Restaurant Visitor Forecasting</b> Predict how many future visitors a restaurant will receive Featured · a month to go ·	\$25,000 1,265 teams	

# 課題① Kickstarter Projects

<https://www.kaggle.com/kemical/kickstarter-projects>

The screenshot shows the Kaggle dataset page for 'Kickstarter Projects'. At the top, there's a banner with the text 'Reviewed Dataset' and '370 voters'. Below the banner, the title 'Kickstarter Projects' is displayed in large green letters, with the subtitle 'More than 300,000 kickstarter projects'. It shows '300K PROJECTS' and was last updated '7 months ago (Version 7)' by 'Mickaël Mouillé'. The main navigation tabs are 'Data', 'Overview', 'Kernels', 'Discussion', and 'Activity'. A prominent blue button says 'New Kernel'. Below the tabs, there's a section titled 'Data (37 MB)' with an 'API' link and a download button labeled 'Download All'. This section includes 'Data Sources' (two CSV files: 'ks-projects-201612.csv' and 'ks-projects-201801.csv'), 'About this file' (which says 'No description yet'), and 'Columns' (a list of project attributes: ID, name, category, main\_category, currency, and deadline). A command line interface (CLI) command is shown above the download buttons: 'kaggle datasets download -d kemical/kickstarter-projects'.

課題①の目的は、あるクラウドファンディングが成功するか(state)を事前に予測するモデルを構築することです。使用する説明変数にご注意ください。

Kaggleにアカウントを作成するとCSVデータをダウンロードできるようになります

## 課題② Car Fuel Consumption

<https://www.kaggle.com/anderas/car-consume>

The screenshot shows the Kaggle website interface. At the top, there's a navigation bar with links for 'Competitions', 'Datasets', 'Kernels', 'Discussion', 'Learn', and 'Sign In'. Below the navigation is a banner for the 'Car Fuel Consumption' dataset, which is described as a 'Reviewed Dataset' by Andreas Wagener, updated 3 months ago (Version 5). The banner features a photograph of a gas station at night. A '12 voters' badge is visible in the top right corner of the banner area. Below the banner, there are tabs for 'Data', 'Overview', 'Kernels', 'Discussion', and 'Activity'. The 'Data' tab is selected. To the right of the tabs are buttons for 'Download (225 KB)' and 'New Kernel'. The main content area is titled 'Data (225 KB)' and contains sections for 'Data Sources', 'About this file', and 'Columns'. The 'Data Sources' section lists three files: 'measurements.csv' (388 x 12), 'gas\_station\_orig.jpg', and 'measurements2.xlsx' (388 x 12). The 'About this file' section provides a brief description of the dataset, stating it was taken from a notebook and can be used to determine the effect of weather, speed, or gas type on car consumption. The 'Columns' section lists the variables: # distance, # consume, # speed, A temp\_inside, # temp\_outside, A specials, A gas\_type, # AC, and # rain. Above the 'Columns' section, there's a command-line interface (CLI) input field showing the command 'kaggle datasets download -d anderas/car-consume'. To the right of the CLI, there are buttons for 'Download All' and a refresh icon.

Kaggleにアカウントを作成すると  
CSVデータをダウンロード  
できるようになります

# 通し課題で具体的にやること と DAY1の宿題

---

1. 自分が取り組む通し課題を1つ選択する
  - Kaggleアカウントを作成し、該当課題のデータをダウンロードする
2. 目的変数と説明変数の関係を確認するためのグラフを作成する（ここからはNotebook上の作業です）
3. 目的変数を説明するのに有効そうな説明変数を見つける
4. DAY1で学んだアルゴリズムを利用する
  - 回帰の場合は線形回帰、分類の場合はロジスティック回帰
  - 質的変数が扱えないアルゴリズムを使う場合は、ダミー変数に置き換える
5. 予測精度または識別精度を確認する
  - 回帰問題の場合は、MSE、RMSE、MAEを求める
  - 分類問題の場合は、混同行列を作成し、Accuracy、Recall、Precisionを求める
6. できたところまでをNotebookでまとめ、KernelsまたはGithubで公開する
  - 公開方法がわからない方は、ipynbファイルを所属チャンネルに貼る

# 通し課題で具体的にやること と DAY2、3の宿題

---

## 8. DAY2、3で学んだことの取り組み

- 交差検証、ホールドアウト法などで汎化性能を確認する
- 欠測値と異常値を確認し、適切に処理する
- DAY2、3で学んだアルゴリズムを利用してモデルをつくり、DAY1宿題提出時の精度と比較する
- 交差検証によるパラメータチューニングを行う
- パラメータチューニング後のモデルによって、精度および結果の評価を行う
- その他、精度の向上ができるような処理に取り組み、精度を上げる
- できたところまでをNotebookでまとめ、宿題として提出する
  - 前回から取り組んだ内容・工夫、精度がどのように変化したかのコメントをNotebookに含めること
  - 15分程度、受講者同士で通し課題の進捗を見せ合う時間を設けます

## 9. DAY4では、DAY3宿題の提出ファイルを元に、最終発表を実施いただく

## 宿題の提出について

---

1. 次回講義日の2日前の17時までに提出をお願いします
2. 提出方法
  1. 検討結果をNotebookにまとめ、KernelsまたはGithubで公開し、Slackの所属チャンネルにそのURLリンクを貼ってください
  2. 上記の操作方法がわからない方は、ipynbファイルをSlackの所属チャンネルに直接ファイルを貼っても良いこととします
  3. 提出したものに対して、講師がコメントを差し上げます
3. 投稿するファイル名は、DayX\_work\_お名前.ipynbでお願いします
  - DAY1の宿題の場合『Day1\_work\_鈴木一郎.ipynb』という形です

# 機械学習モデルを作る前に、グラフを描きましょう

---

- ・より良い機械学習モデルを作るためには、グラフを描くことが重要です
- ・グラフを描くと、説明変数と目的変数の関係性がみえてきます
  - ・例) 家の価格は新築ほど高い
  - ・例) 色が鮮やかならば毒キノコである確率が高い
- ・グラフを描くと、欠損値や異常値に気づけます
- ・グラフを描くと、機械学習モデルの動作を検証することができます
  - ・機械学習モデルがグラフと異なる傾向を出力したとき、すぐにおかしいと気付ける

# Any Questions ?

# Appendix

# 機械学習の文献を読む際に最低限覚えておきたい数学記号

---

# 機械学習の文献を読む際に最低限覚えておきたい数学記号

## ギリシャ文字

ギリシャ文字の一覧を以下に示す。頻出度は作者の独自の判断で定めたものである。

頻出度	大文字	小文字	読み方	使用例
高	A	$\alpha$	アルファ	係数
高	B	$\beta$	ベータ	係数
高	$\Gamma$	$\gamma$	ガンマ	係数, $\Gamma$ 関数
中	$\Delta$	$\delta$	デルタ	微小区間
高	E	$\epsilon$	イプシロン	係数, 微小なもの, 誤差項
低	Z	$\zeta$	ゼータ	
中	H	$\eta$	イータ	係数
高	$\Theta$	$\theta$	シータ	角度, 推定すべきパラメータ
低	I	$\iota$	イオタ	
低	K	$\kappa$	カッパ	係数
高	$\Lambda$	$\lambda$	ラムダ	係数
高	M	$\mu$	ミュー	平均

頻出度	大文字	小文字	読み方	使用例
低	N	$\nu$	ニュー	
低	$\Xi$	$\xi$	グザイ	
低	O	$\circ$	オミクロン	
高	$\Pi$	$\pi$	パイ	確率, $\Pi$ は総乗記号
高	P	$\rho$	ロー	係数
高	$\Sigma$	$\sigma$	シグマ	標準偏差, $\Sigma$ は総和記号や共分散行列など
中	T	$\tau$	タウ	時間
低	$\Upsilon$	$\upsilon$	ユプシロン	
高	$\Phi$	$\phi$	ファイ	係数
中	X	$\chi$	カイ	$\chi^2$ 分布
高	$\Psi$	$\psi$	ブサイ	係数
低	$\Omega$	$\omega$	オメガ	角速度

# 機械学習の文献を読む際に最低限覚えておきたい数学記号

## 記号

## 集合記号

記号 意味 用例

$\hat{\cdot}$  推定値  $\hat{y}, \hat{x}$

$\bar{\cdot}$  平均値  $\bar{y}, \bar{x}$

## 不等号

記号 意味

$\leq$  小なりイコール

$\leq$  小なりイコール,  $\leq$ と同じ

$\geq$  大なりイコール

$\geq$  大なりイコール,  $\geq$ と同じ

$\approx$  近似イコール

記号 意味

{ } 集合

$x_i = \{x_1, x_2, x_3, \dots, x_{100}\}$

$\in$  属する

$x \in \{1, 2, 3, \dots, 100\}$

$\mathbf{R}$  実数

$x \in \mathbf{R}$

$\mathbf{Z}$  整数

$x \in \mathbf{Z}$

$\equiv$  定義

$A \equiv B, A$ を $B$ と定義する

$\cup$  または

$A \cup B$

$\cap$  かつ

$A \cap B$

---

---

# 機械学習の文献を読む際に最低限覚えておきたい数学記号

## 確率

記号	意味	用例
$N(,)$	正規分布	$N(\mu, \sigma^2)$
$\sim$	ある確率変数がある確率分布に従う	$\epsilon \sim N(\mu, \sigma^2)$
$\propto$	比例	$p \propto q$ , pとqは比例している
$p(,)$	同時確率	$p(A, B)$ , AとBが同時に生じる確率
$p( )$	条件付き確率	$p(A   B)$ , Bが生じたという条件つきのもとでのAが生じる確率

# 機械学習の文献を読む際に最低限覚えておきたい数学記号

## 表記に関する慣例

- $\times$ は省略されることが多い  
例、 $a \times x = ax$
- カッコは省略されることがある  
 $\log(x) = \log x$
- 底が省略された対数は、自然対数  
 $\log(x) = \log_e(x) = \ln(x)$   
ln:読み方はエルエヌ,ロンなど
- 総和記号 $\Sigma$ は、記号が省略されることがある

$$\sum_{i=0}^k P_i = \sum_i P_i$$

## 閉区間、開区間

[a,b] : aとbを含む区間、閉区間

(a,b) : a～bの区間だが、aとbを含まない、開区間

[a,b):a～bの区間だが、aは含み、bは含まない

## 絶対値

$$|x|$$

## ノルム

1ノルム

$$\|x\|_1 = |x_1| + |x_2| + |x_3| + \dots + |x_n|$$

2ノルム(ユークリッドノルム。単にノルムという場合はこれ。)

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2}$$

$\|x\|_2$ を $\|x\|$ と略して表現されることもある

2ノルムの2乗

$$\|x\|_2^2$$

## 大文字シグマの使われ方

$\Sigma$ は、総和記号として以外にも共分散行列を意味する記号としても使われる

# 最尤推定法

## 最尤推定法

---

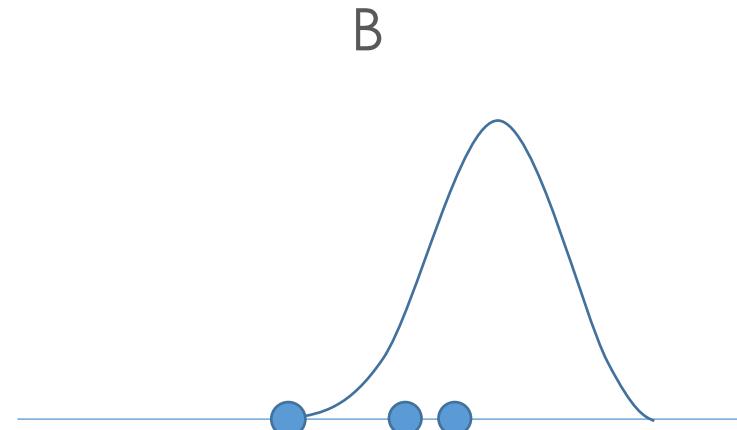
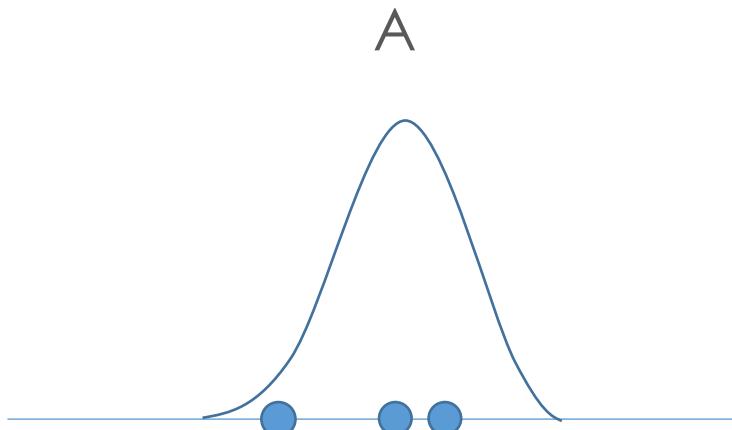
- ・以下のデータが観測されたとする。
- ・このデータに最もよく当てはまる正規分布を求めたい。



# 最尤推定法

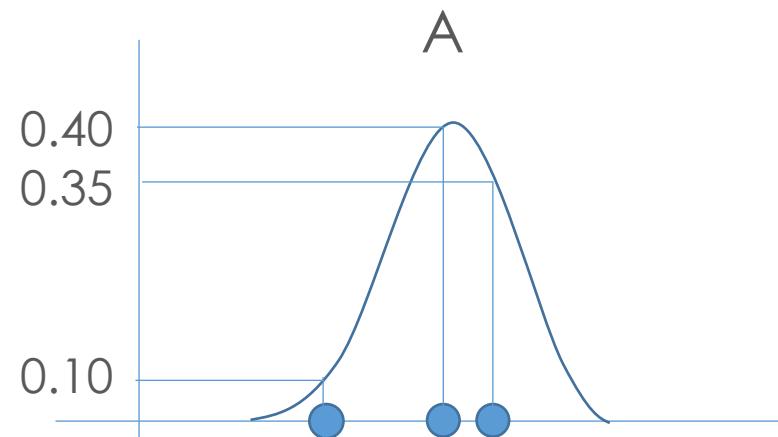
---

- 以下のどちらの正規分布の方がデータによく当てはまっているか？



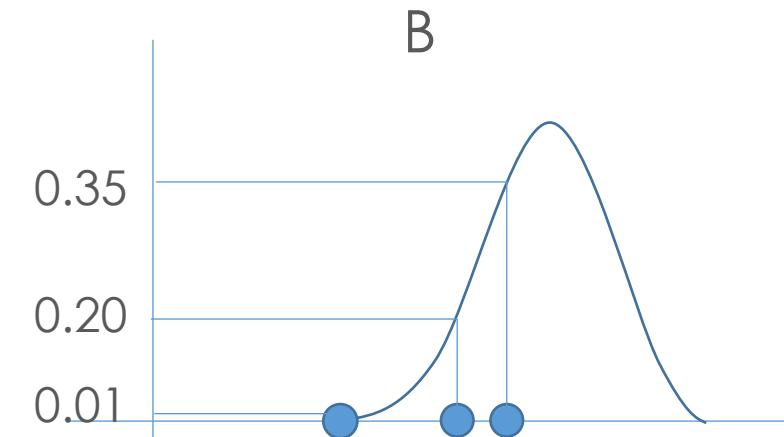
# 最尤推定法

- 正規分布A、正規分布Bのそれぞれについて、尤度を算出する。
- 尤度とは、確率(密度)を掛け算した値。



$$\text{尤度} = 0.40 \times 0.35 \times 0.10 = 0.014$$

&gt;



$$\text{尤度} = 0.35 \times 0.20 \times 0.01 = 0.0007$$

Aの正規分布の方が尤度が大きいので、Aの方がよく当てはまっていると判断する。  
これが最尤法。実際には、解析的に最もよく当てはまる正規分布を求める。

## 情報量、エントロピー、交差エントロピー

---

- 情報量とは、**事象の驚き度合い**を数値化した指標
- 到底起きそうもない事象が起きたことを知った→情報量が大きい
  - 例) 砂漠に雨が降った→到底起きそうにない出来事なので情報量 大
- いつでも起きそうな事象が起きたことを知った→情報量は小さい
  - 例) 東京に雨が降った→いつでも起きそうな出来事なので情報量 小

## 【例】

例えば、東京に雨が降ったという情報よりも、砂漠に雨が降った(確率が低い事象)という情報の方が驚くはずである。

ある日に、東京に雨が降る確率  $p(x)$  を0.3とすれば、その情報量は、

$$I(x) = -\log_2 p(0.3) = 1.737$$

となり、ある日に、砂漠に雨が降る確率  $p(x)$  を0.01とすれば、その情報量は、

$$I(x) = -\log_2 p(0.01) = 6.644$$

となる。

$$I(x) = -\log_2 p(x)$$

$I(x)$  : ある事象  $x$  の情報量 (単位は bits)

$p(x)$  : ある事象  $x$  が起きる確率

# エントロピー（平均情報量）

---

- ・エントロピーとは、情報量を平均化したもので**事象のばらつき具合**を表す
- ・大半は同じ事象しか起きない→エントロピーが小さい
  - ・例) 砂漠の天気→ほとんどが晴れなのでエントロピー 小
- ・観測するたびに事象がばらついている→エントロピーが大きい
  - ・例) 東京の天気→晴れだったり雨だったりとバラバラなのでエントロピー 大

## 【例】

$$H = \sum p(x)I(x) = -\sum p(x) \log_2 p(x)$$

$H$ ：エントロピー

$I(x)$ ：ある事象 $x$ に対する情報量

例えば、ある日の東京の天気の確率が、 $p(\text{晴れ})=0.5$ 、 $p(\text{雨})=0.3$ 、 $p(\text{曇り})=0.2$ とすると、そのエントロピーは、

$$H = -\sum p(x) \log_2 p(x) = 1.485$$

となり、ある日の砂漠の天気の確率が、 $p(\text{晴れ})=0.9$ 、 $p(\text{雨})=0.01$ 、 $p(\text{曇り})=0.09$ とすると、そのエントロピーは、

$$H = -\sum p(x) \log_2 p(x) = 0.516$$

となる。

東京の天気は、確率変数 $x$ がばらついており、砂漠の天気に比べ予測が難しいと言える。

## 交差エントロピー

---

- ・ 交差エントロピーは2つの事柄について、その中で起こりうる事象のばらつきの違いを数値化したもの
- ・ 事象のばらつき方が類似している→交差エントロピーは小さい
  - ・ 例) 横浜の天気と東京の天気 → 天気のばらつき方は似ているため  
交差エントロピー 小
- ・ 事象のばらつき方が全く異なる→交差エントロピーは大きい
  - ・ 例) 砂漠の天気と東京の天気 → 天気のばらつき方が大きく異なるため  
交差エントロピー 大
- ・ 機械学習では予測値の分布と正解ラベルの分布の違いを測るために用いる

## ソフトマックス関数

---

## ソフトマックス関数

- ・ ソフトマックス関数は、複数の入力を正規化し、合計値が1になるようにする関数
- ・ 正規化する前に、指数関数(exp)を計算している。

$$y(k) = \frac{\exp(a_k)}{\sum_{i=1}^n (\exp(a_i))}$$

$k$ :出力ノードの番号