# データ解析演習(2) 分類モデルの構築 高間 康史 下川原 英理, 何 宜欣 システムデザイン学部・情報通信システムコース

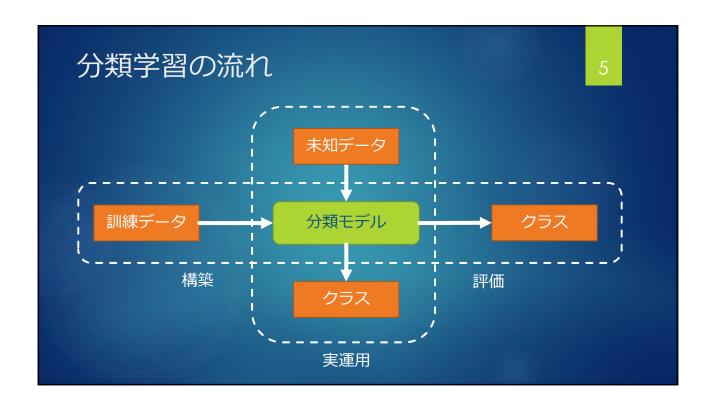
# 講義内容 ) 分類モデルの構築 ) 前回のおさらい ) 手法の分類 ) 最近傍法 ) 1-NN ) K-NN ) ナイーブベイズ ) ベイズの定理 ) サポートベクターマシン (SVM) ) 線形分離可能な場合 ) カーネルトリック ニューラルネットワーク ) 基本アルゴリズム ) 深層学習

# 分類モデル構築とは

3

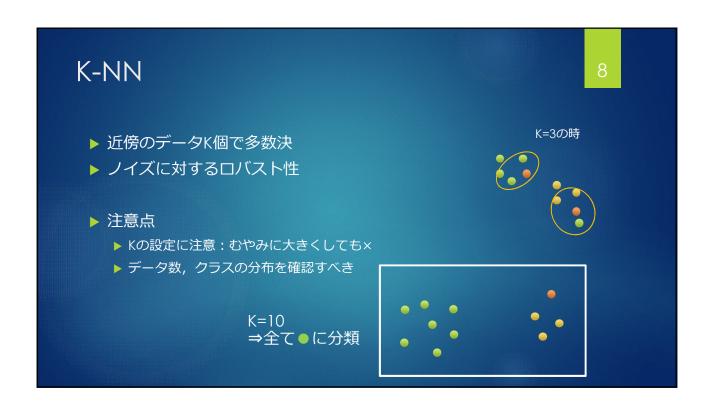
- ▶ 未知のデータを, 既知のクラスに分類
  - ▶ グループ分けという点では、クラスタリングと同様
- ▶ クラスタリングとの違い
  - ▶ **教師あり学習**:訓練データから学習
    - ▶ グループ(クラス)の個数・意味が既知
    - ▶訓練データ(所属クラスが既知)から分類モデル構築
  - ▶ 教師なし学習:訓練データなし ← クラスタリング
    - ▶ グループ(クラスター)の個数・意味が未知
    - ▶ グループの意味は後から考える

# データの話 ▶ 目的変数・従属変数 ※変数→属性と呼ぶ場合もあり ▶ 予測・分類したい対象: クラス ▶ (通常)分類器毎に1つ ▶ 説明変数,独立変数 ▶ 予測・分類に利用する情報 ▶ (通常)複数 ▶ 訓練データ=説明変数+目的変数 ▶ 未知データ=説明変数 X:説明変数 X:説明変数 X:説明変数









# ナイーブベイズ

# Naive Bayes

▶ 各クラスへの所属確率を推定

▶ 確率モデルを訓練データから学習

▶ ベイズの定理

▶ 事前確率: P(X)

▶ 事後確率: *P*(*X*|*Y*)

▶ 尤度: P(Y|X)

▶ 証拠: P(Y)

$$P(X|Y) = \frac{P(X)P(Y|X)}{P(Y)}$$

## 例題

- ▶ ある町で車によるひき逃げ事故があった.
- ▶ その町で走っている車の15%が青色である.
- ▶ 事故の目撃者は、ひき逃げをした車は青色だと証言した.
- ▶ その時間帯・場所において、目撃者は80%正しく識別できることがわかっている。
- ▶ 事故を起こした車が青色であった確率はいくらか

解答

11

- ▶ その町で走っている車の15%が青色である.
  - ▶ X: 車が青色, 1-X: 他の色 ⇒ P(X) = 0.15, P(1-X) = 0.85
- ▶ 事故の目撃者は、ひき逃げをした車は青色だと証言した.
  - ▶ Y: 青色と証言
- ▶ その時間帯・場所において、目撃者は80%正しく識別できることがわかっている.
  - ▶ P(Y|X) = 0.8 ... 実際に青で,正しく証言
  - ▶ P(Y | 1-X) = 0.2 ... 実際は青ではないのに見間違えて証言
  - $\triangleright$  P(Y) = P(Y | X)P(X) + P(Y | 1-X)P(1-X)
- ▶ 事故を起こした車が青色であった確率はいくらか
  - $P(X|Y) = P(X)P(Y|X)/P(Y) = 0.15 \times 0.8 / (0.8 \times 0.15 + 0.2 \times 0.85)$ = 0.41
  - ▶ 証言なし(事前確率): 0.15 ⇒ 証言あり(事後確率): 0.41 に上昇

# 分類モデル構築 例題

- ▶ X: クラス(目的変数). 1変数
- Y: 属性・特徴(説明変数). 通常複数(ベクトル)Y=(曇,高,休日), X=○
- ▶ 確率<u>(事前確率・尤度・証拠)</u>:訓練データから計算
- ▶ Y=(晴, 高, 休日)の時の売上は?
- ▶ 問題:同じYが訓練データに存在しない
- ▶ 解決策:独立性の仮定

$$P(Y|X) = \prod_{i} P(Y_i|X)$$

天気	気温	休/平日	売上
曇	高	休日	0
晴	低	休日	$\circ$
雨	低	休日	0
曇	高	平日	0
晴	高	平日	0
曇	低	休日	$\circ$
雨	高	平日	×
曇	低	平日	×
晴	低	平日	×
曇	低	平日	×

解答

▶ 天気

▶ 売上○: 9 ... 晴: 3, 曇: 4, 雨: 2

▶ 売上×: 7 ... 晴: 2, 曇: 3, 雨: 2

▶ 気温

▶ 売上○: 8 ... 高: 4, 低: 4

▶ 売上×: 6... 高: 2, 低: 4

▶ 休/平日

▶ 売上○: 8 ... 休日: 5, 平日: 3

▶ 売上×: 6 ... 休日: 1, 平日: 5

▶ ラプラス推定:各頻度に1追加 ←ゼロ頻度問題への対処

P(晴,高,休日|○)

= P(晴 $| \circ )P($ 高 $| \circ )P($ 休 $| \circ ) = \frac{345}{988} = \frac{5}{48}$ 

 $P(\mathbf{f}, \mathbf{a}, \mathbf{h}) = \frac{221}{766} = \frac{1}{63}$ 

 $P(\circ |$ 晴, 高, 休日)  $\propto P(\circ)P($ 晴, 高, 休日 $|\circ) = \frac{5}{48} \frac{7}{12} = 0.061$ 

 $P(\times |$ 晴, 高, 休日)  $\propto P(\times)P($ 晴, 高, 休日 $|\times) = \frac{1}{63} \frac{5}{12} = 0.0066$ 

# 離散データと連続データ

▶ 説明変数が離散(Nominal), 連続(Numerical)どちらでも利用可能

Class

Class **Attribute** yes no (0.63) (0.38) P(x)temperature 離散型  $P(y_i|x)$ mild 確率分布 12.0 8.0 [total] humidity high

7.0

7.0

11.0

Attribute (0.63) (0.38)temperature

72.9697 74.8364 5.2304 7.384 weight sum 9 precision 1.9091 1.9091

連続型 確率分布 (正規分布)

※ラプラス推定

[total]

78.8395 86.1111 std. dev. 9.8023 9.2424 weight sum 9

precision 3.4444 3.4444

# サポートベクターマシン SVM: Support Vector Machine

15

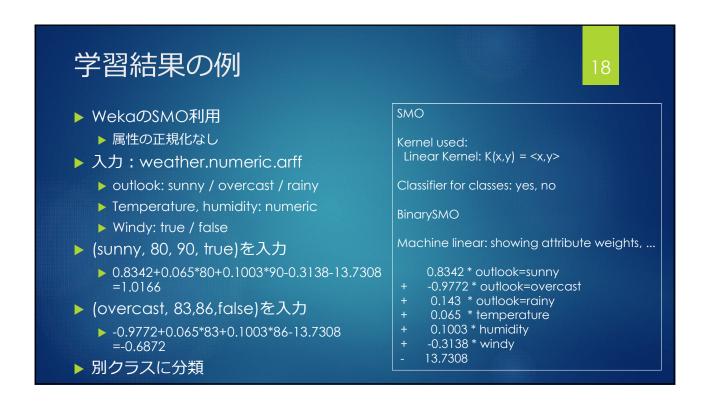
- ▶ 近年主流.多くの場合で高性能
- ▶ クラスを分離する境界(平面)を求める
- ▶ 訓練データの偏りに頑健 境界付近のデータのみをモデル構築に利用
- ▶ カーネルトリック:線形分離不能な場合にも有効
- ▶ WekaではSMO, libSVM

# 線形分離とは

- ▶ 各クラスのデータを分離する平面の存在
  - ▶ 2次元の場合:直線
  - ▶ 4次元以上の場合:超平面
  - $\mathbf{w}^{\mathrm{T}}x + w_0 = 0$
- ▶線形判別関数
  - ▶ 平面のどちら側にあるかでクラスを判別
- ▶ 線形分離不可能な場合もあり
- ▶ とりあえずは線形分離可能な場合について説明











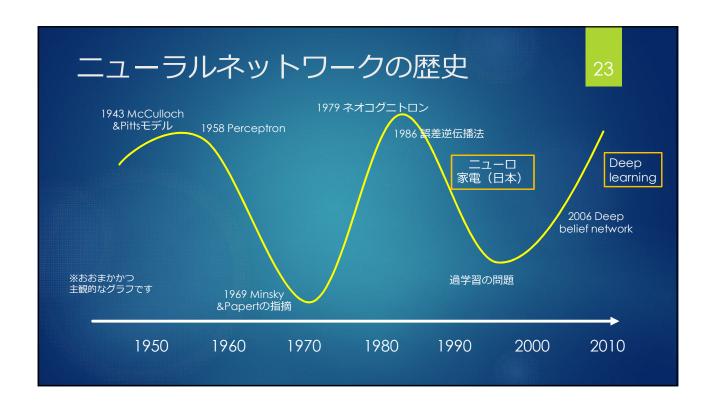
### カーネルトリック

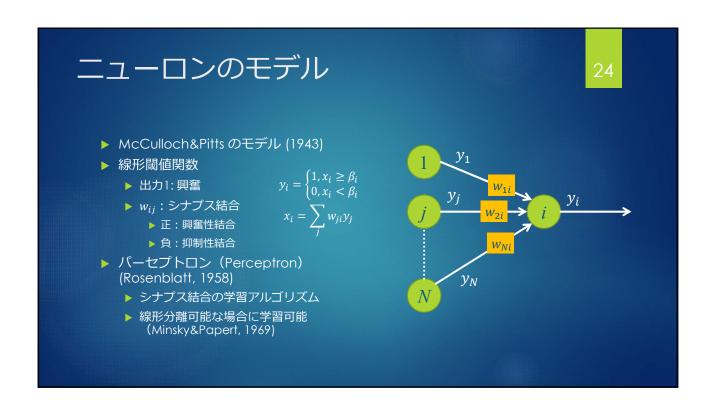
2

- ▶ 元の空間:線形分離不可能 → 高次元空間:線形分類可能
- ▶ 高次元空間で線形識別関数を求める
- ▶ カーネル関数:  $K(x,x') = \phi(x)^{\mathrm{T}}\phi(x')$ 
  - ▶ φ(x): 高次元空間への非線形写像
  - ▶ 識別関数 $g(\phi(x)) = \mathbf{w}^{\mathrm{T}}\phi(x) + w_0 = \sum_{x^i} \alpha_i y^i \phi(x)^{\mathrm{T}}\phi(x^i) + w_0$ =  $\sum_{x^i} \alpha_i y^i K(x, x^i) + w_0$  … 個々のデータの非線形変換は不要!
- ▶ カーネル関数の例
  - ▶ 多項式カーネル:  $K(x, x') = (x^T x' + l)^p$
  - ▶ RBFカーネル:  $K(x,x') = \exp(-\sigma ||x x'||^2)$  (ガウシアンカーネル)

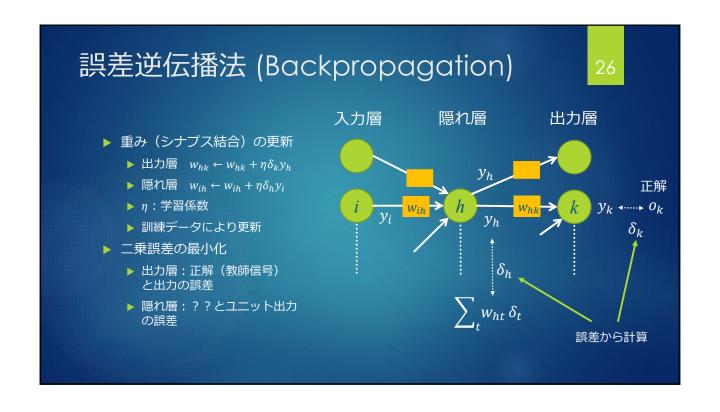
# ニューラルネットワーク Neural Network

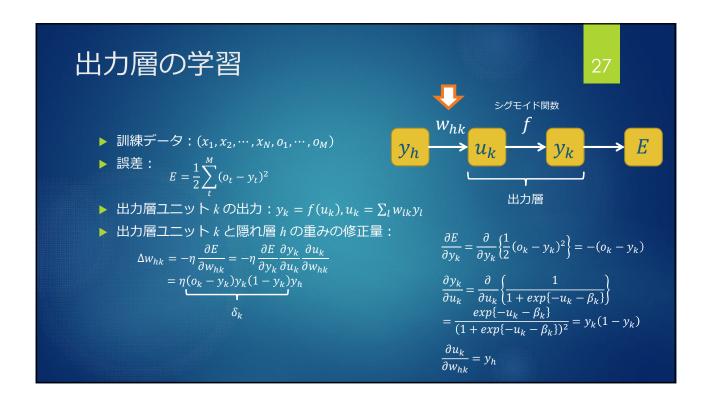
- ▶ 脳神経系をモデルにした情報処理システム
- ▶ 線形分離不能な場合でも機能
- ▶ 学習結果がブラックボックス ⇔ 決定木
- ▶ Weka: MultilayerPerceptionで利用可能

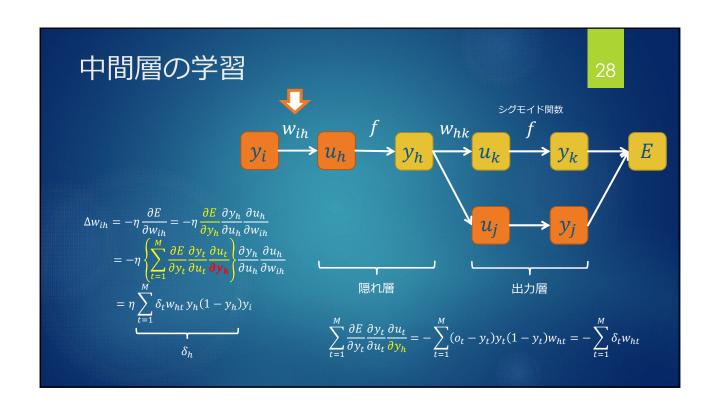


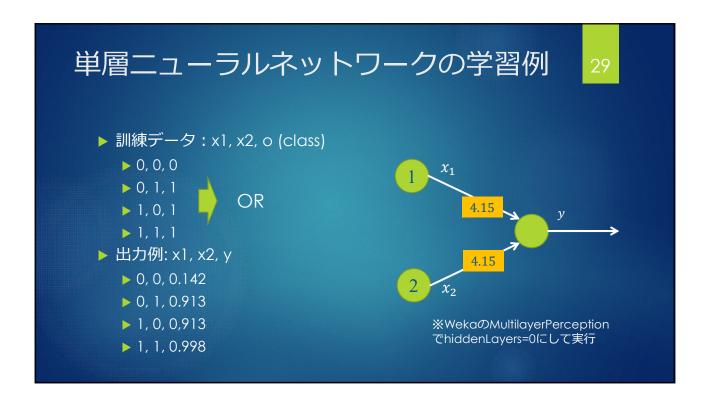


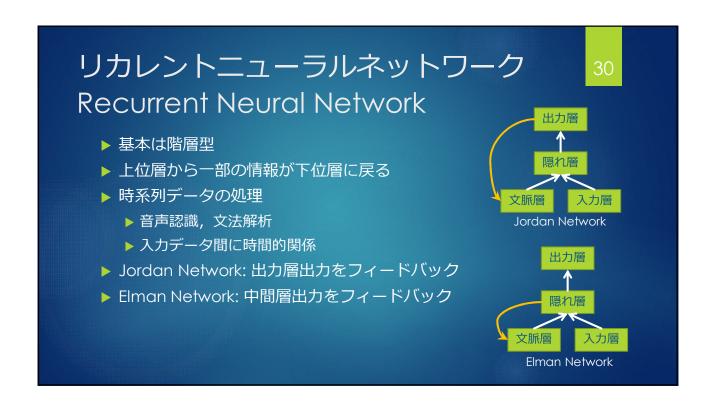


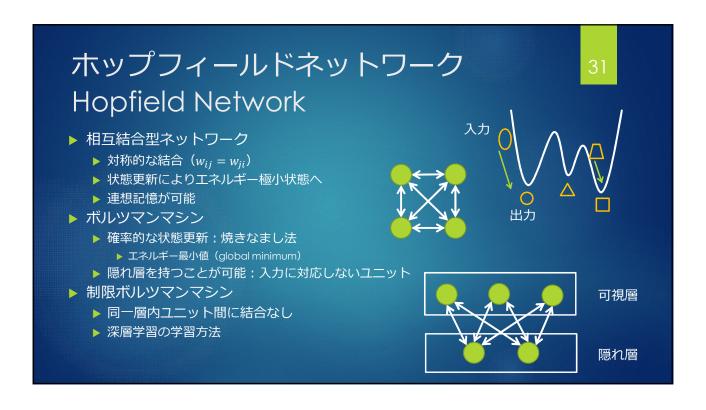




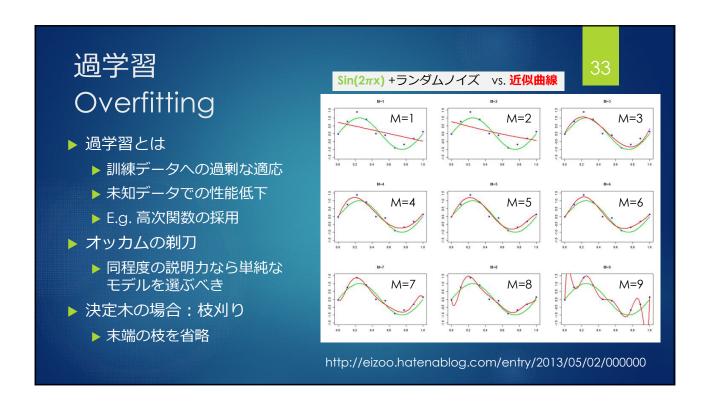












# SVMとの関係

- ▶ どちらも分離超平面の学習
- ▶ ニューラルネットワーク
  - ▶ 訓練データの誤差が小さくなるように決定
  - ▶ 訓練データ依存:過学習,局所最適解
- ► SVM
  - ▶ マージン最大化
  - ▶ 正例・負例の境界の真ん中に設定

# 深層学習(Deep Learning)

35

- ▶ 認識(分類モデル構築):画像,音声
  - ▶ 2012: 画像認識, 音声認識, 化合物の活性予測で最高精度達成(圧倒的). ドメインの非専門家が達成
  - ▶ Microsoft: 音声認識を深層学習ベースに
- ▶ 大規模ニューラルネットワークからの教師なし学習
  - ▶ Google: Youtubeから200x200ピクセルの画像 (動画から1枚ずつ) 1000万枚をランダムに選択して学習
  - ▶ 大規模:9階層, 10億パラメータ, 1000台x 3日間
  - ▶ 人の顔, 体, 猫の顔に反応するニューロン

Q. V. Le et al., Building High-level Features Using Large Scale Unsupervised Learning, ICML12, 2012.









ニューロンが一番反応する画像(数値計算)

## 深層学習の利点

- ▶ 特徴抽出の問題を解決:職人芸から自動化
- ▶ 表現学習
  - ▶ 「観測データから本質的な情報を抽出し表現」する方法を学習
  - ▶ Cf. 主成分分析, データ圧縮, データ要約, クラスタリング
- ▶ ラベルなしデータから基本特徴量の有効な組合せを学習可能
  - ▶ クラスタリング:局所的な表現学習
  - ▶ 深層学習:大域的な表現学習
    - ▶ 階層的:単純な表現 → 総合的・抽象的な特徴に変換

### 深層学習の学習方法

57

▶ 多階層のネットワーク:どうやって学習?

- ▶ プレトレーニング:層毎に学習
  - ▶目的:入力層に近いほど学習が進まない問題の解決
  - ▶ 入力データの集合を再現できるように教師なし学習
    - ▶ オートエンコーダ or 制限ボルツマンマシン
  - ► その後ファインチューニング:全体を教師あり学習 (誤差逆伝播法)
- ▶ 畳みこみニューラルネットワーク
  - ▶ 歴史あり:初期は手書き文字認識
  - ▶ 脳の視覚野がモデル:受容野→単純型細胞→複雑型細胞(抽象的特徴)
  - ▶ 平行移動に対する不変性
  - ▶ プレトレーニング不要

(上)第2層→ (下)第3層
taces
can

H. Lee et al., Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations, ICML09, 2009.

まとめ

- ▶ 今日学んだこと
  - ▶最近傍法
  - ▶ ナイーブベイズ
  - ▶ サポートベクターマシン
  - ▶ ニューラルネットワーク
- ▶ 来调
  - ▶ グループワーク・発表
- ▶ 今日の演習
  - ▶ データ作成と分類モデル構築