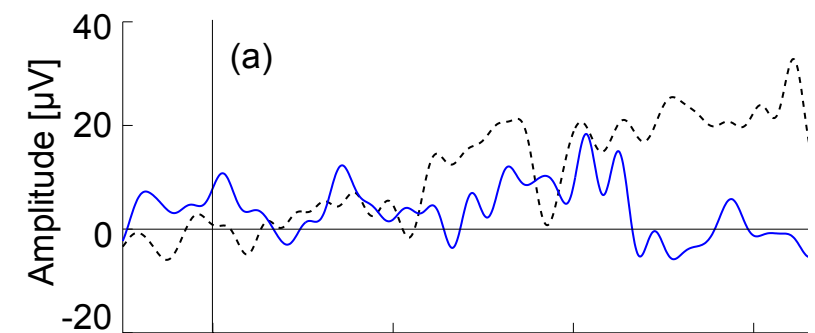
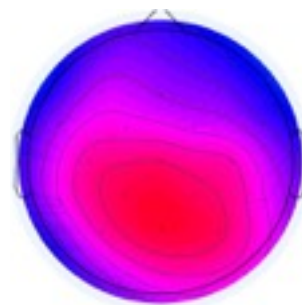


機械学習セミナー

@Life is tech! 大阪オフィス

自己紹介

- 名前：真木勇人 まきはやと (twitter: @mkhyt)
- 所属：奈良先端科学技術大学院大学 (NAIST)
 - ▶ 情報科学研究科 D1 知能コミュニケーション研究室
- 専門
 - ▶ 信号処理、機械学習
 - ▶ 研究テーマ：機械学習を利用した脳情報の分解・解読



今日の目的

- 将来機械学習を使うかもしれないエンジニアに、機械学習の原理、使用上・ビジネス上のポイントを知ってもらう。

機械学習

とは

機械学習 Machine Learning

- いわゆる「人工知能」の基盤技術

音声認識

顔画像検出



手書き文字認識

対話システム (Siri)

amazon

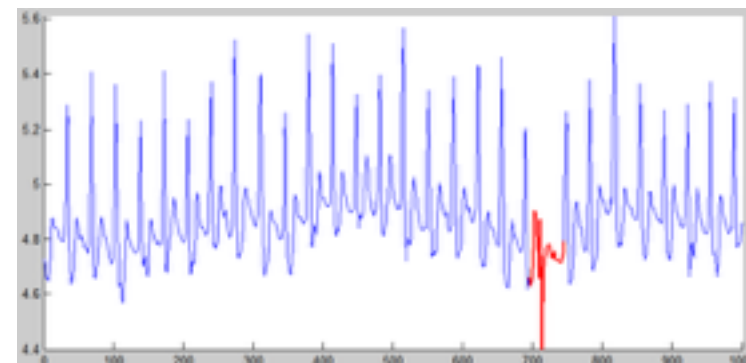
レコメンドシステム



機械翻訳



ユーザークラスタリング

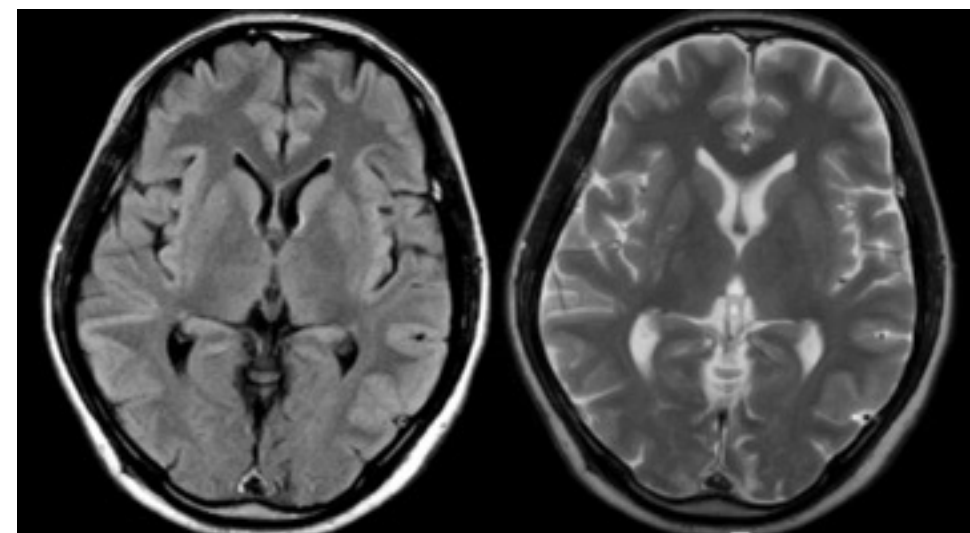


異常検出

将来的な応用



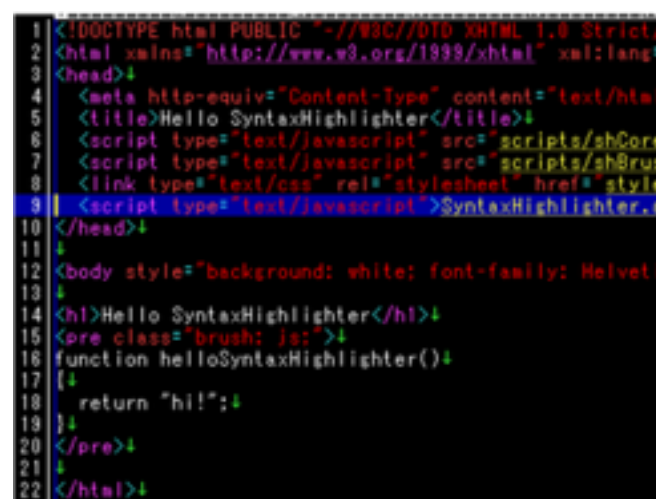
自動運転



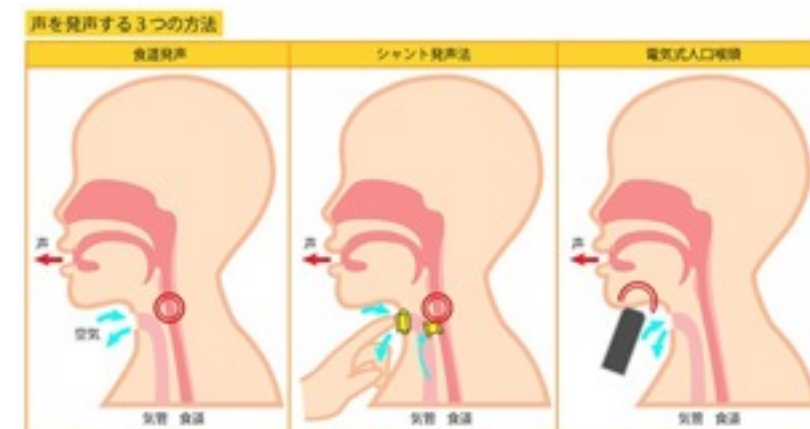
医療画像診断



同時音声翻訳



ソースコード生成



リアルタイム声質変換

機械に

データから法則性(ルール)を

自動的に

学習(発見)させる方法

機械学習の

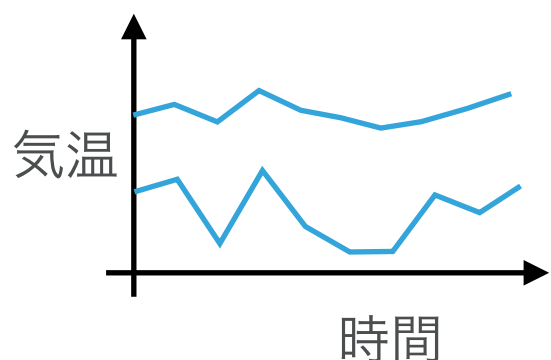
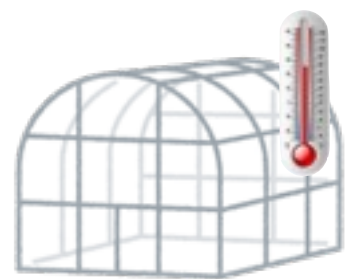
モチベーション

例：ビニールハウスの害虫発生予測

「データサイエンティスト養成読本機械学習入門編」（技術評論社）から改変して引用

- ビニールハウスの内外に温度計・湿度計が1つずつ設置
- 温度と湿度がある条件を満たすと、ビニールハウス内に害虫が発生

生データ



特徴抽出



特徴ベクトル

$X =$

外気温平均
内気温平均
外湿度平均
内湿度平均
当日の季節

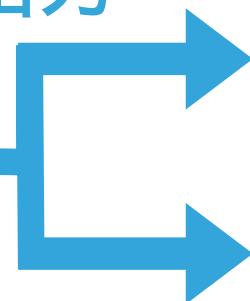
X

入力



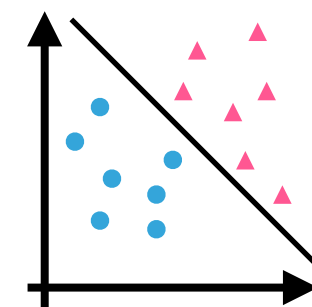
予測モデル
(ルール)

出力



アラート飛ばす

アラート飛ばさない



ルールをどうやって決める？

人手でルールを決めてみよう

- 過去のデータを眺めてみる（架空のデータ）

	外気温(°C)	内気温(°C)	外湿度(%)	内湿度(%)	季節	害虫発生
					夏=0, 冬=1	なし=0, あり=1
data1	33.5	37.1	70.2	72.4	0	1
data2	35.4	41.5	53.2	55.8	0	0
data3	31.8	35.4	63.3	62.0	0	1
data4	24.7	28.0	68.9	70.0	0	0
data5	10.6	25.2	61.1	63.8	1	1
data6	5.1	22.1	44.5	52.9	1	0
data7	6.3	20.4	70.7	75.3	1	0
data8	12.5	23.6	62.7	77.9	1	1

「内気温が30℃以上」ならアラート？

「夏かつ内気温30℃以上」または「冬かつ内気温25℃以上」ならアラート??

「夏かつ内気温30℃以上かつ内湿度60%以上」または「冬かつ内気温25℃以上かつ内湿度%60以上」ならアラート???

人手でルールを決めてみよう

- 過去のデータを眺めてみる（架空のデータ）

	外気温(°C)	内気温(°C)	外湿度(%)	内湿度(%)	季節 夏=0, 冬=1	害虫発生 なし=0, あり=1
data1						1
data2						0
data3						1
data4						0
data5	10.6	25.2	61.1	63.8	1	1
data6						0
data7						0
data8						1

高次元かつ大量のデータに対して人手で
法則性（ルール）を構築するのは困難

ルールの構築を自動化するのが機械学習
(機械にルールを発見させる)

「内気温が

「夏かつ内気温30°C以上」または「冬かつ内気温25°C以上」ならアラート??

「夏かつ内気温30°C以上かつ内湿度60%以上」または「冬かつ内気温25°C以上かつ内湿度%60以上」ならアラート???

機械学習の手法

- **教師あり学習** Supervised Learning

- ▶ 分類 Classification

- ▶ 回帰 Regression

- **教師なし学習** Unsupervised Learning

- ▶ クラスタリング Clustering

- ▶ 次元削減 Dimensionality Reduction

- ▶ 異常検出 Anomaly Detection

教師あり学習

x を入力して、 y を予測する

$$y = f(x)$$

- **トレーニングデータ**（**過去**のデータ）使って、関数 f （予測モデル）を推定する問題（関数近似問題）
- **トレーニングデータ**: **特徴ベクトル**と**正解ラベル**の事例セット

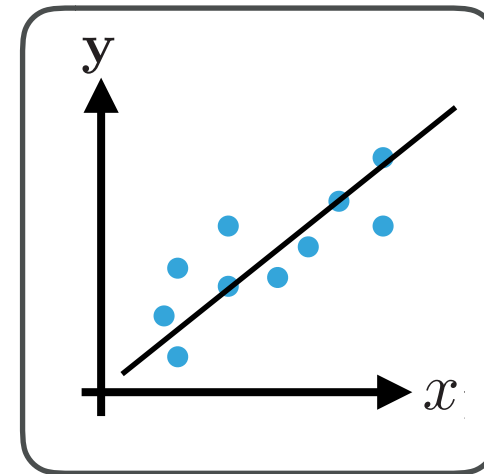
$x = (x_1, x_2) = (\text{気温}, \text{湿度})$, $y = \text{害虫発生あり or 発生なし}$

Day1	$x = (30, 70)$	$y = \text{あり}$
Day2	$x = (24, 65)$	$y = \text{なし}$
⋮	⋮	⋮

回帰と分類

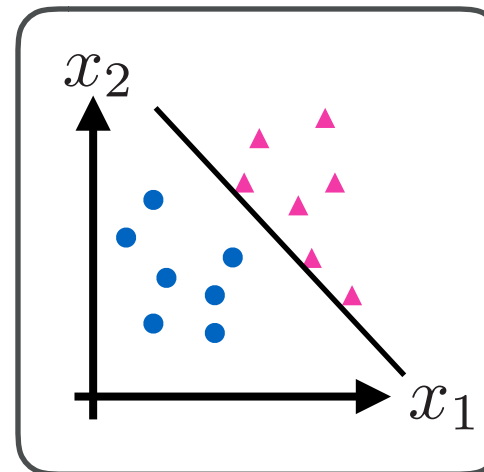
● 回帰

- ▶ 予測値が**数値**である問題
 - ▶ 例1) 気温からテーマパークの来場者数を予想する
 - ▶ 例2) 年齢と喫煙本数から残りの寿命を予想する



● 分類

- ▶ 予測値が**クラス**である問題
 - ▶ 例1) 単語からスパムメールを見分ける
 - ▶ 例2) 手書きの文字を認識する



回帰モデルの学習と予測

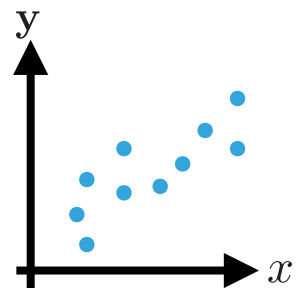


- 例：非雇用率から犯罪発生率を予想する \mathbf{X} =(非雇用率)

y =(犯罪発生率)

トレーニングデータ

	特徴ベクトル	正解ラベル
都市1	$\mathbf{x}^{(tr1)}$	$y^{(tr1)}$
都市2	$\mathbf{x}^{(tr2)}$	$y^{(tr2)}$
	\vdots	\vdots



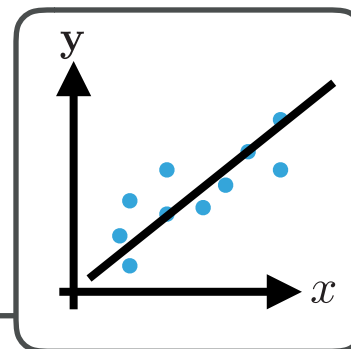
学習アルゴリズム

リッジ回帰
ニューラルネットなど

予測モデル f を推定

$$f(\mathbf{x}) = \underbrace{w_1}_{\text{推定}} \mathbf{x} + \underbrace{w_0}_{\text{推定}}$$

推定



予測モデル

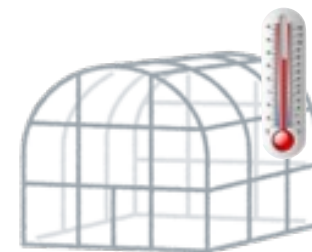
未知のデータ

\mathbf{X}_{new}

予測値

$$y_{\text{new}} = f(\mathbf{x}_{\text{new}})$$

分類モデルの学習と予測

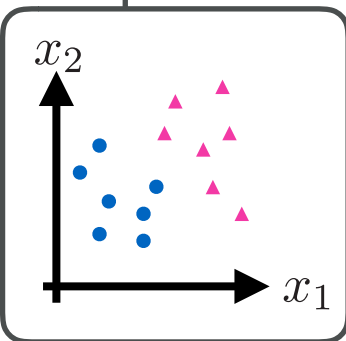


- 例：気温と湿度から害虫の発生を予測する $\mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \text{温度} \\ \text{湿度} \end{pmatrix}$

$y = \text{あり or なし}$

トレーニングデータ

	特徴ベクトル	正解ラベル
Day1	$\mathbf{x}^{(\text{tr1})}$	$y^{(\text{tr1})}$
Day2	$\mathbf{x}^{(\text{tr2})}$	$y^{(\text{tr2})}$
	\vdots	\vdots



学習アルゴリズム

パーセプトロン
SVMなど

予測モデル f を推定

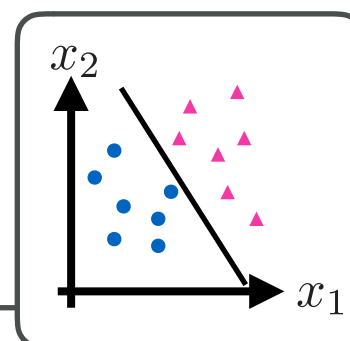
$$f(\mathbf{x}) = \underbrace{w_1 x_1}_{\text{推定}} + \underbrace{w_2 x_2}_{\text{推定}} + \underbrace{w_0}_{\text{推定}}$$

推定

予測モデル

未知のデータ

\mathbf{X}_{new}



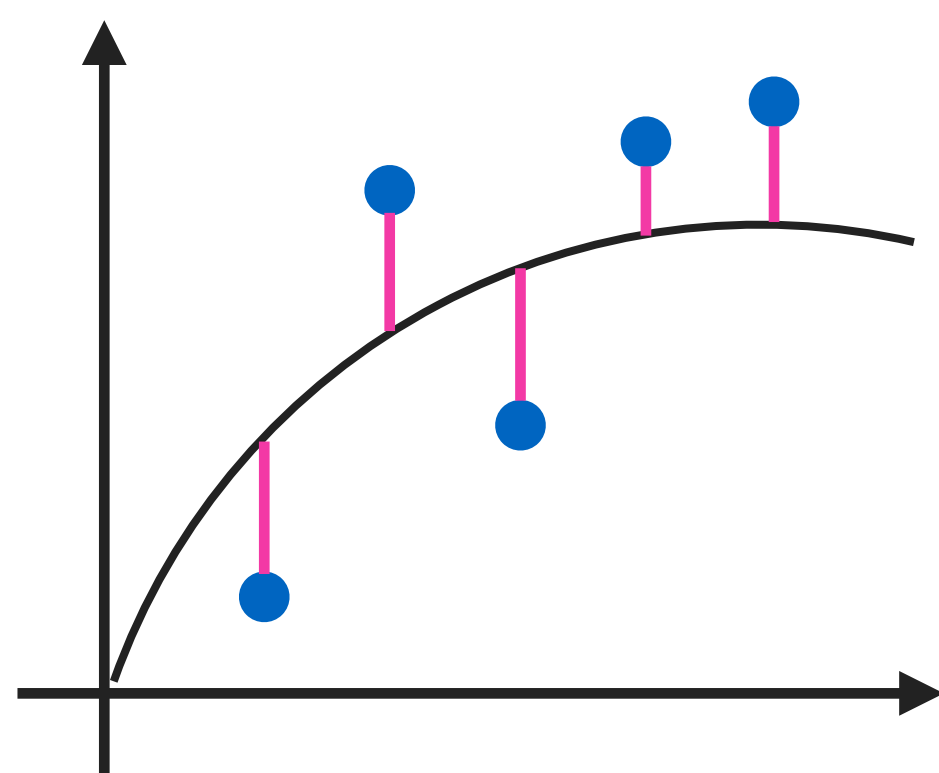
予測値

y_{new}

$$= \begin{cases} \text{あり} & (f(\mathbf{x}_{\text{new}}) > 0) \\ \text{なし} & (f(\mathbf{x}_{\text{new}}) < 0) \end{cases}$$

最小二乗学習

- 多くの機械学習アルゴリズムの原型



— : 予測モデル
● : 正解ラベルの値
— : 誤差

▶ 2乗誤差関数を最小化

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

前处理

正規化

- 年収と年齢から、残りの寿命を予想する

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \text{年収} \\ \text{年齢} \end{pmatrix} \begin{array}{l} \longleftarrow 200\text{万} \sim 3000\text{万} \\ \longleftarrow 18\text{歳} \sim 80\text{歳} \end{array}$$

$$\text{予測モデル: } f(\mathbf{x}) = w_1 x_1 + w_2 x_2 + w_0$$

年収が相対的に大きく影響してしまう

- 平均0、分散1になるように、特徴量ごとに正規化

$$x'_1 = \frac{x_1 - \mu}{\sigma} \quad \begin{array}{l} \mu : x_1 \text{の平均} \\ \sigma : x_1 \text{の標準偏差} \end{array}$$

ダミー変数

- 年収と居住地から、残りの寿命を予想する

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \text{年収} \\ \text{居住地} \end{pmatrix} \begin{array}{l} \longleftarrow 200\text{万} \sim 3000\text{万} \\ \longleftarrow \text{関東、関西、中部} \end{array} \begin{array}{l} \text{数値} \\ \text{カテゴリ} \end{array}$$

$$\text{予測モデル: } f(\mathbf{x}) = w_1 x_1 + w_2 x_2 + w_0$$

- カテゴリ変数を扱えるようにダミー変数を導入

$$f(\mathbf{x}) = w_1 x_1 + w_2 \underline{x_{\text{関東}}} + w_3 \underline{x_{\text{関西}}} + w_4 \underline{x_{\text{中部}}} + w_0$$

該当するところは1、他は0

機械学習を
使いこなす
ために

機械学習のメリット

●メリット

- ▶ アルゴリズムが汎用的、様々な問題に適用可能
- ▶ (うまく学習すれば)人間を上回る精度・速度を実現可能
- ▶ 人間には扱いきれない高次元・大量なデータを取り扱い可能
- ▶ (場合により)コスト削減

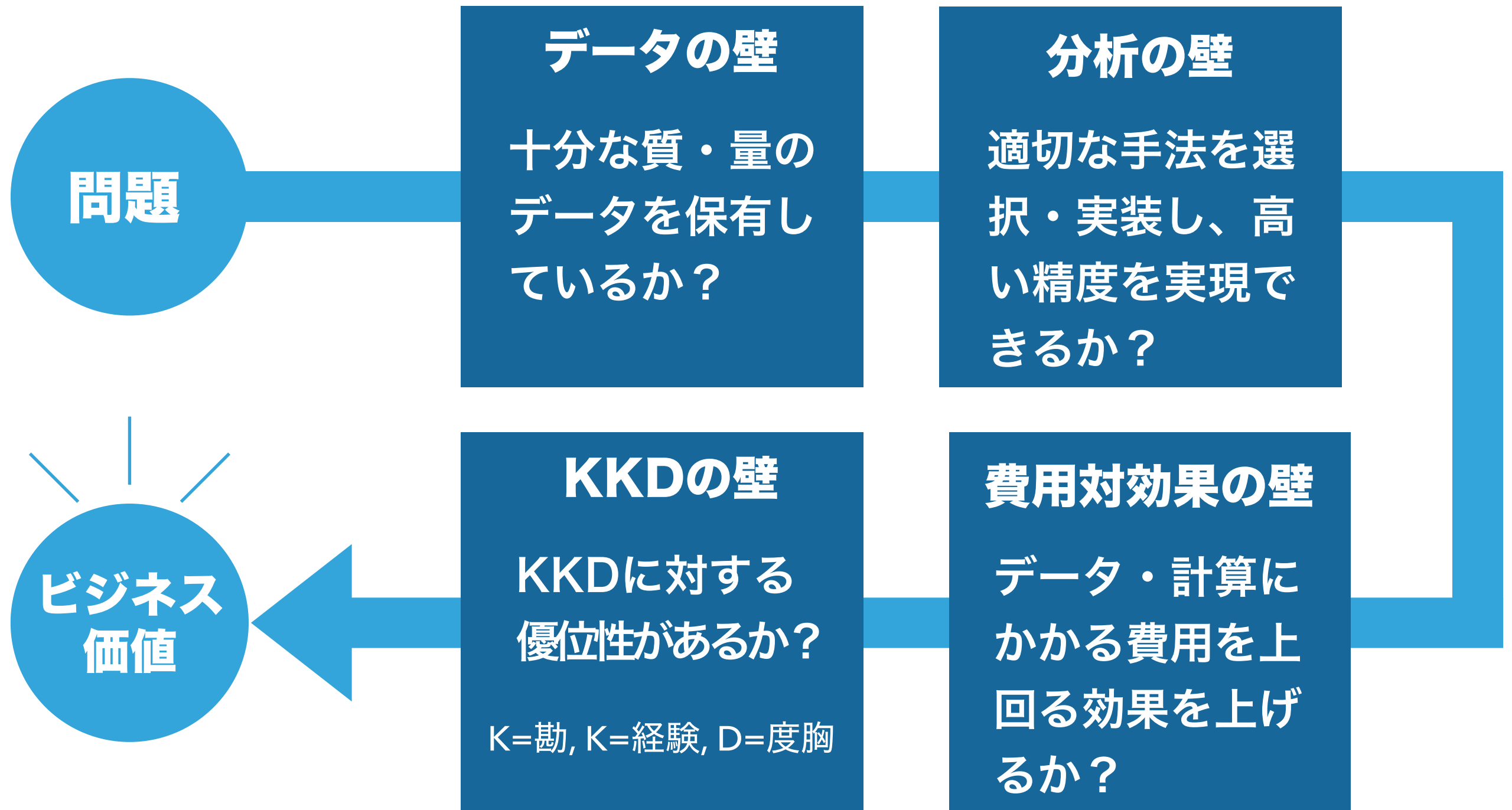
機械学習のデメリット

●デメリット

- ▶ 大量かつ良質なデータが必要
 - ▶ 欠損値、フォーマット不揃い、網羅性
 - ▶ データ前処理ニスト？
- ▶ 計算に長時間または豊富なマシンパワーが必要な場合がある
- ▶ 結果の解釈が容易でない場合がある
- ▶ 特徴量の選択、ハイパーパラメタの調整など、すべて自動になるわけではない

ビジネス利用への4つの壁

- 河本「会社を変える分析の力」講談社 ← めっちゃ良い本



おまけ：人工知能は人間を超えるか

- 答え：問題によりけり
 - 画像認識は人間を超えたといわれている
 - 音声認識は人間の方が遥かに優れている
 - 当分超える見込はないと思う（個人の見解）

応用編

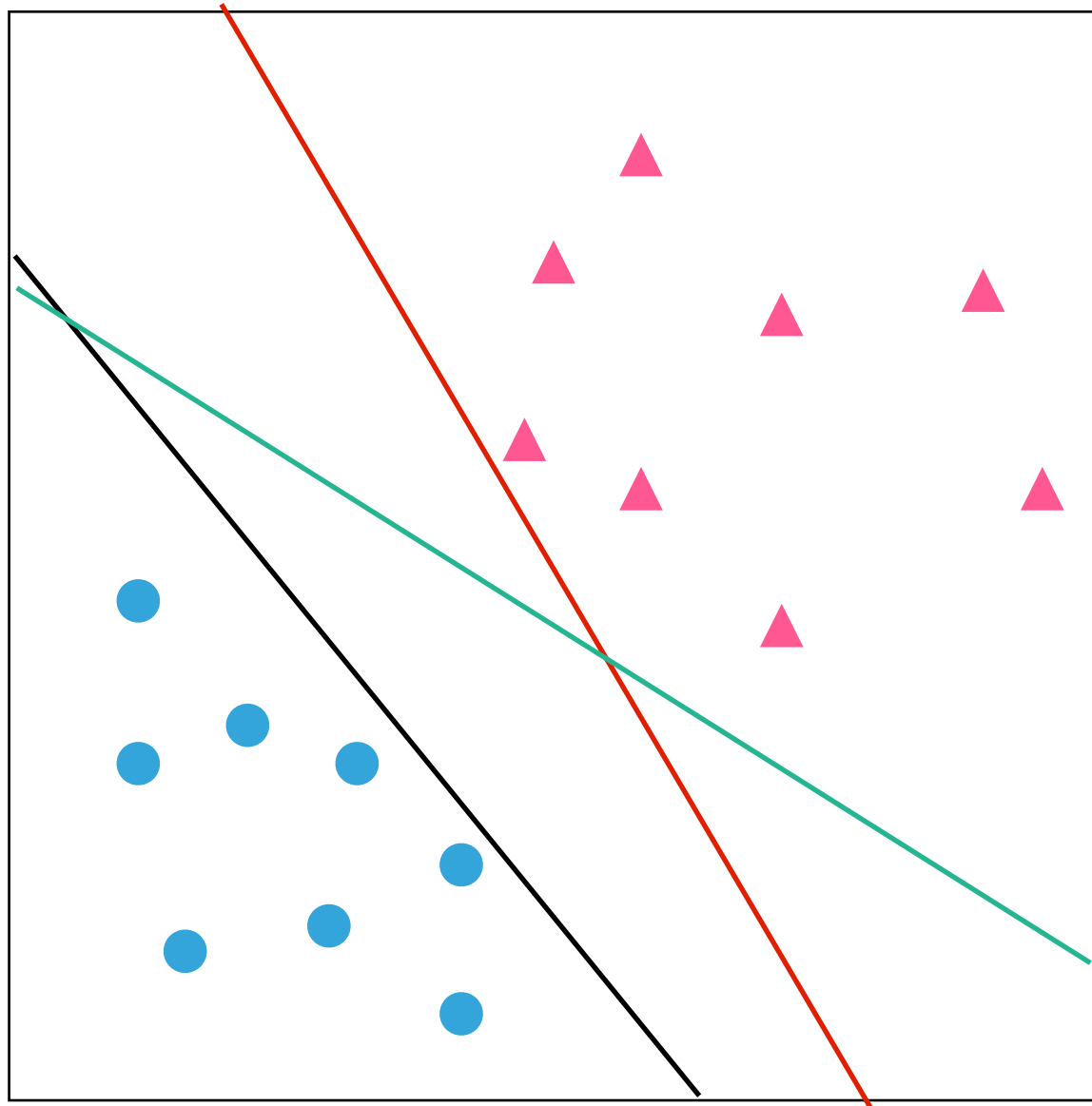
Support Vector Machine

SVMの戦略

▶ マージン最大化

カーネル法による非線形化

パーセプトロン

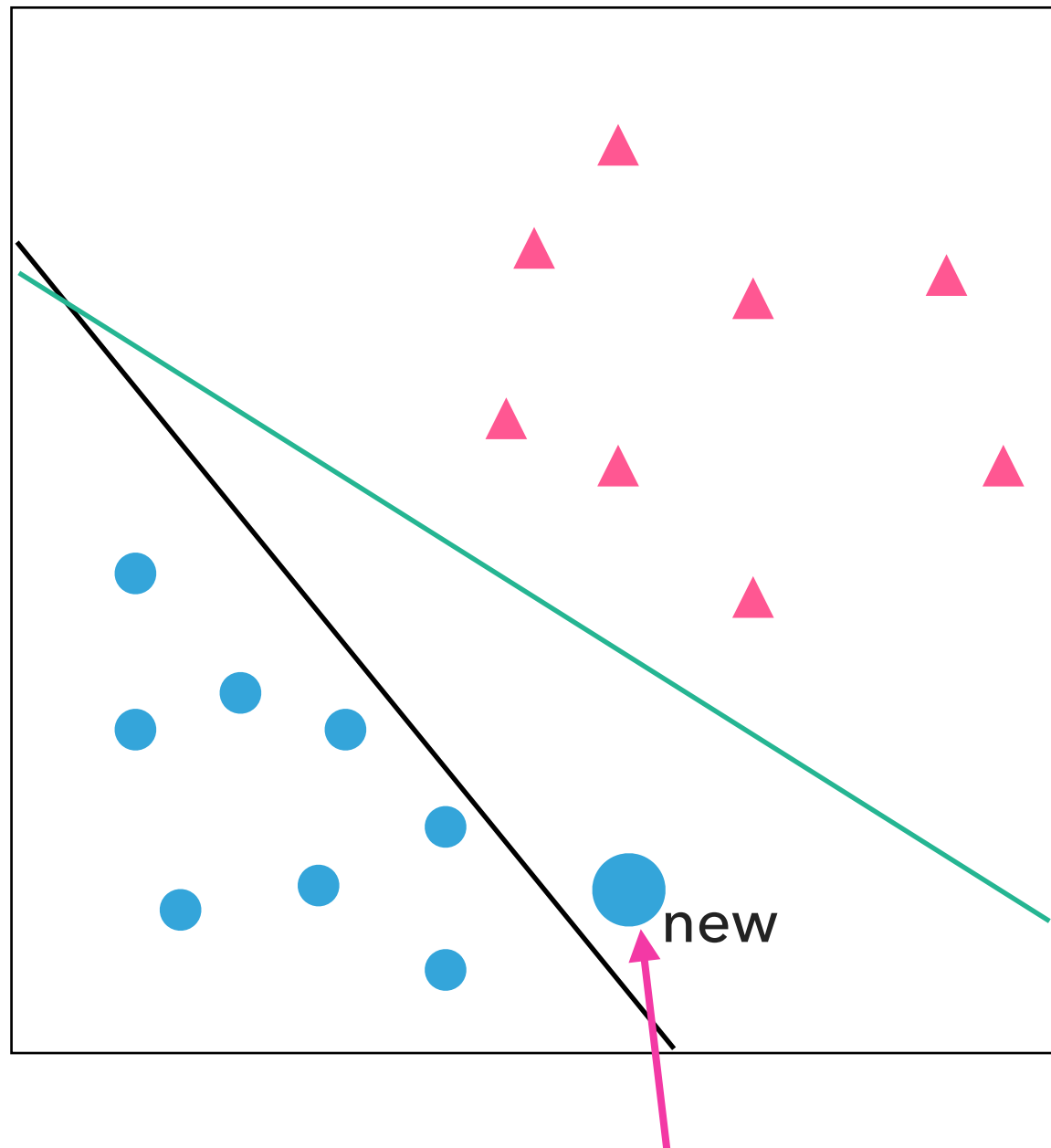


データが線形分離可能なら、必ず決定境界を見つけ出す



どの決定境界に収束するか不確定(初期値に依存)

パーセプトロン



黒い決定境界だと誤分類！

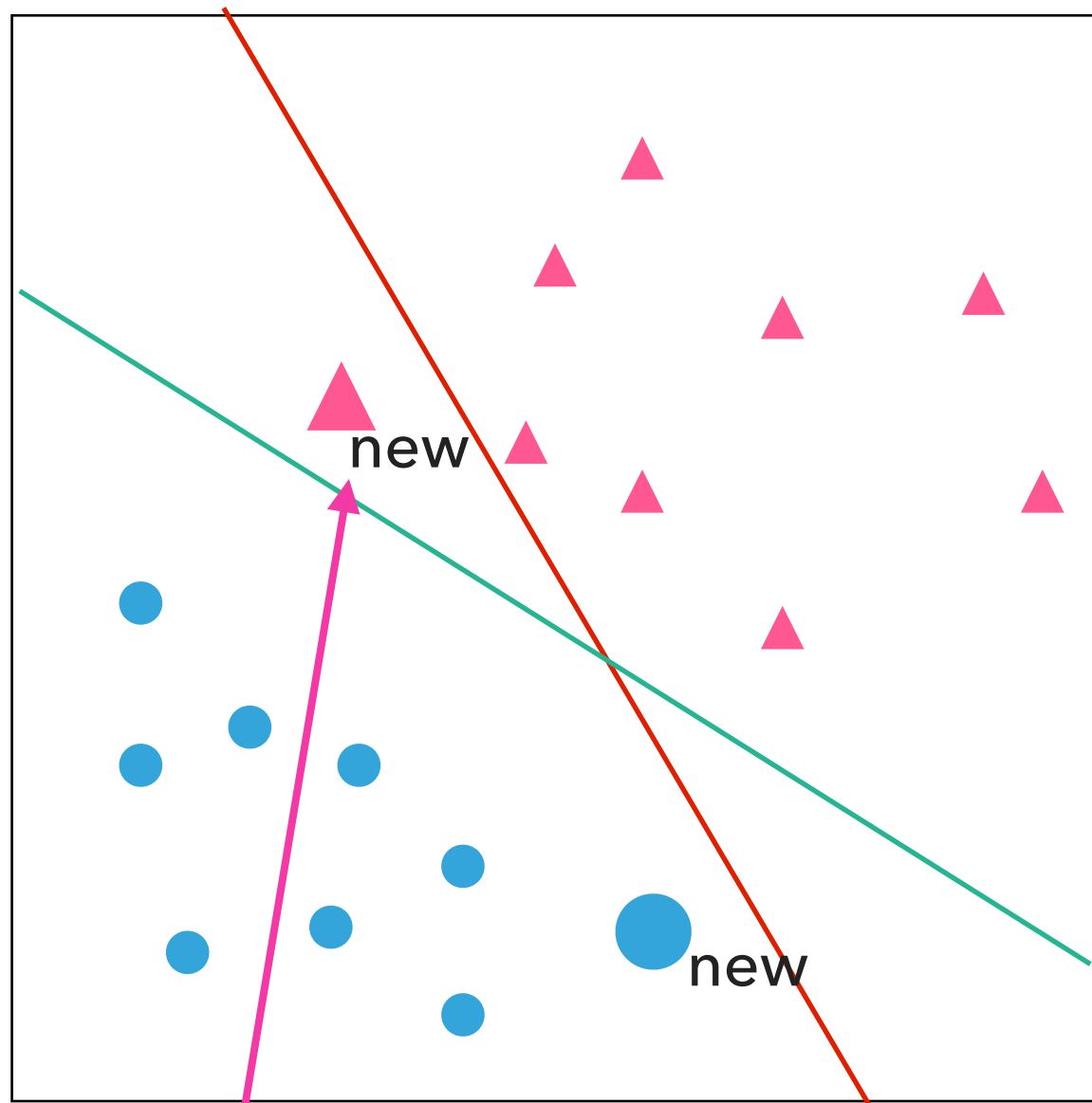


データが線形分離可能なら、必ず決定境界を見つけ出す



どの決定境界に収束するか不確定(初期値に依存)

パーセプトロン



赤い決定境界だと誤分類！

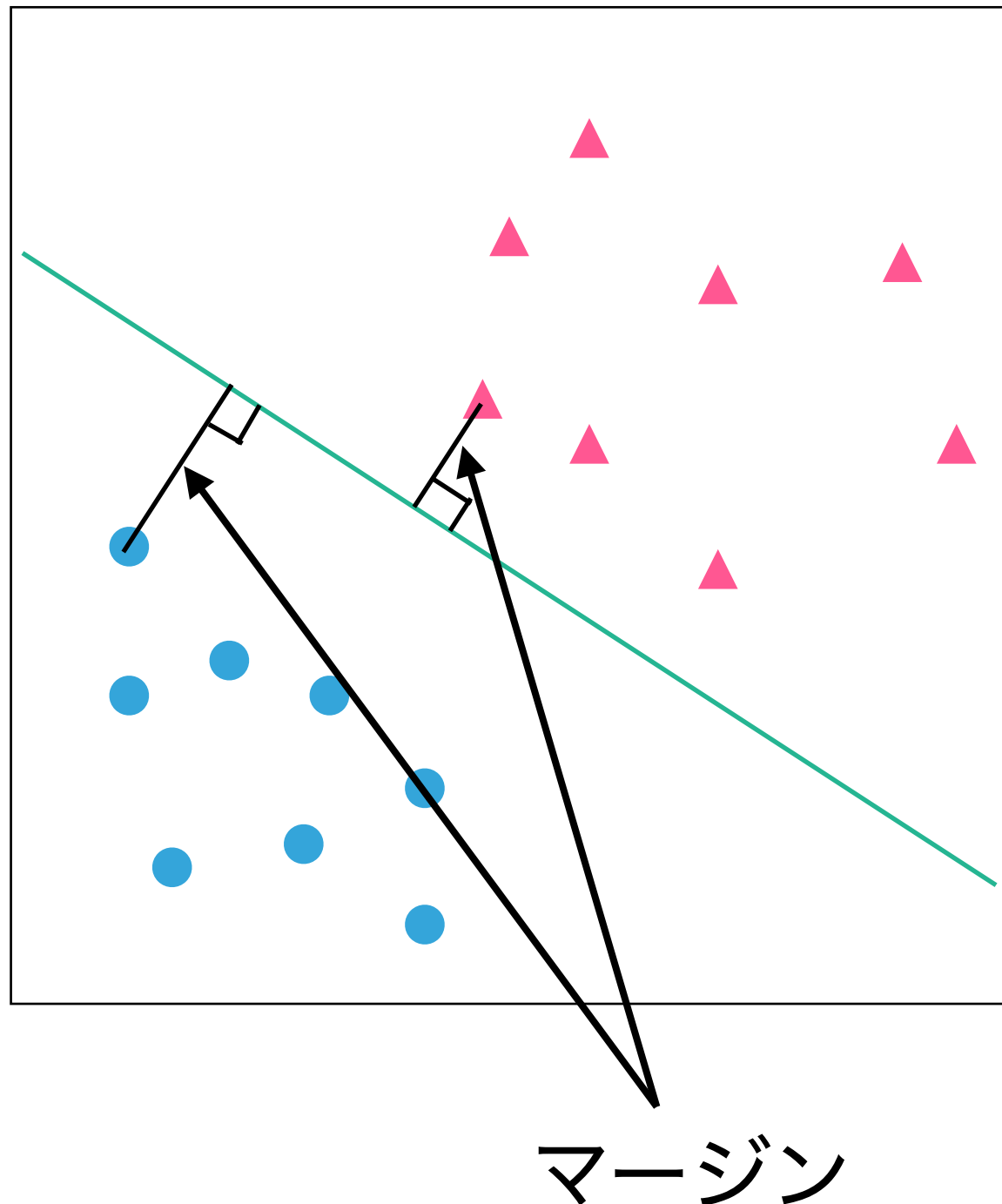


データが線形分離可能なら、必ず決定境界を見つけ出す



どの決定境界に収束するか不確定(初期値に依存)

パーセプトロン



- 😊 データが線形分離可能なら、必ず決定境界を見つけ出す
- 💡 どの決定境界に収束するか不確定(初期値に依存)
- 😊 マージン最大化！
- 💡 線形分離のみ

SVMの戦略

マージン最大化

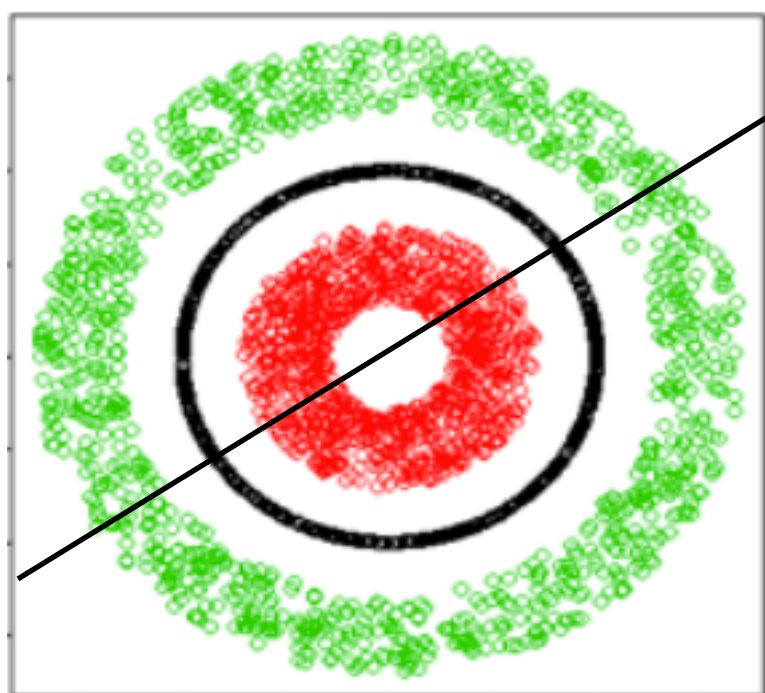
▶ カーネル法による非線形化

やや上級者向け



高次元空間への写像

線形分離不能

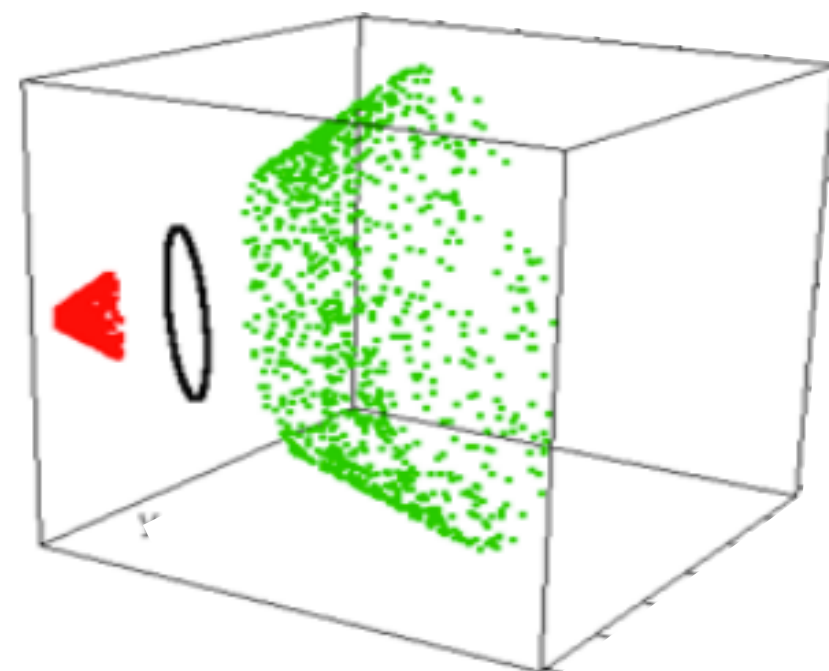


3次元空間へ写像



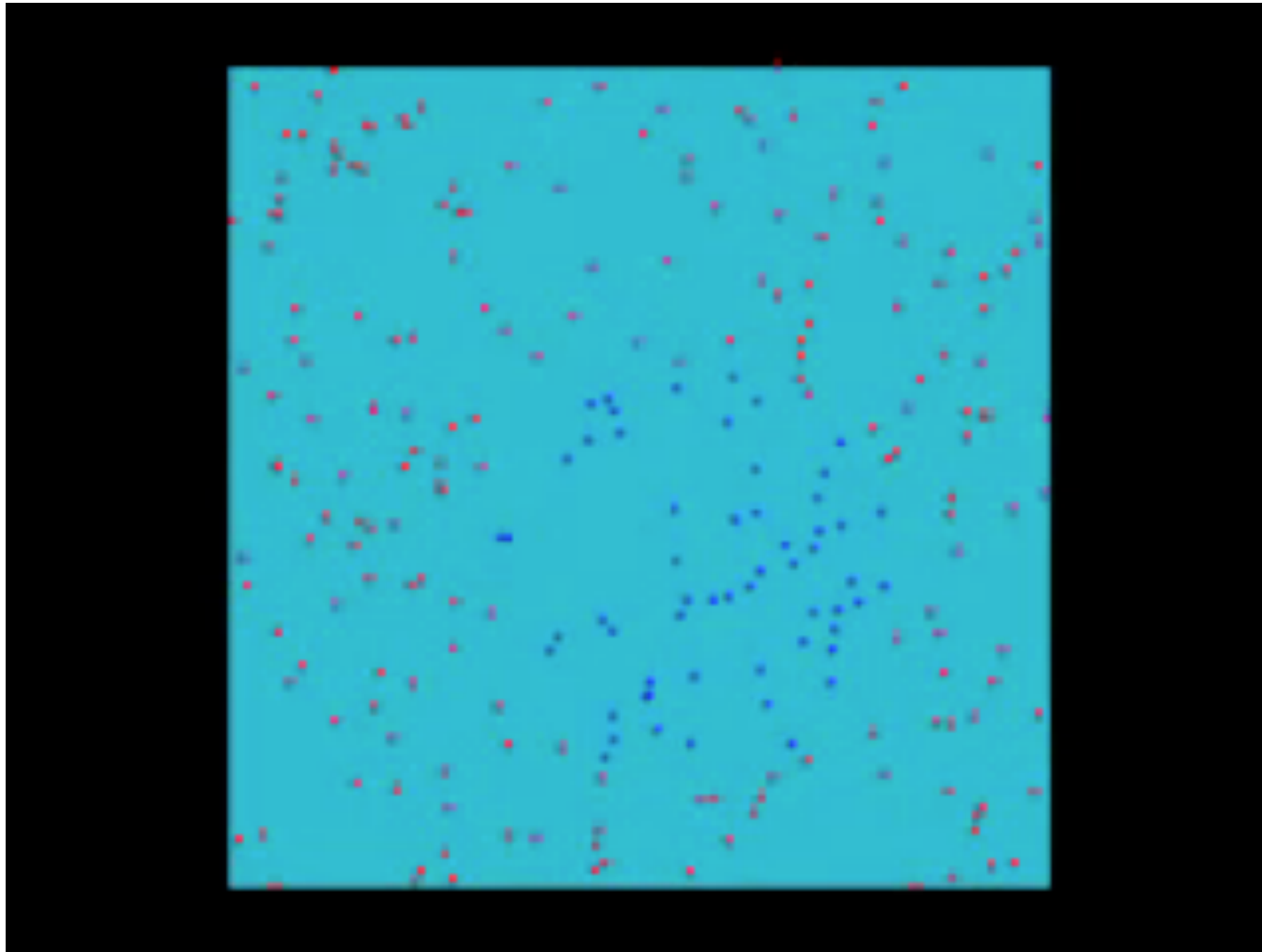
写像 ϕ

線形分離可能



$$\mathbf{x} = (x_1, x_2)$$

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1, x_2)$$



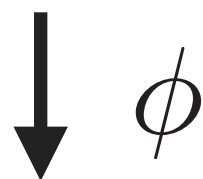
<https://www.youtube.com/watch?v=3liCbRZPrZA&feature=youtu.be>

計算量の爆発（次元の呪い）

- 2次以下の特徴を抽出する写像

$$\mathbf{x} = (x_1, \dots, x_{1000})$$

1000次元



$$\phi(\mathbf{x}) = (x_1^2, \dots, x_1 x_2, \dots, 1)$$

501501次元



表現力 計算量

低次元

低い

少ない

高次元

高い

多い



高次元の表現力を実現しつつ、計算は低次元で行なうアイディア

→ カーネル法！



カーネルトリック

- SVMの定式化

$$\begin{aligned} \text{max.} \quad & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y^{(i)} y^{(j)} \left(\phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)}) \right) + \sum_i \lambda_i \\ \text{s.t.} \quad & \sum_i \lambda_i y^{(i)} = 0 \quad \lambda_i \geq 0 \quad (\lambda \text{ はラグランジュ乗数}) \end{aligned}$$

カーネルトリック

- SVMの定式化

高次元化した特徴ベクトル

$$\max. \quad -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y^{(i)} y^{(j)} \left(\phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)}) \right) + \sum_i \lambda_i$$

$$\text{s.t.} \quad \sum_i \lambda_i y^{(i)} = 0 \quad \lambda_i \geq 0 \quad (\lambda \text{ はラグランジュ乗数})$$

内積さえ計算できればOK
 $\phi(\mathbf{x})$ にアクセスする必要無し

内積の性質を満たす関数K (カーネル関数) で置き換えて計算

$$K \left(\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \right) = \left(\phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)}) \right)$$

カーネル関数 1/2

多項式カーネル $K\left(\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}\right) = \left(\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} + r\right)^d$

↑ 変形

d次以下のすべての項を持つ特徴ベクトルの内積

カーネル関数 1/2

多項式カーネル $K(\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}) = (\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} + r)^d$

変形

計算は低次元！

高次元化した $\phi(\mathbf{x})$ ではなく
 \mathbf{X} についての式なので

表現力は高次元！

$\phi(\mathbf{x})$ の内積に等しいので

d次以下のすべての項を
持つ特徴ベクトルの内積

カーネル関数 2/2

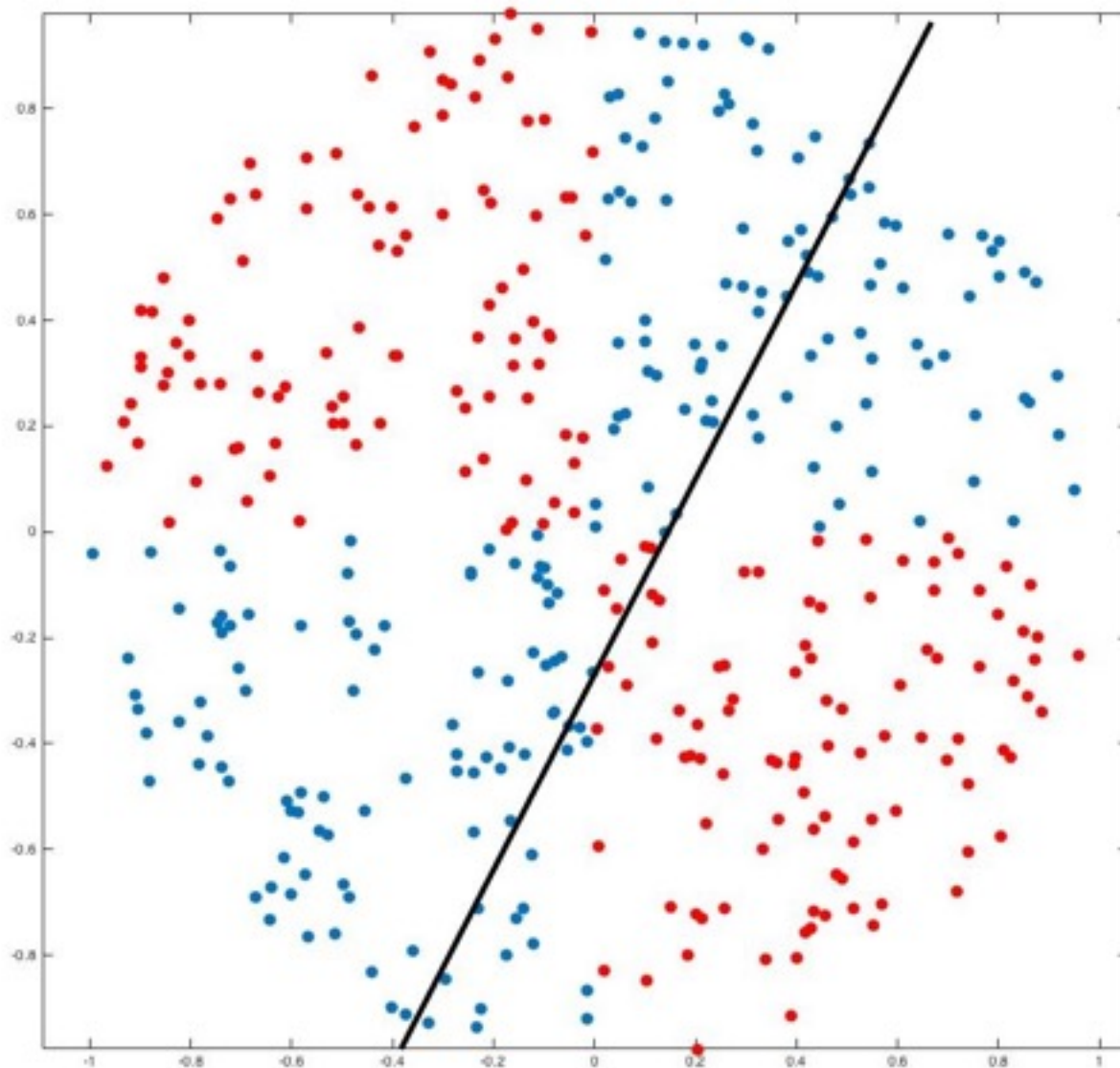
RBFカーネル $K \left(\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \right) = \exp \left(-s \left| \mathbf{x}^{(i)} - \mathbf{x}^{(j)} \right|^2 \right)$

↑ 変形

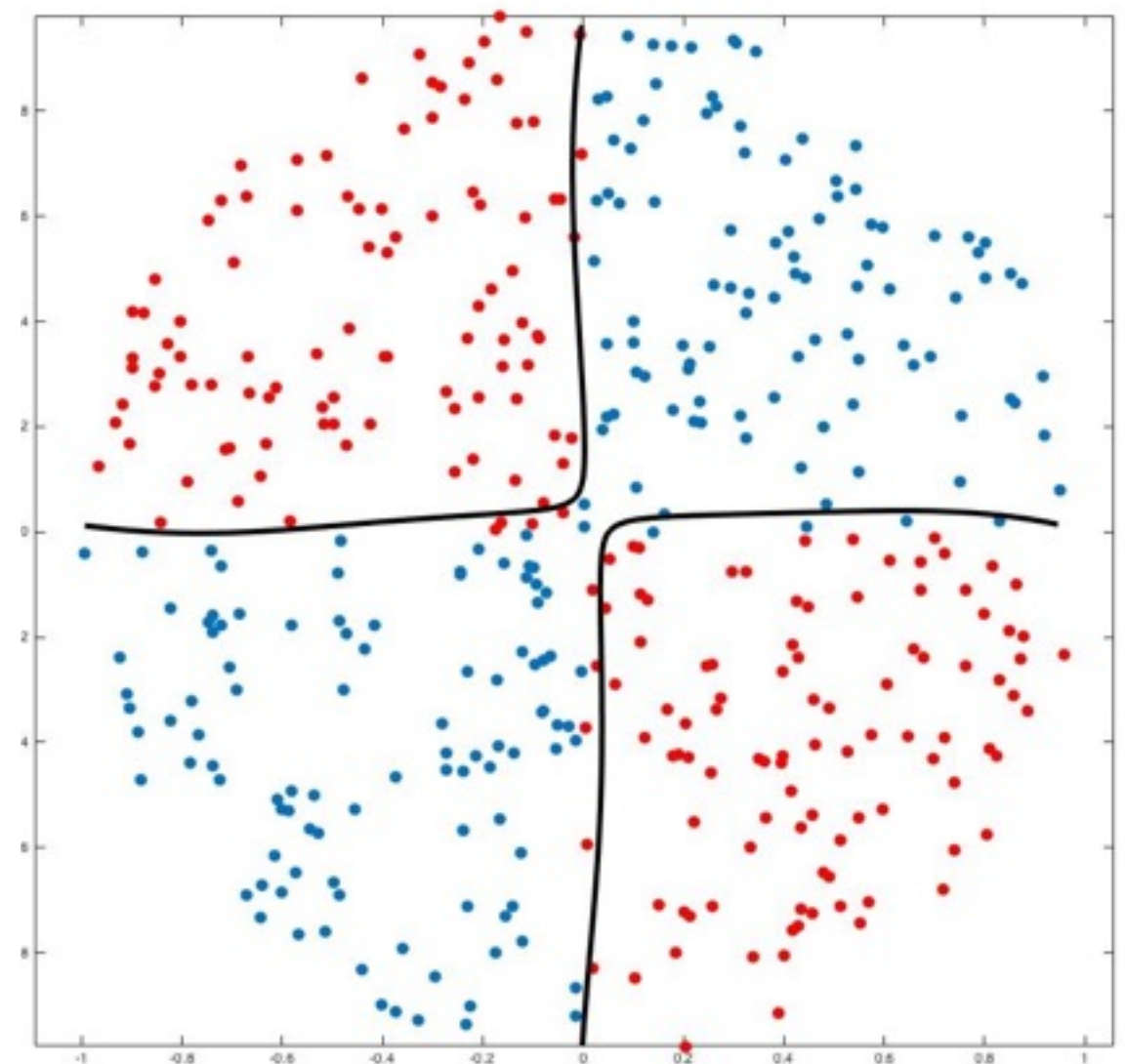
∞次元の特徴ベクトルの内積

カーネル法による決定境界

カーネルなし
(線形カーネル)



RBFカーネル



参考文献

- 比戸 他, データサイエンティスト養成読本 機械学習入門編, 技術評論社
- 河本, 会社を変える分析の力, 講談社
- 杉山, イラストで学ぶ機械学習 最小二乗法による識別モデル学習を中心に, 講談社
- Bishop, Pattern Recognition and Machine Learning, Springer
- 高村, 自然言語処理のための機械学習入門, コロナ社

実践パート

Kaggleに挑戦して
みよう

Kaggle

- 世界最大のデータサイエンティストコミュニティ
- データ解析のコンペティションを多数開催
 - 賞金が出る
 - 他企業の協賛コンペ多数、ジョブマッチング

kaggle

Host


Competitions

Scripts


Jobs

Community ▾


Active Competitions

**Western Australia Rental Prices** ⓘ
Predict rental prices for properties across Western Australia


5.9 days
89 teams
\$100,000

**The Allen AI Science Challenge**
Is your model smarter than an 8th grader?


2 months
383 teams
\$80,000

**The Winton Stock Market Challenge**
Join a multi-disciplinary team of research scientists

2 months
102 teams
\$50,000

**Rossmann Store Sales**
Forecast sales using store, promotion, and competitor data

19 days
2875 teams
1825 scripts
\$35,000

**Prudential Life Insurance Assessment**
Can you make buying life insurance easier?

2 months
114 teams
82 scripts
\$30,000

2 months

タイタニック生存者予想



- Kaggleのチュートリアルコンペ
- タイタニック搭乗者のプロフィールから、その人が生きて帰ったかどうかを予想する
- トレーニングデータ：891人分
- テストデータ：418人分

データに含まれる情報

- Pclass: 搭乗者のクラス (1st, 2nd, 3rd)
- Name, Sex, Age, Fare(料金)
- SlibSp: 同乗した兄弟または配偶者の数
- Parch: 同乗した親または子供の数
- Ticket: チケット番号
- Cabin: 客室
- Embarked: 出発港 (Cherbourg, Queenstown, Southampton)

Pythonライブラリ

- **numpy, scipy**: 数値計算ライブラリ
- **pandas**: データ解析ライブラリ
- **scikit-learn**: 機械学習ライブラリ
- **matplotlib**: グラフ描写ライブラリ
- **IPython**: 対話型シェル
- 1つずつ入れるとめんどいので、**Anaconda**おすすめ
- Kaggleのサイト上でも動かすことができる

コードを書いてみよう

- STEP1

scikit-learnの使い方を調べ、SVMで学習と予測をおう
できたら、出力部分のコメントアウトを解除して実行、
Kaggleに提出して精度を確認

- STEP2

「Fare」と「Age」をそれぞれ正規化した変数「NorFare」
と「NorAge」を作ろう
できたら、FareとAgeの代わりに特徴ベクトルに追加しよう
再度実行、Kaggleに提出して精度を確認

以上

コンタクト : @mkhyt on twitter