

人工知能A

Topic 8: PAC learning, SVM, boostingなど

8

SVMとBoosting

- PAC学習の概要、SVMとアンサンブル学習のBaggingやBoosting (AdaBoost) の基礎について学ぶ。
 - よく見かけるAdaBoostはPAC学習の研究が基本になっています。
- 講義の内容
 - 数学の復習など少し
 - PAC学習（ここは考え方だけ）、
 - Support Vector Machine (SVM)
 - Boosting (Adaboost)
- より詳細は各自で自習するか、大学院で必要に応じて習得する。ここではその準備。

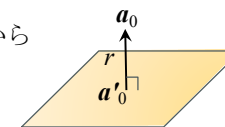
2

3

数学の復習と準備

数学の復習(1)

- n 次元ベクトル空間の（超）平面
$$w_1x_1 + \dots + w_nx_n = b$$
と平面上にない点 $\mathbf{a}_0 = (a_{01}, \dots, a_{0n})$ を考える。 \mathbf{a}_0 から平面に下ろした垂線との交点を \mathbf{a}'_0 、長さを r とする。
- $\mathbf{x} = (x_1, \dots, x_n)$, $\mathbf{w} = (w_1, \dots, w_n)$ とおけば、この平面は $f(\mathbf{x})=0$ 。ただし $f(\mathbf{x}) = \mathbf{w}\mathbf{x}^T - b$ とおいた。
- \mathbf{w} の単位ベクトルは、 $\frac{\mathbf{w}}{\|\mathbf{w}\|}$ となる。ただし、 $\|\mathbf{w}\| = \sqrt{w_1^2 + \dots + w_n^2}$
- したがって、 $\mathbf{a}_0 = r \frac{\mathbf{w}}{\|\mathbf{w}\|} + \mathbf{a}'_0$ と書ける。 $f(\mathbf{a}_0) > 0$ のときは $r > 0$ 。逆方向の場合は、 r が負の数と考えればよい。
- $f(\mathbf{a}_0) = \mathbf{w}\mathbf{a}_0^T - b = \mathbf{w}\mathbf{a}'_0{}^T - b + \frac{r}{\|\mathbf{w}\|} \mathbf{w}\mathbf{w}^T$ となり、赤の部分は0 (平面上なので) 青は、 $\|\mathbf{w}\|^2$ となるので、 $r = f(\mathbf{a}_0)/\|\mathbf{w}\|$ と書ける。



4

数学の復習(2)

- 関数のベクトル $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))^T$ ただし、

$$f_k: \mathbb{R}^n \rightarrow \mathbb{R} \text{ for } 1 \leq k \leq m$$

を考える。 $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ である。

$m \times n$ 行列 \mathbf{A} と m 列のベクトル \mathbf{b} が存在して、

$$\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}^T + \mathbf{b}$$

と表せるとき、 $\mathbf{f}(\mathbf{x})$ をアフィン関数 (affine function) という

- つまり線形関数です (まれですが、 $\mathbf{b} = \mathbf{0}$ のとき、特に線形関数と呼ぶこともあります)。

- 実数の $n \times n$ 行列 $\mathbf{A} = (a_{ij})$ が半正定値 (positive semi-definite or nonnegative definite) とは、 $\mathbf{0}$ でない任意のベクトル \mathbf{c} に対し、 $\mathbf{c}\mathbf{A}\mathbf{c}^T \geq 0$ となること。

- 例: $\begin{pmatrix} 9 & 6 \\ 6 & 4 \end{pmatrix}$

$$\therefore (x \ y) \begin{pmatrix} 9 & 6 \\ 6 & 4 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = 9x^2 + 12xy + 4y^2 = (3x + 2y)^2 \geq 0$$

5

数学の準備

- このあとの「数学の準備」は紹介だけ。
- もしかするとオートマトンとかラムダ計算などで知っている人がいるかもしれません。
- 数学科とか応用数理学科の学生であれば、数学基礎論の入門部分でやると思います。

6

数学の準備(1)

- 帰納的関数 (recursive functions) とは、(直感的には) プログラムとして書ける関数のこと。
- 再帰的関数、計算可能関数とも言う。
- 同様な概念が、チューリング機械やラムダ計算 (つまりプログラム言語の Scheme や LISP) でも定義でき、これらは同義であることが知られている。
- 原始帰納的関数 (primitive recursive functions)
 - 原始再帰的関数とも言う。
 - ここでは制御構造が loop のみで書かれたプログラムでかける関数のことと考えてよい。

7

数学の準備(2)

- n 次元実数空間上の (プログラムとして書ける) 実数値関数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ に対し $C_+(f) = \{\mathbf{x} \in \mathbb{R}^n | f(\mathbf{x}) > 0\}$, $C_-(f) = \{\mathbf{x} \in \mathbb{R}^n | f(\mathbf{x}) < 0\}$ と \mathbb{R}^n を分割することを考える。

- 上記の f から容易に以下のようなプログラムを作成できる。

$$g(\mathbf{x}) = \begin{cases} 1 & \text{for } \mathbf{x} \in C_+(f) \\ -1 & \text{for } \mathbf{x} \in C_-(f) \end{cases}$$

- このような $g(\mathbf{x})$ を判別関数 (分類関数) と呼ぶ。(=0 の部分は適当にどちらかの分割領域に組込でもよい)
- 特に f を超平面 $f(\mathbf{x}) = \mathbf{w}\mathbf{x}^T - b$ とすると、 $C_{\pm}(f)$ はこの面で分割された領域を表す。
- なお $\text{sgn}(h(\mathbf{x})) = \begin{cases} 1 & \text{for } h(\mathbf{x}) \geq 0 \\ -1 & \text{for } h(\mathbf{x}) < 0 \end{cases}$ と定義 ($h(\mathbf{x}) = 0$ のときの定義は状況によって)

8

9 PAC学習

まずだいたいのお話から。SVM（特にソフトマージンSVM）の理論的基礎やAdaboostなどが生まれる背景には、計算論（つまりNP-Completeとかいうやつ）にもとづく学習理論がありました。

計算論に基づいた帰納推論と学習理論

- お話として歴史的背景から。
 - 先に参考文献
 - [Gold67] E. M. Gold, Language identification in the limit, *Information and Control*, 1967.
 - [Blum75] L. Blum and M. Blum, Toward a mathematical theory of inductive inference, *Information and Control*, 1975.
 - [Valiant84] L. G. Valiant. 1984. A theory of the learnable. *Commun. ACM* Vol.27, 11, 1984.
- 「帰納推論・帰納学習 (inductive reasoning/learning, induction)」とは、いくつかの具体例から一般的な法則や理論を導き出す過程である。
- 学習の一形式と考えられる。

10

計算論に基づいた帰納推論

- まずは何をしようとしているか。
 - C, Java, pythonなどでプログラムされた関数 $F(x)$ がある。
$$F: \mathbb{R} \rightarrow \mathbb{R}$$
入力と出力の値 $(x_i, F(x_i))$ をいくつか観測して、そのプログラムの中身を知ることはできないか？
 - 簡単のため入力は一変数、出力は $+1, -1$ のみとします。
- プログラムで書かれた関数 $F(x)$ のことを帰納的関数と呼びました。言い換えると。
 - ± 1 の値をとる帰納関数 $F(x): \mathbb{R} \rightarrow \mathbb{R}$ の入出力のペアをいくつか観測して、帰納関数の定義を知りたい。

11

計算論に基づいた帰納推論

- 極限による同定（帰納推論の形式的な定義）
 - \mathcal{R} を $F: X \rightarrow \{-1, 1\}$ ($F \in \mathcal{R}$)なるプログラム (帰納的関数) の集合とする。 X は共通の定義域。
 - プログラム $F \in \mathcal{R}$ が生成する学習データの無限の列 $(x_1, y_1), (x_2, y_2), \dots$ を仮定する。ここで $F(x_i) = y_i$, ただし $y_i = -1$ or 1 . $y_i = 1$ の学習データを正例、 $y_i = -1$ のデータを負例と呼ぶこともある。つまり教師あり学習。
 - 学習データの無限の列が正例のみのも考えられる（負例はあるのだが、あえて与えない）。

12

計算論に基づいた帰納推論

- 極限による同定（帰納推論の形式的な定義）
 - 定義 (Identification in the limit, Gold 1964) :
 $\forall F \in \mathcal{R}$ に対し、正の整数 $n > 0$ に対し、アルゴリズム（プログラム） \mathcal{A} に、 $(x_1, y_1), \dots, (x_n, y_n), \dots$ を順に入力すると \mathcal{A} は帰納的関数 $h_1, h_2, \dots \in \mathcal{R}$ を順に出力し、
$$h_n(x_j) = y_j \text{ for } 1 \leq j \leq n \text{ となるとする。}$$

このとき十分大きい自然数 N が存在し、 $\forall n \geq N$ なら \mathcal{A} は同じ $h_n = h_N$ を出力し続けるとき \mathcal{A} は \mathcal{R} を極限で同定可能（学習可能）という。

13

計算論に基づいた帰納推論

- 当初は計算論というよりは、言語学習（幼児はどうして言語を学習できるのか）の内容が中心であった [Gold67]
 - たとえば、幼児は親から言葉の正例だけが与えられて、負例はあまり与えない（たとえばめちゃくちゃな言葉を使って「これは日本語ではない」というような教育はしていない）。
 - それなのに自然言語を獲得できる。それを理論的に証明できるか？（答えはNOであった。負例が与えられればある程度はYES）。
 - 当時の言語学はチョムスキー学派に権威があり、その言語階層に従った。
 - 言語学において自然言語は、一部の例外を除き「文脈自由文法」の範囲であるらしいことが1980年ごろから通説となった。
 - チョムスキー階層（形式言語）は、特に計算理論の関連性が重要視され、現在では計算機科学の分野でのみ多く利用されている。
- [Blum75]が、この研究が「プログラムによる学習可能性と計算量」と繋がることを示し、計算論的学習理論の研究が進む（現在も）。

14

計算論に基づいた帰納推論

- 学習としては完全に正しくなるという条件が厳しかった。
 - 当然、訓練データものどのように現れる（与えられる）か分からない。
 - データに誤差はないのか？
 - 極限で同定と言われても。。。
 - \Rightarrow PAC learning

15

計算論的学習（PAC学習）(1)

- PAC learning = Probably approximately correct learning
- 何をしようとしているか。
 - C, Java, pythonなどで書かれた関数 $F(x)$ がある。
$$F: \mathbb{R} \rightarrow \{-1, 1\}$$

 \mathbb{R} （整数や有限集合でもよい）上の確率分布があり、その確率に基づいて x_1, \dots, x_n, \dots が生成され、それに従って入力と出力の組（ただし $F(x) = \pm 1$ ）
$$(x_1, F(x_1)), \dots, (x_n, F(x_n)), \dots$$

をいくつか観測して、 $F(x)$ と ほとんど 同じとなるプログラムを高い確率で生成できるか？
 - 「ほとんど」とは高確率で正しいこと。頻繁に発生する x について $F(x)$ は正しく、あまり発生しない x については正しくないかも。

16

計算論的学習 (PAC学習) (2)

- PAC学習では標本空間 (起こり得る事象) X (前では \mathbb{R}) 上の確率分布 \mathcal{D} を仮定し、標本 S はその分布に従って選択される。
- \mathcal{D} に従って $x \in X$ が選択される確率を P_x と置く。 $\sum_{x \in X} P_x = 1$ である。
- X 上の判別関数の集合 H (前の \mathcal{R} に相当) を仮説集合という。ここから正解を探す。
- PAC学習では学習アルゴリズム \mathcal{L} の成功を示すパラメータとして信頼度 δ と近似度 ε を導入する (ただし、 $0 \leq \delta, \varepsilon \leq 1$) 。
- 学習データの表現は前述と同じ。その真の判別関数を $h^* \in H$ と表す。 \mathcal{L} は H から候補 h を出力する。なお、 $h^* \in H$ と仮定する。
- h^* と h の誤差 (差分、距離) を

$$d_X(h^*, h) = \sum_{x \in X, h^*(x) \neq h(x)} P_x$$

と定義する。全体での誤差を測る。

17

計算論的学習 (PAC学習) (3)

- 帰納的関数のクラス H を PAC 学習可能とは、判別関数 $h^* \in H$ と X 上の任意の確率分布 \mathcal{D} に対して、多項式時間の学習アルゴリズム \mathcal{L} が存在し、任意の $0 \leq \delta, \varepsilon \leq 1$ について

$$\Pr(d_X(h^*, h) \leq \varepsilon) \geq 1 - \delta$$

となる $h \in H$ を出力し停止すること。[Valiant84]

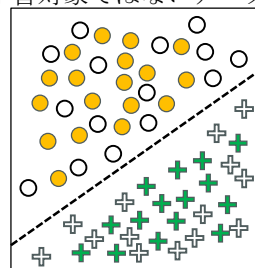
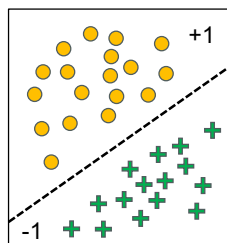
- したがってあまり与えられないデータについては正解できなくても仕方ない (ε)。[逆に言えば正解して欲しい部分は、ちゃんと学習
→ Boosting (Adaboost)]
- データも確率的に与えられるので、望ましいデータ系列ばかりではないので、このような h をいつも出力できるわけではない ($1 - \delta$)。[なら、数回試したり並列化して多数決を取ればいいのかも (アンサンブル学習)]

18

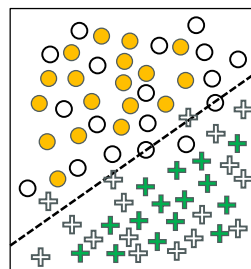
計算論的学習 (PAC学習) (4)

- 学習対象ではないデータも分類 (極限による同定)

- 学習サンプル (S) と区分



白抜きの点を含めて X



- 多少は誤差があっても、ほとんどが合っていればよい (PAC)

19

計算論的学習 (PAC学習) (5)

- ここでいくつかの疑問：
 - 正解に必要な訓練データって最低いくつぐらい？ また訓練データに正解しても X 全体に対する正解の割合は分からないのでは？ (今後これは省略)
 - データが複雑で大量でも区別できるの？ その目安は？
 - 精度を上げる (一致する確率を高くする) 学習アルゴリズムは難しい。それなら、誤差が大きめの弱い学習器をいくつか合わせて、精度の高い強力な学習器はできるか？

20

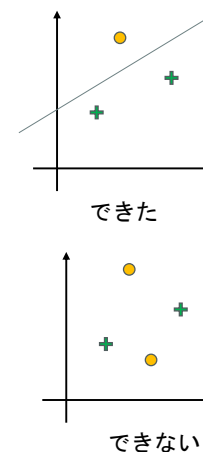
VC次元

(Vapnik-Chervonenkis dimension)

VC次元、次元の呪い

VC次元

- PAC学習の疑問2
 - 訓練データが与えられたとして、いつも区別（つまり判定する関数を (PAC) 学習）できるのか？
 - 前の式では $d_X(h^*, h) - d_S(h^*, h)$ は、 $|H|$ に依存（標本での一致度と全体での一致度の差）。でも H が無限なら？ (S は標本)
 - 直線（超平面）で区別することを考えよう。
 - 2次元：3点なら常にできる（一つの直線に乗るような場合は除く）
 - 4点以上は無理なことがある。
 - では3次元では？一般に n 次元では？
 - （超）平面で分ける場合には、 $n+1$ の点まで分類できることが分かっている。



VC次元

- X 上の判別関数からなる仮説集合を H と表し、 X の部分集合（サンプル）を S とおく。 S を $h \in H$ の値が +1 か -1 で判別することを考える。 $|S| = N$ とおく。
 - $C_+(h) = \{x \in X \mid h(x) = 1\}$ とおく。
 - S を 2 分割する方法は 2^N 通り。 S の全ての部分集合の集まりを 2^S と書くことがある。したがって $|2^S| = 2^N$ 。
 - もし H に含まれる判別関数で S の全ての分割方法が表せるとき H は S を **細分化できる (shatter)** という。これは、 $\Pi_H(S) = \{C_+(h) \cap S \mid h \in H\}$ とおいたとき、 $\Pi_H(S) = 2^S$ となることである。
 - H の VC 次元、 $VCdim(H)$ を以下のように定義する。
 $VCdim(H) = \max\{|S| \mid S \subset X, S \text{ は } H \text{ によって細分化できる}\}$
 最大値であることに注意。
- つまり、全ての分割が表せる標本の最大数。

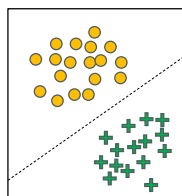
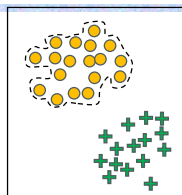
VC次元の例と訓練データ数の目安

- H を 2 次元実数空間 \mathbb{R}^2 上の平面で区分する関数とする。このとき $VCdim(H) = 3$ である。
- 一般に H を n 次元実数空間 \mathbb{R}^n 上の超平面 (n 変数の 1 次式) の集合とすると、 $VCdim(H) = n + 1$ となる。
- 参考：
 - $VCdim(H) = d$ とする。 $h^* \in H$ のサンプル S に矛盾しない $h \in H$ を出力するアルゴリズムが存在するなら、 h^* は PAC 学習可能で、必要訓練データ数 m は

$$m \geq \max\left(\frac{4}{\delta} \log \frac{2}{\delta}, \frac{8d}{\epsilon} \log \frac{13}{\epsilon}\right)$$
 であることが知られている。Log の低は 2。

VC次元の大小

- VC次元は大きい方（よってデータの次元が高い方）が判別能力が上がるのでよい？
=> そうともいえない。
- 学習の汎化（一般化）の観点
 - 訓練データ（教師信号）だけではなく、そこにないデータもある程度適切に区分して欲しい。
 - 次元が上がるとパラメータの組み合わせは指数的に増えるので、複雑なことも高次元の複雑な関数で表せる。あまりにも厳密となり、訓練データ以外のデータはうまく判定できなくなる（右図：どちらがよさそう？）=>次元の呪い（curse of dimensionality. いろいろな意味がある）
 - データ含まれる誤差（雑音）も学習する。
 - 過学習 (overfitting) の問題



26

VC次元の大小

- 判別は可能でVC次元がなるべく小さい空間で判別する。
 - 過学習を減らし、確率的に正解に近づける。
- まずは超平面で判別することを考える

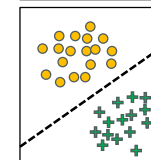
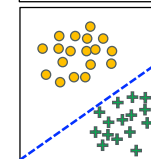
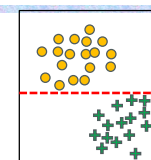
28

29 SVM(1) 線形分類 ハードマージン

SVMは、1960年代に考案されました。ここではその当時の線形分類に関して説明します。特に、マージン最大化とVC次元の最小化の概念について。

ハードマージンと区分する境界線

- 直感的な説明
 - 例外なく分離可能とする（ハードマージン）
 - 右図の赤青黒の区分線は、どれがよい？
 - 赤線はその近くまでサンプルデータがあり、境界がアンバランス。マージンが狭い
 - 青線は+によりすぎている。マージンが偏っている
 - 黒は両者の真ん中当たりの広い空間で分類している（マージンが共に大きく、バランスがとれている）
- 黒点線の境界線のようなマージンを決定する超平面の決定は？



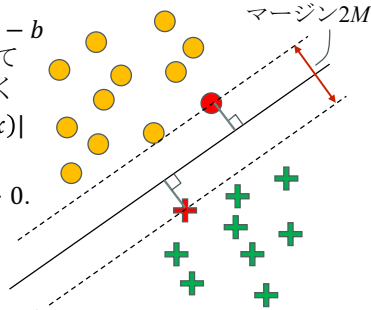
30

マージン（余白）最大化 (1)

- 分類できる平面のうち、その平面までの距離の最小値が、両者とも最大化（バランス）するように平面を決める。
- 両者の点からの距離の最小値を最大化すればよい。（一方から離れすぎると、他方のサンプルからの距離が近くなり、最大化とならない）。
- サンプルを $S \subset \mathbb{R}^n$ 、平面を $f(\mathbf{x}) = \mathbf{w}\mathbf{x}^T - b$ とおく。定数倍して $\|\mathbf{w}\| = 1$ と仮定して良い。 $f(\mathbf{x})$ による判別関数を $h(\mathbf{x})$ とおく

$$\max_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \min_{\mathbf{x} \in S} \frac{|f(\mathbf{x})|}{\|\mathbf{w}\|} = \max_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \min_{\mathbf{x} \in S} |f(\mathbf{x})|$$
 であり、さらに
 $f(\mathbf{x})$ は正確に分類できるので $f(\mathbf{x})h(\mathbf{x}) > 0$ 。
 特に

$$f(\mathbf{x})h(\mathbf{x}) - M \geq 0$$
- ここで M はマージン。マージンを決定する点線上の点をサポートベクトルとよぶ。



32

マージン（余白）最大化 (2)

- この値 M が 1 となるように線形変換（拡大・縮小）する。
 \mathbf{w} を $\hat{\mathbf{w}} = \mathbf{w}/M$, b を $\hat{b} = b/M$, $\hat{f}(\mathbf{x}) = f(\mathbf{x})/M$ と変換。 $\hat{f}(\mathbf{x})$ の最小値は 1 なので、

$$\max_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \min_{\mathbf{x} \in S} \frac{|f(\mathbf{x})|}{\|\mathbf{w}\|} = \max_{\hat{\mathbf{w}} \in \mathbb{R}^n, \hat{b} \in \mathbb{R}} \frac{1}{\|\hat{\mathbf{w}}\|} \min_{\mathbf{x} \in S} |\hat{f}(\mathbf{x})| = \max_{\hat{\mathbf{w}} \in \mathbb{R}^n, \hat{b} \in \mathbb{R}} \frac{1}{\|\hat{\mathbf{w}}\|}$$
 ただし (s.t.)

$$\hat{f}(\mathbf{x})h(\mathbf{x}) - 1 \geq 0 \quad \text{for } \forall \mathbf{x} \in S$$
- $\hat{f}(\mathbf{x})$ と $f(\mathbf{x})$ は同じ平面。 $\hat{\mathbf{w}}, \hat{b}$ を単純に \mathbf{w}, b と表すと、以下の $\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}$ を求めればよい（ $\|\mathbf{w}\|$ より $\|\mathbf{w}\|^2 = \mathbf{w}\mathbf{w}^T$ の方が容易）

$$\min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 \quad (7.1)$$
 s.t. $f(\mathbf{x})h(\mathbf{x}) = |f(\mathbf{x})| \geq 1 \quad \text{for } \forall \mathbf{x} \in S$
 ただし、 $f(\mathbf{x}) = \mathbf{w}\mathbf{x}^T - b$ 。上記は 2 次計画問題となる。
- 式(7.1)は min だが、マージン最大化の制約であることに注意。

33

34 SVM (2) 線形分類 ソフトマージン

プログラムを書くときはライブラリなどを使うだけになると思いますが、それには簡単ではありませんが自然な流れと、厳密な証明があります。

数学の準備(1)

- ラグランジュの未定乗数法のやや一般化（再掲）
- $f(x_1, \dots, x_n), g_k(x_1, \dots, x_n)$ は関数 ($1 \leq k \leq m$)。 $\mathbf{x} = (x_1, \dots, x_n)$ とおく。条件 $g_k(x_1, \dots, x_n) = 0$ のもと、 $f(x_1, \dots, x_n)$ を最小化（最大化）する解 $\mathbf{c} = (c_1, \dots, c_n)$ は、

$$L(x_1, \dots, x_n, \lambda_1, \dots, \lambda_m) = f(x_1, \dots, x_n) - \sum_{k=1}^m \lambda_k g_k(x_1, \dots, x_n)$$

としたとき、

$$\frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1} = \dots = \frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_n} = \frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} = \frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_m} = 0 \quad (2)$$

の解（つまり停留点、極値）でもある。

➔ これをさらに一般化

35

数学の準備(2)

- $f(x_1, \dots, x_n), g_k(x_1, \dots, x_n), h_j(x_1, \dots, x_n)$ は関数 ($1 \leq k \leq m, 1 \leq j \leq l$)。 $\mathbf{x} = (x_1, \dots, x_n)$ とおく。条件 $g_k(x_1, \dots, x_n) = 0, h_j(x_1, \dots, x_n) \geq 0$ のもと、 $f(x_1, \dots, x_n)$ を最小化する解 $\mathbf{c} = (c_1, \dots, c_n)$ を求めるとする。このときラグランジュ関数とは、

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\xi}) = f(\mathbf{x}) - \sum_{k=1}^m \lambda_k g_k(\mathbf{x}) - \sum_{j=1}^l \xi_j h_j(\mathbf{x})$$

と定義する。

- なお、 $h_k(x_1, \dots, x_n)$ にマイナスをつけて、不等号の向きを変えると、上記の式を以下のように表すこともある。

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\xi}) = f(\mathbf{x}) + \sum_{k=1}^m \lambda_k g_k(\mathbf{x}) + \sum_{j=1}^l \xi_j h_j(\mathbf{x})$$

36

数学の準備(3)

- Kuhn-Tucker theorem (クーン・タッカーの定理)

$f(x_1, \dots, x_n), g_k(x_1, \dots, x_n), h_j(x_1, \dots, x_n)$ は関数 ($1 \leq k \leq m, 1 \leq j \leq l$)。 $\mathbf{x} = (x_1, \dots, x_n)$ とおく。条件 $g_k(x_1, \dots, x_n) = 0, h_j(x_1, \dots, x_n) \geq 0$ のもと、 $f(x_1, \dots, x_n)$ を最小化する解 $\mathbf{c} = (c_1, \dots, c_n)$ を求めるとする。このとき $\mathbf{x} = \mathbf{w}^*$ が最適解となる必要十分条件は、ラグランジュ関数 $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\xi})$ に対し、

$$\frac{\partial L(\mathbf{w}^*, \boldsymbol{\lambda}, \boldsymbol{\xi})}{\partial \mathbf{x}} = 0,$$

$$\frac{\partial L(\mathbf{w}^*, \boldsymbol{\lambda}, \boldsymbol{\xi})}{\partial \boldsymbol{\lambda}} = 0,$$

$$\xi_j h_j(\mathbf{w}^*) = 0, h_j(\mathbf{w}^*) \geq 0, \xi_j \geq 0$$

となる $\boldsymbol{\lambda}, \boldsymbol{\xi}$ の解が存在することである。

37

ソフトマージン

- $\forall \mathbf{x} \in S$ について $f(\mathbf{x})h(\mathbf{x}) \geq 1$ とする仮定は強すぎ
 - データの雑音。実世界ではそれほど綺麗なデータはない。
- $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ とし $(\mathbf{x}_k, y_k) = (\mathbf{x}_k, h(\mathbf{x}_k)) \in \mathbb{R}^n \times \{1, -1\}$ と表す。上記の仮定を以下のように緩和する

$$y_k f(\mathbf{x}_k) \geq 1 - \zeta_k, \quad \zeta_k > 0 \text{ for } 1 \leq \forall k \leq m$$
 - この場合、「マージンの間」や「マージンを超えて異なるクラスの点」も存在することになる。
 - そうすると最大マージンって何？という疑問。
 - 上記の式は、 $f(\mathbf{x}) = \mathbf{w}\mathbf{x}^T - b = 1, f(\mathbf{x}) = -1$ の間をマージンとしている。
 - 分類にエラー（誤分類）があるのは $\zeta_k > 1$ のときなので、 $\sum \zeta_k \leq K$ (K は正の整数) ならエラーとなる点は高々 K 個となる。これも小さい方がよい。

38

ソフトマージンによる分類 (1)

- (7.1)を含め、これらを合わせて以下の最適化問題（2次計画問題）として表す。

$$\min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}, \zeta \in \mathbb{R}^m} \left(\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{k=1}^m \zeta_k \right)$$

s.t.

$$y_k (\mathbf{w}\mathbf{x}_k^T + b) \geq 1 - \zeta_k, \\ \zeta_k \geq 0 \quad \text{for } 1 \leq \forall k \leq m$$

- ここで C は正則化係数 (regularization parameter) という。マージンはなるべく大きい方がよいが、エラー（誤分類）は少ない方がよい。マージンの大きさと誤分類のバランスをとる。
 - C が大きいとエラーが大きく作用し、誤分類を許しにくくなる。（特に C が極めて大きいときは事実上のハードマージンの分類となる）。
 - 逆に C が小さいと、エラーを許容し、マージンを大きく取る。
 - C の値は実験的に求める（→ 交差検定 (cross validation)）

39

ソフトマージンによる分類 (2)

- 前の式からラグランジュ関数 L を定義。以下の最適化問題 (2次計画問題) として表す (Kuhn-Tucker theorem)

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\xi}') = \frac{1}{2} \mathbf{w} \mathbf{w}^T + C \sum_{k=1}^m \xi_k - \sum_{k=1}^m \xi_k (\xi_k - 1 + y_k (\mathbf{w} \mathbf{x}_k^T + b)) - \sum_{k=1}^m \xi'_k \xi_k$$

かつ $\boldsymbol{\xi}, \boldsymbol{\xi}' \geq \mathbf{0}$

- このうち $\boldsymbol{\xi}, \boldsymbol{\xi}'$ で微分して0となるものは後回し。まずは
- $\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{k=1}^m \xi_k y_k \mathbf{x}_k = \mathbf{0} \Rightarrow \mathbf{w} = \sum_{k=1}^m \xi_k y_k \mathbf{x}_k$
 - たとえば、 $\frac{\partial L}{\partial w_i} = w_i - \sum_{k=1}^m \xi_k y_k x_{ki}$
- $\frac{\partial L}{\partial b} = \sum_{k=1}^m \xi_k y_k = 0$
- $\frac{\partial L}{\partial \xi} = C - \xi - \xi' = 0$ 。これらを代入。 ➡ この式から上の斜線

40

ソフトマージンによる分類

- $L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\xi}')$ を最大化する $\boldsymbol{\xi}$ を求める。

$$\begin{aligned} \max_{\boldsymbol{\xi}} L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\xi}') \\ &= \max_{\boldsymbol{\xi}} \left(\frac{1}{2} \sum_{k=1}^m \xi_k y_k \mathbf{x}_k \sum_{j=1}^m \xi_j y_j \mathbf{x}_j^T + \sum_{k=1}^m \xi_k - \sum_{k,j=1}^m \xi_k y_k \xi_j y_j \mathbf{x}_j \mathbf{x}_k^T \right) \\ &= \max_{\boldsymbol{\xi}} \left(\sum_{k=1}^m \xi_k - \frac{1}{2} \sum_{k,j=1}^m \xi_k y_k \xi_j y_j \mathbf{x}_j \mathbf{x}_k^T \right) \end{aligned}$$

s.t. $0 \leq \xi_k \leq C, \sum_{k=1}^m \xi_k y_k = 0$

- 注: $\sum_{k=1}^m \xi_k y_k = 0$ から b の項は消える。 $C - \xi = \xi' \geq 0$ から $C \geq \xi (\geq 0)$.
- この問題を [サポートベクターマシンの双対問題](#) といい、これを解く。

41

ソフトマージンによる分類

- $\boldsymbol{\xi}$ の値によりサンプルを分類できる (正式にはもう少し議論は必要だが詳細は省く。多分、最適化などの大学院の講義で)。
- $0 \leq \xi_k < C$ のときは、 $y_k (\mathbf{w} \mathbf{x}_k^T + b) \geq 1$ であり、そのサンプルを正しく分類できている。特に、 $\boldsymbol{\xi} = \mathbf{0}$ の点は境界線上にあり (スライド「マージン (余白) 最大化 (1)」の点線。 $y_k (\mathbf{w} \mathbf{x}_k^T + b) = 1$)、この境界線を定める点となる。
- $\xi_k = C$ のときは $y_k (\mathbf{w} \mathbf{x}_k^T + b) < 1$ となりマージンを超え (つまりマージン内だったり、相手側のマージンの向こう側にある。この点については、正しく分類できない。
- RにSVMのパッケージがあるので、興味があれば、各自調べてみてください。

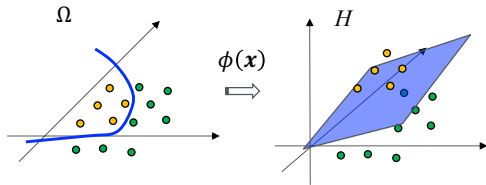
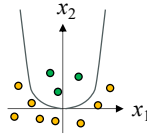
42

43

SVM (3) 非線形への拡張

高次元への写像

- 線形の超平面での分類は限度がある。
- 高次元に写像（下図）。高次元であればVC次元も大きくなり埋め込むことができる。
 - 右図： $(x_1, x_2) \rightarrow (x_1, x_1^2, x_2) = (y_1, y_2, y_3)$
 - $x_1 - x_1^2 + x_2 = 0$ （符号が逆かも...）
 - 次元が高くなれば、計算コストも高くなる（これも次元の呪い）
 - 次元の呪い（過学習）も考慮すべき。何らかの制限が必要だが、最大マージンという制限が次元の呪いへの対応（の一部）となっている。



44

高次元への写像

- アイデア：非線形の空間を高次元の空間へ \mathbb{R}^n の要素を埋め込み、その空間で線形分離を試みる。
- $f(\mathbf{x}) = \mathbf{w}\mathbf{x}^T - b$ に $\mathbf{w} = \sum_{k=1}^m \xi_k y_k \mathbf{x}_k$ （スライド39）を代入すると、

$$f(\mathbf{x}) = \sum_{k=1}^m \xi_k y_k \mathbf{x}_k \mathbf{x}^T - b$$

となり、超平面もまたその判別関数も内積の形式になる。

- 事象の空間 $X \subset \mathbb{R}^n$ とサンプル $S \subset \mathbb{R}^n$ とする。高次元への写像

$$\phi(\mathbf{x}): X \rightarrow \mathbb{R}^{n'}$$

- 高次元に写像しても、内積で計算ができるような変換なら、同様に計算できる。これをカーネルトリックと呼ぶ。
- つまり、 $\phi(\mathbf{x}_i)\phi(\mathbf{x}_j)^T = K(\mathbf{x}_i, \mathbf{x}_j)$ と高次元での内積だけから計算できるとき簡単になる可能性がある。

45

正定値カーネル

- 定義
 - 集合 Ω 、関数 $K: \Omega \times \Omega \rightarrow \mathbb{R}$ が正定値カーネルとは以下の2条件を満たすこと。

- $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$
- Ω 上の n 個の任意の要素 $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ に対し $n \times n$ 行列

$$\begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_n, \mathbf{x}_1) & \cdots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

は正定値行列である。

- 正定値カーネルを利用して高次元に埋め込む。
 - $\phi(\mathbf{x}) = K(\cdot, \mathbf{x})$ と定義すれば良い。このような写像を使って高次元に埋め込む。

46

正定値カーネル

- 前出の $L(\mathbf{w}, b, \boldsymbol{\zeta}, \boldsymbol{\xi}, \boldsymbol{\xi}')$ をは以下のようにする。

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\zeta}, \boldsymbol{\xi}, \boldsymbol{\xi}') &= \sum_{k=1}^m \xi_k - \frac{1}{2} \sum_{k,j=1}^m \xi_k y_k \xi_j y_j \phi(\mathbf{x}_j) \phi(\mathbf{x}_k)^T \\ &= \sum_{k=1}^m \xi_k - \frac{1}{2} \sum_{k,j=1}^m \xi_k y_k \xi_j y_j K(\mathbf{x}_j, \mathbf{x}_k) \end{aligned}$$

- と同様な形式となる。
 - ただし、上記を得るにはヒルベルト空間（内積が定義されているベクトル空間）の性質を使う。難しくはないが、長くなるので、ここでは省略する。

47

SVMの利点と欠点

利点

- 分類の能力は高い。高速。
- 線形の場合は学習結果の解釈が比較的容易。
- 最大マージンという性質で、過学習に対応している。
- 自動的に線形・非線形の特徴抽出ができる
- 複雑な非線関数も近似的に分類可能

欠点

- 線形分離だけでは能力に限界はある。
- 学習結果の解釈は容易ではない。特に非線形の場合
- 非線形ではデータ数が多くなると（たとえば10000以上）急速に遅くなる。

48

練習問題

課題8-1：

- ソフトマージンとハードマージンについて、違いを述べよ。

課題8-2：

- 次元をあげるとVC次元があがり、判別能力があがるが、学習面からみると、判別能力の向上は必ずしも望ましくはない。その理由について議論せよ。

49

50

アンサンブル学習と Adaboost

Ensemble learning and adaptive boosting

アンサンブル学習

- 考え方：精度の低い弱い学習器（弱学習器）を組み合わせ、精度の良い学習器を作りたい（PAC学習の疑問3）。これをアンサンブル学習という
- ここで想定する弱学習器：
 - 多くの場合は、線形分離やあまり深くない決定木などを使う。また、訓練データの一部のみからの学習結果もある。
 - 決定木の場合の極端な例としては深さ1の決定木（木ほど高くないので決定株とも言われる）に制限する方法もある。
- ただし研究の発端は、弱PAC学習アルゴリズムという概念を定義して、これを通常のPAC学習アルゴリズムに理論的に押し上げることが狙って研究されたのが（あとで述べる）ブースティングの始まり。

51

参考：弱PAC学習

- 帰納的関数のクラス H を弱PAC学習可能とは、ある定数 $\epsilon_0 < \frac{1}{2}$ と δ_0 が存在して、任意の判別関数 $h^* \in H$ と X 上の任意の確率分布 \mathcal{D} に対して、多項式時間の学習アルゴリズム \mathcal{L} が存在し、ついて
$$\Pr(d_X(h^*, h) \leq \epsilon_0) \geq 1 - \delta_0$$
となる $h \in H$ を出力し停止すること。[Kearns and Valiant 89]
- (PAC学習との比較) 帰納的関数のクラス H をPAC学習可能とは、判別関数 $h^* \in H$ と X 上の任意の確率分布 \mathcal{D} に対して、多項式時間の学習アルゴリズム \mathcal{L} が存在し、任意の $0 \leq \delta, \epsilon \leq 1$ について
$$\Pr(d_X(h^*, h) \leq \epsilon) \geq 1 - \delta$$
となる $h \in H$ を出力し停止すること。[Valiant84]

52

主なアンサンブル学習の例

- バギング (Bagging, L.Breiman 1996)
 - m 組のサンプルを抽出し、それぞれのサンプルから m 個の弱学習器を作る。並列化が可能。
 - 組み合わせは、例えば多数決 (voting)とか平均値をとるなど。
 - ランダムフォレスト (前出) は、広い意味で一例
- ブースティング (Boosting, Kearns and Valiant 1989)
 - Adaboost (Freund and Schapire, 1997): 弱学習器で誤差のある訓練データの重みを大きくし、次に(重点的に)学習する。[つまり弱点を強化]。これを繰り返して学習器の修正を積み重ねる。逐次的。
 - Gradient Boosted Decision Tree (GBDT): 弱学習器として決定木を利用。誤差に基づいてサンプルの重みを変える考え方は同じ。逐次的。
- ブースティングの代表的アルゴリズムであり、よく使われるAdaBoost (adaptive boosting) について概説する。

53

もう一度確認

- 前のスライドをもう一度確認 (+ちよつと修正)
 - 訓練データ (教師信号) は、 X 上のある確率分布 \mathcal{D} に従ってその要素が選択される。
 - \mathcal{D} に従って、 $x \in X$ が選択される確率を $\mathcal{D}(x)$ ($= P_x$)と置く。確率分布なので、 $\sum_{x \in X} \mathcal{D}(x) = 1$ である。
 - データ正負を決める真の判別関数を $h^* \in H$ と表す。仮説集合 H からアルゴリズム \mathcal{L} は候補 h を出力する。
 - h^* と h の誤差 (差分、距離) を
$$d_X(h^*, h) = \sum_{x \in X, h^*(x) \neq h(x)} \mathcal{D}(x)$$
と定義する。距離を d_X と書くことにした。
 - 学習アルゴリズムに与えられる入力 $(x_1, y_1) \dots (x_m, y_m)$ ここで $x_i \in X, y_i = h^*(x_i) \in \{-1, 1\}$ 。なるべく正確に分類したい。

54

Boosting

- データ X 上の初期分布 \mathcal{D}_1 を等確率とする ($\mathcal{D}_1(x_j) = 1/m$)とする。これに従ってサンプルを選択。 $t = 1$ 。
 1. \mathcal{D}_t に基づくサンプルからアルゴリズム \mathcal{A} が弱学習器が h_t を出力する。
 2. $\epsilon_t = d_X(h^*, h_t)$ を求める。
 3. ϵ_t に基づいて「 h_t の重要度」と「新しい X 上の分布 \mathcal{D}_{t+1} 」を定義して1へ戻る。
 4. 上記を一定回数 (T 回) 繰り返し、得られた出力 h_1, h_2, \dots, h_T を組み合わせる最終出力を $H(x)$ とする。
- アンダーラインの部分がboostingの方法ごとに違う。
- アルゴリズム \mathcal{A} には学習アルゴリズムがくる。

55

Adaboost

- 訓練データ X 上の初期分布 \mathcal{D}_1 を等確率とする。つまり、 $\mathcal{D}_1(\mathbf{x}_j) = 1/m$ とする。 $t = 1$ 。
 1. \mathcal{D}_t に基づくサンプルからアルゴリズム \mathcal{A} が弱学習器が h_t を出力
 2. $\varepsilon_t = d_X(h^*, h_t)$ // 誤差 (エラー) を計算
 3. $\alpha_t = \frac{1}{2} \ln \frac{1-\varepsilon_t}{\varepsilon_t}$ をエラーに基づく h_t の重要度とする。
 4. $\mathcal{D}_{t+1}(\mathbf{x}_j) = \frac{\mathcal{D}_t(\mathbf{x}_j) \exp(-\alpha_t y_j h_t(\mathbf{x}_j))}{Z_t}$ と定義。ここで Z_t は $\sum_{\mathbf{x} \in X} \mathcal{D}_{t+1}(\mathbf{x}) = 1$ となるように正規化するもの。
 5. $H(\mathbf{x}) = \text{sgn} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$ // これは重み付きの多数決。

56

Adaboost

- Adaboostの計算は何をしているのか？(1)
 - まず \mathcal{D}_{t+1} から。 $\sum_{\mathbf{x}_j \in X, y_j h_t(\mathbf{x}_j) = -1} \mathcal{D}_t(\mathbf{x}_j) \exp(-\alpha_t y_j h_t(\mathbf{x}_j))$
 $= \sum \mathcal{D}_t(\mathbf{x}_j) \exp(\alpha_t) = \sum \mathcal{D}_t(\mathbf{x}_j) \sqrt{\frac{1-\varepsilon_t}{\varepsilon_t}} = \sqrt{\varepsilon_t(1-\varepsilon_t)}$
 - $\sum_{\mathbf{x}_j \in X, y_j h_t(\mathbf{x}_j) = 1} \mathcal{D}_t(\mathbf{x}_j) \exp(-\alpha_t y_j h_t(\mathbf{x}_j))$
 $= \sum \mathcal{D}_t(\mathbf{x}_j) \exp(-\alpha_t) = \sum \mathcal{D}_t(\mathbf{x}_j) \sqrt{\frac{\varepsilon_t}{1-\varepsilon_t}} = \sqrt{\varepsilon_t(1-\varepsilon_t)}$
 - 正しかった部分とエラーの部分が等確率 (1/2) となるように分布を決めている。
 - これは h_t がランダムになるようにわざと再定義して、次の学習している。
 (一方、例えば $\varepsilon_t = 0.9$ とすると、これは精度良くエラーを出している
 ので-1倍すれば、精度良い学習器になる。上記のようにするのが、前の
 学習器にとっては不利なものとなる)

57

Adaboost

- Adaboostの計算は何をしているのか？(2)
 - $\alpha_t h_t(\mathbf{x}) = \frac{1}{2} \ln \frac{1-\varepsilon_t}{\varepsilon_t} \times h_t(\mathbf{x})$
 - 直感的には ε_t (逆に $1-\varepsilon_t$) が小さいものほど重みを大きくして強く反映している。
 - $\varepsilon_t = 0.5$ で $\alpha_t = 0$ となる。ランダムに相当し、情報がないということ。
 - 仮に $\varepsilon_t > 0.5$ とすると、
 $\alpha_t h_t(\mathbf{x}) = \frac{1}{2} \ln \frac{1-\varepsilon_t}{\varepsilon_t} \times h_t(\mathbf{x}) = \frac{1}{2} \ln \frac{1-(1-\varepsilon_t)}{1-\varepsilon_t} \times (-h_t(\mathbf{x}))$
 となり、 $-h_t(\mathbf{x})$ を使う。これならエラーは $1-\varepsilon_t < 0.5$
 - それまでの弱点を補強する形で、アルゴリズム \mathcal{A} が新しい学習器を見つける。

58

Adaboostの利点と欠点

利点

- 実装が容易。処理も高速
- パラメータが少ない (T ぐらい)
- 弱学習器との独立性が高く、うまく活用できる。
- 汎化能力は高い
- 過学習 (overfitting) が他の方法に比べて起こりにくいと言われる

欠点

- 解の最適性に難あり
- ノイズデータの影響を受けやすい。過学習になることも。
- 弱学習器に対してデータが複雑な場合、判別がうまくいかなかったり、逆に過学習になることもある。

60

Ensemble Learning文献

- Breiman, L. “Bagging Predictors,” *Machine Learning* (1996) 24: 123.
<https://doi.org/10.1023/A:1018054314350>
- M.J. Kearns and L. Valiant. “Cryptographic limitations on learning Boolean formulae and finite automata,” *Proc. of the 21th ACM STOC*, pages 433–444, 1989.
- Yoav Freund and Robert E Schapire. 1997. “A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting,” *J. Comput. Syst. Sci.* 55.