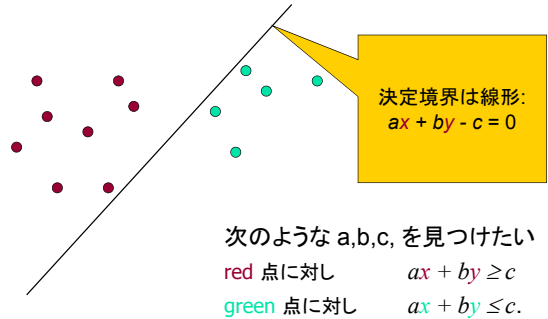


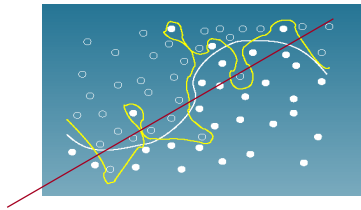
情報意味論(11) サポートベクターマシン

理工学部管理工学科
櫻井彰人

基礎的復習: 線形判別関数



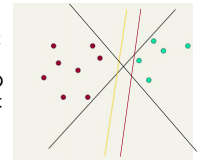
これも復習: 複雑な境界は?



Christopher Manning のスライドから

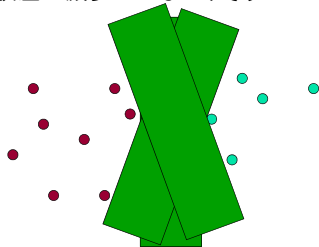
どの超平面を選ぶべきか?

- a, b, c にはいくつもの可能性あり
- 見つけたどれもが最良なわけではない
[何か「よさ」の基準を設ける必要がある]
 - パーセプトロン学習アルゴリズムではどうであったか?
- サポートベクターマシンは「最良」のものをみつける。
 - 超平面とそれに近い「困難点」との距離を最大化する
 - 直感的解釈: 決定境界に近いところに(別のクラスの)点がなければ、決定の不確かさは少なからう



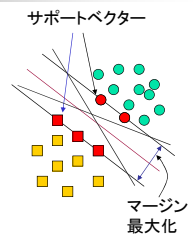
直感的解釈をもう一つ

- 分離境界を幅のある帯に置き換えてみよう。この幅が狭いときには、選択範囲がせばまり、汎化誤差の減少につながりそう



サポートベクターマシン (SVM)

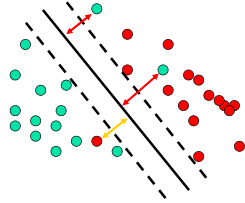
- SVM は、分離超平面周囲のマージンを最大化する。
 - ラージマージン分類器ともいう
- 決定関数はサポートベクターと呼ばれる訓練データによって完全に定まる。
- 2次計画問題である
- 広範囲の問題に対してうまくいく方法であると考えられている



ラージマージン分類器

線形分離可能でないならば

- 誤りを許す
 - コストを払って、本来あるべき場所に動かす
- ただ、超平面はどちらのクラスからも遠ざける

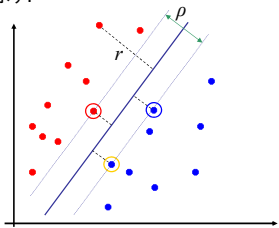


最大マージン: 定式化

- w : 決定超平面への垂線ベクトル
 - x_i : i 番目のデータ点
 - y_i : 属するクラス (+1 or -1) 注: 1/0 ではない
 - 分類器: $\text{sign}(w^T x_i + b)$
 - そのとき x_i の関数マージン: $y_i (w^T x_i + b)$
 - 勿論 w を大きくすればマージンは増大する、そこで、
- (訓練データ全体の関数マージンは、上記の値の最大値)

幾何的マージン

- データ点から分離超平面までの距離 $r = \frac{w^T x + b}{\|w\|}$
- 分離超平面に最も近い点がサポートベクター。
- 分離超平面のマージン ρ は相異なるクラスのサポートベクターがどの程度分離しているかを示す。



線形 SVM を数学的に

- 全ての点が超平面から関数値で 1 離れていると仮定しよう。そうであれば次の2つの制約が訓練データ集合 $\{(x_i, y_i)\}$ から得られる

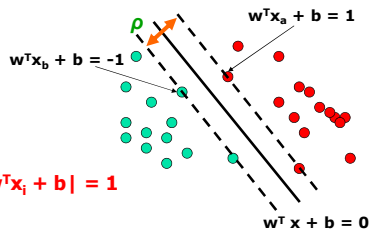
$$w^T x_i + b \geq 1 \quad \text{if } y_i = 1$$

$$w^T x_i + b \leq -1 \quad \text{if } y_i = -1$$

- サポートベクターに対しては、上記不等式は等式となる; そうすると、各データの超平面からの距離は $r = \frac{w^T x + b}{\|w\|}$ であるから、マージンは次の値となる: $\rho = \frac{2}{\|w\|}$

線形サポートベクターマシン

- 超平面 $w^T x + b = 0$
- 制約: $\min_{i=1, \dots, n} |w^T x_i + b| = 1$
- 書換えると: $w^T(x_a - x_b) = 2$
 $\rho = \|x_a - x_b\|_2 = 2 / \|w\|_2$



線形サポートベクターマシン

- 次の2次計画問題が得られる:

次のような w と b を見出せ:
 $\rho = \frac{2}{\|w\|}$ は最大であり、全ての $\{(x_i, y_i)\}$ につき
 $w^T x_i + b \geq 1$ if $y_i = 1$; $w^T x_i + b \leq -1$ if $y_i = -1$

- よりよい定式化 ($\min \|w\| = \max 1 / \|w\|$):

次のような w と b を見出せ:
 $\Phi(w) = \frac{1}{2} w^T w$ は最小であり、すべての $\{(x_i, y_i)\}$ につき
 $y_i (w^T x_i + b) \geq 1$

最適化問題の解法

次のような w と b を見出せ
 最小化: $\Phi(w) = \frac{1}{2} w^T w$;
 全ての $\{(x_i, y_i)\}$ につき: $y_i (w^T x_i + b) \geq 1$

- 線形制約のもとでの2次関数の最適化
- 2次計画問題は、よく知られた数理計画問題の一つ。多くの解法が知られている
- 解法にあたっては、ラグランジュ乗数 a_i を主問題の各制約に割付けた双対問題を構成する:

次のような $a_1 \dots a_N$ を見出せ
 最大化: $Q(a) = \sum a_i - \frac{1}{2} \sum \sum a_i a_j y_i y_j x_i^T x_j$;
 (1) $\sum a_i y_i = 0$
 (2) $a_i \geq 0$, 任意の a_i

主問題のラグランジアンは

$$L(w, b, a) = \frac{1}{2} w \cdot w - \sum_i a_i (y_i (w \cdot x_i + b) - 1)$$

従って、

$$\frac{\partial L(w, b, a)}{\partial w} = w - \sum_i a_i y_i x_i, \quad \frac{\partial L(w, b, a)}{\partial b} = \sum_i a_i y_i$$

となるゆえ、停留点は、

$$w = \sum_i a_i y_i x_i, \quad 0 = \sum_i a_i y_i$$

これらを主問題に戻せば

$$\begin{aligned} L(w, b, a) &= \frac{1}{2} w \cdot w - \sum_i a_i (y_i (w \cdot x_i + b) - 1) \\ &= \frac{1}{2} \sum_{i,j} a_i a_j y_i y_j x_i \cdot x_j - \sum_{i,j} a_i a_j y_i y_j x_i \cdot x_j + \sum_i a_i \\ &= \sum_i a_i - \frac{1}{2} \sum_{i,j} a_i a_j y_i y_j x_i \cdot x_j \end{aligned}$$

最適化問題の解法

- 解の形は:

$$w = \sum a_i y_i x_i, \text{ かつ } a_k \neq 0 \text{ なるすべての } x_k \text{ につき } b = y_k - w^T x_k$$

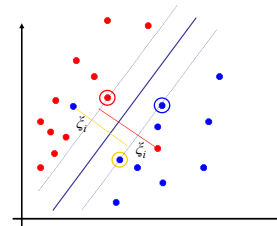
- 各非零の a_i は、対応する x_i がサポートベクターであることを示す。
- 識別関数は次のようになる:

$$f(x) = \sum a_i y_i x_i^T x + b$$

- 当該式は新規点とサポートベクトル x_i の内積であることを注意。
- また、当該最適化問題を解くには、訓練データのすべての組合せに関する内積 $x_i \cdot x_j$ の計算が含まれていることを注意しておく。

ソフトマージン分類器

- もし訓練データが線形分離可能でなければ、スラック変数 ξ_j を用いて分類が難しい点やノイズがのった点の誤分類を許すようにする。



ソフトマージン分類器

- 以前の定式化:

次のような w と b を見出せ
 最小化: $\Phi(w) = \frac{1}{2} w^T w$;
 全ての $\{(x_i, y_i)\}$ について: $y_i (w^T x_i + b) \geq 1$

- スラック変数を含む、新しい定式化:

次のような w と b を見出せ
 最小化: $\Phi(w) = \frac{1}{2} w^T w + C \sum \xi_i$;
 全ての $\{(x_i, y_i)\}$ について: $y_i (w^T x_i + b) \geq 1 - \xi_i$, かつ
 全ての i について: $\xi_i \geq 0$

- パラメータ C は過学習を制御する方法と見ることができる。

ソフトマージン分類器 - 解

- ソフトマージン分類器の双対問題:

次のような $a_1 \dots a_N$ を見出せ:
 最大化: $Q(a) = \sum a_i - \frac{1}{2} \sum \sum a_i a_j y_i y_j x_i^T x_j$; ただし
 (1) $\sum a_i y_i = 0$
 (2) 全ての a_i につき $0 \leq a_i \leq C$

- スラック変数 ξ_i もラグランジュ乗数も、双対問題には表れていない!
- 再び、非零の a_i に対応する x_i はサポートベクターである。
- 当該双対問題への解は:

$$w = \sum a_i y_i x_i$$

$$b = y_k (1 - \xi_k) - w^T x_k \text{ where } k = \operatorname{argmax}_k a_k$$

明示的には w がなくても分類できる!

$$f(x) = \sum a_i y_i x_i^T x + b$$

主問題のラグランジアンは

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_i \xi_i - \sum_i \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - (1 - \xi_i)) - \sum_i \nu_i \xi_i$$

従って、

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i, \quad \frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial b} = \sum_i \alpha_i y_i, \quad \frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \xi_i} = C - \alpha_i - \nu_i$$

となるゆえ、停留点は、

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i, \quad 0 = \sum_i \alpha_i y_i, \quad 0 = C - \alpha_i - \nu_i$$

これらを主問題に戻せば

$$L(\mathbf{w}, b, \mathbf{a}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

ただし、条件は、

$$0 = \sum_i \alpha_i y_i, \quad 0 \leq \alpha_i \leq C \quad (\text{for all } i)$$

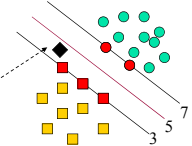
SVMを用いた分類

- 所与の新点 (x_1, x_2) に対し、その超平面への垂直射影の高さを計る (score としよう):
 - 2次元の場合: $\text{score} = w_1 x_1 + w_2 x_2 + b$.
 - すなわち: $\text{score} = \mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$
 - 信頼限度 t を定めよう.

score > t : yes

score < -t : no

それ以外: 判定放棄



線形 SVM: まとめ

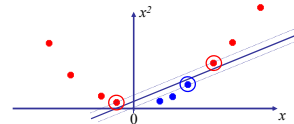
- 分類器は、分離超平面 *separating hyperplane*.
- 最も重要な訓練データ点がサポートベクターとなる; それが当該超平面を決める.
- 2次計画問題を解けば、どの点 \mathbf{x}_i がサポートベクターで非零のラグランジュ乗数 α_i に対応するかが分かる.
- 当該問題の双対問題においても解法においても、訓練データ点は、内積の中しにしか現れない:

次のような $\alpha_1, \dots, \alpha_N$ を見せ:
 最大化: $Q(\mathbf{a}) = \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$, 但し
 (1) $\sum \alpha_i y_i = 0$
 (2) すべての α_i につき: $0 \leq \alpha_i \leq C$

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

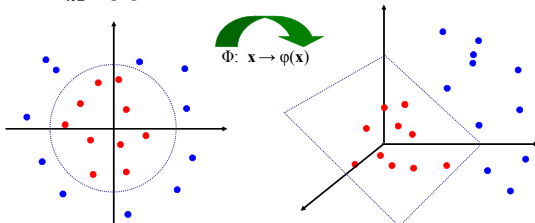
非線形 SVM

- 線形分離可能なデータに対しては、少々ノイズがあっても、うまくいく:
- しかし、データ集合が線形分離可能でなかったらどうしよう?
- 例えば... データをより高次元の空間に写像したらどうだろうか?



非線形 SVM: 特徴空間

- 一般的なアイデア: もともとの特徴空間は、いつでも、ある高次元特徴空間に写像すれば、線形分離可能となる:



高次元空間への写像: 問題点

- 計算時間:
 - データが1001個あれば、(非線形関数で)1000次元に写像すれば、必ず、線形分離できる。
 - しかし、非線形関数の計算は時間がかかるのに、データ1個につき1000回の計算(共通部分を多くして計算するにせよ)が必要では、大変な計算量となる
 - ⇒ カーネル関数の利用(カーネルトリック)による、計算量の大幅な削減
- 汎化能力:
 - データが1001個あれば、(非線形関数で)1000次元に写像すれば、必ず、線形分離できる。
 - それは、すなわち、どんな教師データであっても、それを実現できるということ。すなわち、過学習!!
 - ⇒ この問題の解決こそ、ラージマージン分類器の本領

高次元空間への非線形写像

- データ x から高次元空間 F への(非線形)写像 $\phi(x)$ を考える. $\phi: R^N \rightarrow F$

主問題のラグランジアンは $L(w, b, a) = \frac{1}{2} w \cdot w - \sum_i \alpha_i (y_i (w \cdot x_i + b) - 1)$

$$L(w, b, a) = \frac{1}{2} w \cdot w - \sum_i \alpha_i (y_i (w \cdot x_i + b) - 1)$$

双対問題: 拘束条件付き最大化 $L(w, b, a) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(x_i) \cdot \phi(x_j)$

$$L(w, b, a) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(x_i) \cdot \phi(x_j)$$

$\sum_i \alpha_i y_i = 0$ and $\forall i \alpha_i \geq 0$

カーネルトリック “Kernel Trick”

- 注目すべきは、 $\Phi(x)$ は、 $\Phi(x) \cdot \Phi(y)$ というように、内積でしか表れない。

$$L(w, b, a) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(x_i) \cdot \phi(x_j)$$

- そこで、もし、 $K(x, y) = \Phi(x) \cdot \Phi(y)$ となる、簡単な関数 K があれば、計算が非常に楽になる。
 - 特に、 $K(x, y)$ が x, y の関数であるとさらに楽になる

Mercerの定理

- 関数 K が内積の形で書ける:

$$K(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y)$$

ための必要十分条件は、 K が対称かつ半正定値であることである。i.e.,

$$K(x, y) = K(y, x)$$

$$\iint K(x, y) f(x) f(y) dx dy \geq 0 \quad \text{for any } f$$

なお、 $\phi_i(x)$ は $K(x, y)$ の固有関数となる。i.e.,

$$\int K(x, y) \phi_i(x) dx = \lambda_i \phi_i(y)$$

よく使われるカーネル関数

- 線形カーネル $K(x, y) = x^T y$
 - 多項式カーネル $K(x, y) = (x^T y + 1)^p$ or $(x^T y)^p$
 - RBFカーネル $K(x, y) = \exp(-\|x - y\|^2 / 2\sigma^2)$
 - MLP $K(x, y) = \tanh(\beta_0 x^T y + \beta_1)$ ← 半正定値ではない
- 例: 2次元ベクトル $\mathbf{x} = [x_1, x_2]$ に対し $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$ とおく
このとき、次の式が成立する $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$:
- $$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$$
- $$= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2}$$
- $$= [1 \ x_{i1}^2 \ \sqrt{2} \ x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} \ x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}]$$
- $$= \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$
- ただし $\phi(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} \ x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]$

SVM: 汎化能力の推定

- 汎化能力最大(新規データに対して最も正確)の分類器がほしい。
- 良い汎化性能を得るための糸口は?
 - 訓練データを大きくする
 - 訓練データに対する誤りを小さくする
 - 容量/分散 (モデル記述パラメータ数, モデルの表現能力) を小さくする
- SVM では、これらの量に基づいて、新規データに対する誤差限界を明示的に示すことができる。

容量/分散: VC 次元

- 理論的なリスク限界:

$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{\frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l}}$$
- Risk = 平均誤り率
- α - 当該モデル (ν パラメータで決まる)
- R_{emp} - 経験リスク, l - 観測数, h - VC 次元, 当該式は確率 $(1-\eta)$ で正しい
 - VC (Vapnik-Chervonenkis) 次元/容量: shatterできる点の最大数
 - ある点集合が shatterできるとは、その任意のラベル付けを当該分類器が行えること。
- 重要な理論的性質。しかし、実際にはあまり使われない

例: 超平面のVC次元

- d 次元空間に n 個の点があり、それらは、red か green とラベルが付けられていると仮定する。 n を (d の関数として) どれだけ大きくとれば、red 点と green 点が線形分離でなくなる例が作れるか?
- 例, $d=2$ に対しては $n \geq 4$.



スケッチ: マージン最大化の理論的な正当化

- Vapnik は次のことを証明した:
マージンが ρ 以上の線形判別器クラスの VC 次元 h は、次の上界をもつ

$$h \leq \min \left\{ \left\lceil \frac{D^2}{\rho^2} \right\rceil, m_0 \right\} + 1$$

ただし D は訓練事例をすべて囲い込む最小の超球の直径、そして m_0 は (事例の表現空間の) 次元である。

- 直感的に、これは空間の次元 m_0 にかかわらず、マージン ρ を最大化することにより、VC 次元を最小化することができる。
- こうして、分類器の複雑度は、次元数に関わりなく小さく保つことができる。

Vapnik 1982. Estimation of Dependences Based on Empirical Data. Springer Series in Statistics (Springer Series in Statistics) Springer-Verlag New York, Inc.

SVM の性能

- SVM は、最良の性能を持つと考える人は多い。
- 統計的な有意性は明確な場合もそうでない場合も。
- SVM と同程度の性能をもつ手法は他にもある。
- 例: regularized logistic regression (Zhang & Oles)
 - Tong Zhang, Frank J. Oles: Text Categorization Based on Regularized Linear Classification Methods. Information Retrieval 4(1): 5-31 (2001)
- 比較研究の例: Yang & Liu
 - Yiming Yang, Xin Liu: A re-examination of text categorization methods, 22nd Annual International SIGIR (1999).

評価例: 古典的な Reuters データ

- 非常によく使われたデータセット
- 21578 documents
- 9603 training, 3299 test articles (ModApte split)
- 118 categories
 - 一つの article は複数の category に属しうる
 - 118 個の2値分類
- 1 document 当たりの category 数
 - 1.24
- 10 categories のみ大きい (全 118 categories)

大きめの categories
(#train, #test)

| | |
|----------------------------|-----------------------|
| • Earn (2877, 1087) | • Trade (369, 119) |
| • Acquisitions (1650, 179) | • Interest (347, 131) |
| • Money-fx (538, 179) | • Ship (197, 89) |
| • Grain (433, 149) | • Wheat (212, 71) |
| • Crude (389, 189) | • Corn (182, 56) |

Reuters Text Categorization data set (Reuters-21578) document 例

```
<REUTERS TOPICS="YES" LEWISSPLIT="TRAIN" CGISPLIT="TRAINING-SET" OLDID="12981" NEWID="798">
```

```
<DATE> 2-MAR-1987 16:51:43.42</DATE>
```

```
<TOPICS><D>livestock</D><D>hog</D></TOPICS>
```

```
<TITLE>AMERICAN PORK CONGRESS KICKS OFF TOMORROW</TITLE>
```

```
<DATELINE> CHICAGO, March 2 - </DATELINE><BODY>The American Pork Congress kicks off
```

tomorrow, March 3, in Indianapolis with 160 of the nations pork producers from 44 member states determining industry positions on a number of issues, according to the National Pork Producers Council, NPPC.

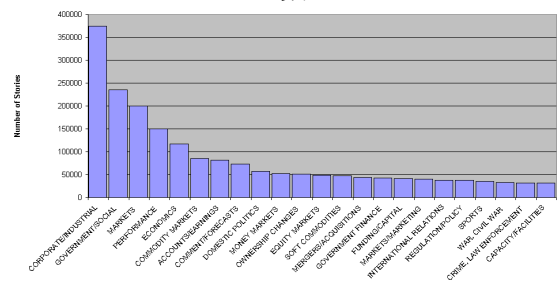
Delegates to the three day Congress will be considering 26 resolutions concerning various issues, including the future direction of farm policy and the tax law as it applies to the agriculture sector. The delegates will also debate whether to endorse concepts of a national PRV (pseudorabies virus) control and eradication program, the NPPC said.

A large trade show, in conjunction with the congress, will feature the latest in technology in all areas of the industry, the NPPC added. Reuter

```
&#3;</BODY></TEXT></REUTERS>
```

New Reuters: RCV1: 810,000 文書

Reuters RCV1 の頻出トピック



Category
<http://about.reuters.com/researchandstandards/corpus/statistics/index.asp>

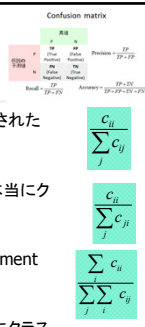
評価尺度

- Recall: クラス i の document 中、正しく i に分類されたものの割合:

$$\frac{C_{ii}}{\sum_j C_{ij}}$$
- Precision: クラス i に分類された document 中、本当にクラス i に属するものの割合:

$$\frac{C_{ii}}{\sum_j C_{ji}}$$
- Accuracy (1- error rate) 正しく分類された document の割合:

$$\frac{\sum_i C_{ii}}{\sum_j \sum_i C_{ij}}$$
- Micro-average: クラス間の平均。各クラスの値にクラスの要素数を重みとして平均。各要素 (document) の値の平均。
- Macro-average: クラス間の平均。各クラスの値を重みを等しいとして平均。普通は用いない。



Dumais et al. 1998: Reuters – Break-Even Performance

| | Rocchio | NBayes | Trees | LinearSVM |
|-------------|---------|--------|-------|-----------|
| earn | 92.9% | 95.9% | 97.8% | 98.2% |
| acq | 64.7% | 87.8% | 89.7% | 92.8% |
| money-fx | 46.7% | 56.6% | 66.2% | 74.0% |
| grain | 67.5% | 78.8% | 85.0% | 92.4% |
| crude | 70.1% | 79.5% | 85.0% | 88.3% |
| trade | 65.1% | 63.9% | 72.5% | 73.5% |
| interest | 63.4% | 64.9% | 67.1% | 76.3% |
| ship | 49.2% | 85.4% | 74.2% | 78.0% |
| wheat | 68.9% | 69.7% | 92.5% | 89.7% |
| corn | 48.2% | 65.3% | 91.8% | 91.1% |
| Avg Top 10 | 64.6% | 81.5% | 88.4% | 91.4% |
| Avg All Cat | 61.7% | 75.2% | na | 86.4% |

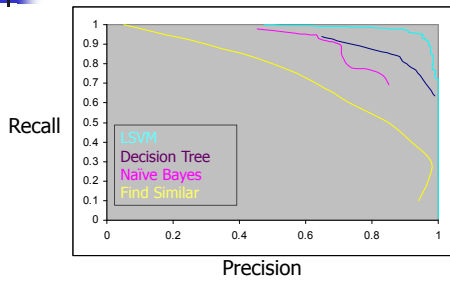
Break Even: Recall = Precision

Recall: = $TP/(TP+TN)$; % 当該カテゴリ中そのカテゴリに属すると判定したもの

Precision: = $TP/(TP+FP)$; % そのカテゴリに属するとして本当にそのカテゴリに属するもの

S. T. Dumais, J. Platt, D. Heckerman, and M. Sahami. Inductive learning algorithms and representations for text categorization. In CIKM-98. Proceedings of the Seventh International Conference on Information and Knowledge Management, 1998.

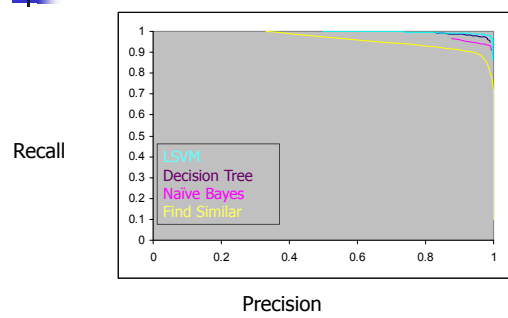
Precision vs. Recall - Category "Grain"



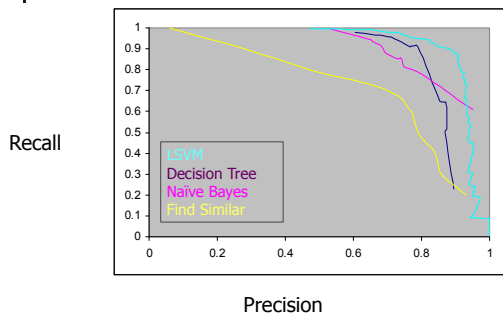
Recall: = $TP/(TP+TN)$; % 当該カテゴリ中そのカテゴリに属すると判定したもの

Precision: = $TP/(TP+FP)$; % そのカテゴリに属するとして本当にそのカテゴリに属するもの

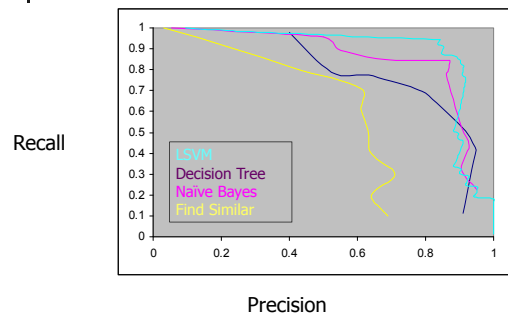
Precision vs. Recall - Category "Earn"



Precision vs. Recall - Category "Crude"



Precision vs. Recall - Category "Ship"



カーネルによる違い (Joachims)

| | Bayes | Rocchio | C4.5 | k-NN | SVM (poly) degree $d =$ | | | | | SVM (rbf) width $\gamma =$ | | | |
|-----------|-------------|-------------|-------------|-------------|----------------------------|------|-------------|-------------|-------------|-------------------------------|------|-------------|-------------|
| | | | | | 1 | 2 | 3 | 4 | 5 | 0.6 | 0.8 | 1.0 | 1.2 |
| earn | 95.9 | 96.1 | 96.1 | 97.3 | 98.2 | 98.4 | 98.5 | 98.4 | 98.3 | 98.5 | 98.5 | 98.4 | 98.3 |
| acq | 91.5 | 92.1 | 85.3 | 92.0 | 92.6 | 94.6 | 95.2 | 95.2 | 95.3 | 95.0 | 95.3 | 95.3 | 95.4 |
| money-fx | 62.9 | 67.6 | 69.4 | 78.2 | 66.9 | 72.5 | 75.4 | 74.9 | 76.2 | 74.0 | 75.4 | 76.3 | 75.9 |
| grain | 72.5 | 79.5 | 89.1 | 82.2 | 91.3 | 93.1 | 92.4 | 91.3 | 89.9 | 93.1 | 91.9 | 91.9 | 90.6 |
| crude | 81.0 | 81.5 | 75.5 | 85.7 | 86.0 | 87.3 | 88.6 | 88.9 | 87.8 | 88.9 | 89.0 | 88.9 | 88.2 |
| trade | 50.0 | 77.4 | 59.2 | 77.4 | 69.2 | 75.5 | 76.6 | 77.3 | 77.1 | 76.9 | 78.0 | 77.8 | 76.8 |
| interest | 58.0 | 72.5 | 49.1 | 74.0 | 69.8 | 63.3 | 67.9 | 73.1 | 76.2 | 74.4 | 75.0 | 76.2 | 76.1 |
| ship | 78.7 | 83.1 | 80.9 | 79.2 | 82.0 | 85.4 | 86.0 | 86.5 | 86.0 | 85.4 | 86.5 | 87.6 | 87.1 |
| wheat | 60.6 | 79.4 | 85.5 | 76.6 | 83.1 | 84.5 | 85.2 | 85.9 | 83.8 | 85.2 | 85.9 | 85.9 | 85.9 |
| corn | 47.3 | 62.2 | 87.7 | 77.9 | 86.0 | 86.5 | 85.3 | 85.7 | 83.9 | 85.1 | 85.7 | 85.7 | 84.5 |
| microavg. | 72.0 | 79.9 | 79.4 | 82.3 | 84.2 | 85.1 | 85.9 | 86.2 | 85.9 | 86.4 | 86.5 | 86.3 | 86.2 |
| | | | | | combined: 86.0 | | | | | combined: 86.4 | | | |

Fig. 2. Precision/recall-breakeven point on the ten most frequent Reuters categories and microaveraged performance over all Reuters categories, k-NN, Rocchio, and C4.5 achieve highest performance at 1000 features (with $k = 30$ for k-NN and $\beta = 1.0$ for Rocchio). Naive Bayes performs best using all features.

T. Joachims, *Text Categorization with Support Vector Machines: Learning with Many Relevant Features*. Proceedings of the European Conference on Machine Learning (ECML), Springer, 1998

Yang&Liu: SVM vs 他手法

Table 1: Performance summary of classifiers

| method | miR | miP | miF1 | maF1 | error |
|--------|-------|-------|-------|-------|--------|
| SVM | .8120 | .9137 | .8599 | .5251 | .00365 |
| KNN | .8339 | .8807 | .8567 | .5242 | .00385 |
| LSF | .8507 | .8489 | .8498 | .5008 | .00414 |
| NNet | .7842 | .8785 | .8287 | .3765 | .00447 |
| NB | .7688 | .8245 | .7956 | .3886 | .00544 |

miR = micro-avg recall; miP = micro-avg prec.;
miF1 = micro-avg F1; maF1 = macro-avg F1.

LLSF: Linear Least Square Fit

Nnet: 中間素子数 64

SVM kernel: linear (他より良かった)

特徴数

NNet: 1000

NB: 2000

KNN: 2415

LLSF: 2415

SVM: 10000

まとめ

- サポートベクターマシン (SVM) は
 - サポートベクターに基づいて超平面を決める
 - Support vector = 判定境界付近のクリティカルな点
 - 線形 SVM は線形分類器
 - カーネル: 高次元へ写像するが、その内積は低次元の内積で簡単に計算できる
 - リスクの上界 (リスク = テストデータでの期待誤り)
 - (邪魔な属性が多いときの) 分類器としてベスト?
 - 数1000も属性があるときは、安定的に強い
 - ポピュラー: SVMlight がきっかけ?
 - 速くて無料 (研究目的には)
 - 他にもいくつか: TinySVM, libsvm, ...

参考

- A Tutorial on Support Vector Machines for Pattern Recognition (1998) Christopher J. C. Burges
- S. T. Dumais, Using SVMs for text categorization, *IEEE Intelligent Systems*, 13(4):21-23, Jul/Aug 1998
- S. T. Dumais, J. Platt, D. Heckerman and M. Sahami. 1998. Inductive learning algorithms and representations for text categorization. *Proceedings of CIKM '98*, pp. 148-155.
- A re-examination of text categorization methods (1999) Yiming Yang, Xin Liu 22nd Annual International SIGIR
- Tong Zhang, Frank J. Oles: Text Categorization Based on Regularized Linear Classification Methods. *Information Retrieval* 4(1): 5-31 (2001)
- Trevor Hastie, Robert Tibshirani and Jerome Friedman, "Elements of Statistical Learning: Data Mining, Inference and Prediction" Springer-Verlag, New York.
- 'Classic' Reuters data set: <http://www.daviddlewis.com/resources/testcollections/reuters21578/>
- T. Joachims, *Learning to Classify Text using Support Vector Machines*. Kluwer, 2002.

SVM回帰 (Support Vector Regression)

SVM は

線形モデルを用いて分類を行う手法 (の一つ):

$$y(x) = w^T \phi(x) + b$$

そのSVMを用いて回帰を行う

ある見方: コスト関数 = 誤差関数 + 正則化項

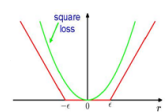
線形回帰で、次の誤差関数の最小化を図る:

$$\frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 + \frac{\lambda}{2} \|w\|^2$$

この二乗誤差部分を ϵ -insensitive 誤差関数で置き換える:

$$C \sum_{n=1}^N E_{\epsilon}(y(x_n) - t_n) + \frac{1}{2} \|w\|^2$$

ϵ -insensitive 誤差関数の例:



$$L_{\epsilon}(y) = \begin{cases} 0 & \text{for } |f(x) - y| < \epsilon \\ |f(x) - y| - \epsilon & \text{otherwise} \end{cases}$$

Gunn, *Support Vector Machines for Classification and Regression*

スラック変数の導入

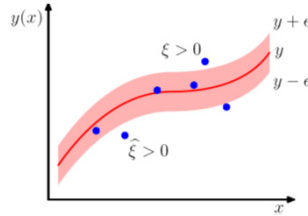
目標値がこの ε -tube の中にあるためには:

$$y_n - \varepsilon \leq t_n \leq y_n + \varepsilon$$

スラック変数を導入して、 ε -tube の外側に点があることを許すには:

$$t_n \leq y(x_n) + \varepsilon + \xi_n$$

$$t_n \geq y(x_n) - \varepsilon - \xi_n^-$$



Support Vector Regression の最適化問題

Minimize:

$$C \sum_{n=1}^N (\xi_n + \xi_n^-) + \frac{1}{2} \|w\|^2$$

Subject to:

$$\begin{aligned} \xi_n &\geq 0 & t_n &\leq y(x_n) + \varepsilon + \xi_n \\ \xi_n^- &\geq 0 & t_n &\geq y(x_n) - \varepsilon - \xi_n^- \end{aligned} \quad \text{かつ}$$

ラグランジュ関数

最小化問題:

$$L = C \sum_{n=1}^N (\xi_n + \xi_n^-) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \mu_n^- \xi_n^-) - \sum_{n=1}^N a_n (\varepsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N a_n^- (\varepsilon + \xi_n^- - y_n + t_n)$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{n=1}^N (a_n - a_n^-) \phi(x_n)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N (a_n - a_n^-) = 0$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n + \mu_n = C$$

$$\frac{\partial L}{\partial \xi_n^-} = 0 \Rightarrow a_n^- + \mu_n^- = C$$

双対方程式

最大化:

$$W(a, a^-) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - a_n^-)(a_m - a_m^-) k(x_n, x_m) - \varepsilon \sum_{n=1}^N (a_n + a_n^-) + \sum_{n=1}^N (a_n - a_n^-) t_n$$

$$0 \leq a_n \leq C$$

$$0 \leq a_n^- \leq C$$

予測:

$$y(x) = \sum_{n=1}^N (a_n - a_n^-) k(x, x_n) + b$$

定数項bの決め方

Karush-Kuhn-Tucker (KKT) 条件: サポートベクターに関する

$$a_n (\varepsilon + \xi_n + y_n - t_n) = 0$$

$$a_n^- (\varepsilon + \xi_n^- - y_n + t_n) = 0$$

$$(C - a_n) \xi_n = 0$$

$$(C - a_n^-) \xi_n^- = 0$$

$$b = t_n - \varepsilon - w^T \phi(x_n) = t_n - \varepsilon - \sum_{m=1}^N (a_m - a_m^-) k(x_n, x_m)$$

カーネル法

カーネルにはいろいろ

- 「カーネル」という用語は、一般には(対称や正定値とは限らない)2変数関数に使われることも多い。
例)ノンパラメトリックな確率密度推定

$$p(x) = \frac{1}{N} \sum_{i=1}^N g(x - x_i)$$

などに用いる密度関数 $g(x)$ も「カーネル」と呼ばれる(関係はあるのですが)。

例)積分変換の核関数

- 最近では、正定値カーネルのことを単に「カーネル」と呼ぶことが多いので(Mercer's theoremのせいかな?)、注意が必要。

カーネル法の動機

- 線形関数の学習には良い性質が多い
 - 最適解は一つ
 - 高速な学習アルゴリズムが存在する
 - 統計的な解析がうまくできる
- しかし、大きな問題がある
 - 学習能力が不十分(仮説空間の複雑性が不足)

歴史を紐解くと

- 例の Minsky & Pappert "Perceptrons" でその弱点を明らかにした
- ニューラルネットワークは、この弱点を、線形関数に対し非線形の活性化関数を施すことによって、克服した
 - 学習能力の低さを解決し、学習アルゴリズムが広く利用できることを示した
 - しかし学習速度の遅さ、極小点の多さという問題に遭遇した

そこで、カーネル法

- カーネル法は、線形関数に拘るのだが、それを高次元空間で行おうというもの:

$$\phi: x \in X \rightarrow \phi(x) \in F$$

- 期待しているのは、特徴空間が、入力空間よりはるかに高い次元を持てば、表現能力(学習能力)が増大するであろうということ。

例

- 次の写像を考える

$$(x_1, x_2) \rightarrow (x_1^2, x_1x_2, x_2x_1, x_2^2)$$

- この特徴空間で線形制約を考える。例えば:

$$ax_1^2 + bx_2^2 = c$$

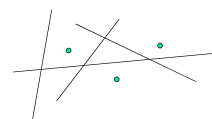
- (この場合)楕円になる - i.e. すなわち、元の入力空間でみると、非線形な形となる。

特徴空間の表現能力

- 表現能力(従って、学習能力)はその次元に比例する - 例:

定理: m 次元空間の一般位置に $m+1$ 個の点を与えられると、それらのどのような2値分類も線形閾値関数で表現できる

- 2次元:



関数の形

- というわけで、カーネル法では、特徴空間では線形関数を用いる:

$$x \rightarrow \langle \mathbf{w}, \phi(x) \rangle$$

- 回帰課題のときには、これが回帰関数となる。
- 分類課題の場合には、この関数値に閾値関数を施す必要がある(単に、0 or 1 にするためです)

$$x \rightarrow \text{sgn}(\langle \mathbf{w}, \phi(x) \rangle + b)$$

高次元空間での問題

- 学習能力を非常に高くすることは、実は容易である。

$$\phi: x \in X \rightarrow \phi(x) \in F$$

- そして過学習が発生する:
例えば、N個の点を N-1次元空間に写像
⇒どんな分類も実現できる
⇒どんな分類も表現できるということは、汎化能力がないということである
- 長いベクトルを扱うわけであるから、計算量も馬鹿にならない

対処法

- 過学習への対応: マージン最大化
- 計算量への対応: カーネルトリック

ところで、SVMでは

- 最適化問題は

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{subject to: } y_i (\langle \mathbf{w}, \phi(x_i) \rangle + b) \geq 1 - \xi_i \\ \xi_i \geq 0$$

計算量を考えてみると

- 2次カーネルでは

$$(x_1, x_2) \rightarrow (x_1^2, x_1x_2, x_2x_1, x_2^2)$$

例えば 30x30 = 900 ピクセル(グレイレベル)の画像に対しては、約405000次元になってしまう

- 重複を省いて数えた ${}_{900}C_2$
- この特徴空間で計算するのは、あまりに馬鹿らしい

双対表現を考える

- 荷重ベクトルの空間を制限したらどうだろうか。
例えば、次のように、訓練データの線形結合で表現できるもののみを考える(パーセプトロアルゴリズムを思い出す):

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \phi(x_i)$$

実は一定の条件のもと、SVMでは、最適値に対して、これは常に成立する (Representer Theorem)

- このとき、未知例に対する内積は、

$$\langle \mathbf{w}, \phi(x) \rangle = \left\langle \sum_{i=1}^m \alpha_i \phi(x_i), \phi(x) \right\rangle = \sum_{i=1}^m \alpha_i \langle \phi(x_i), \phi(x) \rangle$$

双対変数を学習する

- この α_i は双対変数と呼ばれる
- 荷重ベクトル w が訓練データの(変換の)張る空間(未知データも含むしように)直交する成分を持っていたとしても、その成分は、予測には何の影響も与えない(前のスライド)ゆえ、「ある解(w)があれば、それと予測能力が同じで、訓練データの張る空間に属する解がある」
 - representer theorem という
 - もし未知データが、訓練データの張る空間に属さないときはどうなるか？実はこの問いは意味がない。訓練データの張る空間に属さない成分については、訓練データに情報が無いので、答えようがないからである。
- 荷重ベクトルを直接学習する代わりに、この双対変数の学習を行えばよいことになる

SVM の双対問題

- 先ほどの双対問題は、SVM最適化問題の双対問題としても得られる:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle \\ \text{subject to:} \quad & \sum_{i=1}^m y_i \alpha_i = 0 \\ & 0 \leq \alpha_i \leq C \quad \text{for all } i \end{aligned}$$

ここで、カーネルを用いると

- 内積しか使っていないことに注意
- 仮に、次のようにして計算が簡単にできたら:
計算が: 簡単 大変
$$\kappa(x, z) = \langle \phi(x), \phi(z) \rangle$$
 Kの計算が簡単だという仮定のもと
- そうすると、訓練時にも未知データに対しても、(極めて高次元の)特徴空間での計算を実質的には行う必要がなくなる

その例

- 例によって次の例を考えると

$$\phi: (x_1, x_2) \rightarrow (x_1^2, x_1 x_2, x_2 x_1, x_2^2)$$

- こんな具合:

$$\begin{aligned} \langle \phi(x), \phi(z) \rangle &= \langle \mathbf{xx}^t, \mathbf{zz}^t \rangle && \text{行列の同一位置の要素の積和} \\ &= \text{tr}(\mathbf{xx}^t \mathbf{zz}^t) \\ &= \langle \mathbf{x}, \mathbf{z} \rangle^2 \end{aligned}$$

計算の効率性

- つまり、例えば、先ほどの 約405000次元のベクトルの場合、実は 900 次元のベクトルで内積をとり、それを二乗するだけでよいことになる
- 一気に敷衍すれば、無限次元の特徴空間も扱うことができることになる。例えばガウスカーネル:

$$\kappa(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

無限次元の特徴関数を用いてその内積 $\forall x, y \in X, \lim_{n \rightarrow \infty} \kappa_n(x, y) = \kappa(x, y)$ で表現可能なことを示すのは少々面倒。

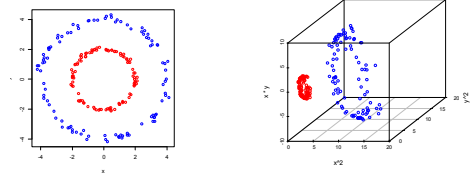
非線形化法としてのカーネル法

非線形データ解析方法としてのカーネル法

- 「線形判別」における「判別」を非線形にするという視点からカーネル法を見てきたが、「線形」関数ではなく、訓練データ空間での非線形関数の推測する方法として、カーネル法を考えることができる。
- (何らかの)線形手法において、線形関数を非線形関数に置き換える方法としてのカーネル法を考えよう。すなわち

$$\mathbf{w} \cdot \mathbf{X} \Rightarrow f(\mathbf{w}, \mathbf{X})$$

例えば(といってもSVMの例)



$$(x_1, x_2) \rightarrow (x_1^2, x_2^2, x_1x_2)$$

活躍するのは

- 再生性
 - 関数空間 H に属する (非線形) 関数の値が内積で計算できる
- カーネルトリック
 - 高次元 (無限次元) 関数空間 H において、内積が容易に計算できる。
 - さまざまな線形データ解析手法が H 上で適用可能
SVM, Kernel 主成分分析, Kernel 正準相関分析
 - $\langle \Phi(x_i), \Phi(x_j) \rangle_H = k(x_i, x_j)$... 正定値カーネル
データ x_i と x_j の類似度
(まったく似ていない時 0)

$$f(x) = \langle f, \Phi(x) \rangle_H \quad \begin{array}{l} f \text{ も } \Phi(x) \text{ も } H \text{ の要素。} \\ f(\cdot), \Phi(x)(\cdot) \text{ と書いた} \\ \text{方が良からう} \end{array}$$

カーネルによる非線形化

- 線形項 $\mathbf{w} \cdot \mathbf{X}$ の代わりに $f(\mathbf{X})$
- 再生性
 $f(\mathbf{X}_i) = \langle f, \Phi(\mathbf{X}_i) \rangle_{H_k}$
- 多くの場合 $f = \sum_{i=1}^N \alpha_i \Phi(\mathbf{X}_i) = \sum_{i=1}^N \alpha_i k(\cdot, \mathbf{X}_i)$
- $\langle f, \Phi(\mathbf{X}_i) \rangle_{H_k} = \sum_{j=1}^N \alpha_j \langle \Phi(\mathbf{X}_j), \Phi(\mathbf{X}_i) \rangle = \sum_{j=1}^N \alpha_j k(\mathbf{X}_j, \mathbf{X}_i)$
- Gram行列を用いて計算する

Representer Theorem

Theorem H を RKHS とし、 $k: X \times X \rightarrow \mathbb{R}$ をそのカーネル (コンパクト領域上の対称な半正定値関数) とする。任意の関数 $L: \mathbb{R}^n \rightarrow \mathbb{R}$ と任意の非減少関数 $\Omega: \mathbb{R} \rightarrow \mathbb{R}$ を考える。もし

$$J^* := \min_{f \in H} J(f) := \min_{f \in H} \left\{ \Omega(\|f\|_H^2) + L(f(x_1), \dots, f(x_n)) \right\} \quad (1)$$

が well-defined であれば、ある $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ が存在して

$$f(\cdot) = \sum_{i=1}^n \alpha_i k(x_i, \cdot) \quad (2)$$

が $J(f) = J^*$ を達成する。さらに、もし Ω が増加関数であれば、 $J(f)$ を最小にするものはどれも (2) 式で表現される。

RKHS: Reproducing kernel Hilbert space

PCA

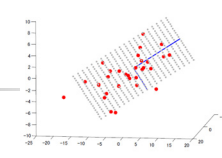
- 主成分分析 (PCA)
 - m 次元データ X_1, \dots, X_N
 - 分散最大化する方向 (単位ベクトル \mathbf{a}) の発見

$$\text{Var}[\mathbf{a}^T X] = \frac{1}{N} \sum_{i=1}^N \left(\mathbf{a}^T (X_i - \frac{1}{N} \sum_j X_j) \right)^2 = \mathbf{a}^T V \mathbf{a}$$

$$\text{分散共分散行列 } V = \frac{1}{N} \sum_{i=1}^N \tilde{X}_i \tilde{X}_i^T$$

V の単位固有ベクトル u_1, \dots, u_m ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$)

第 p 主成分の軸 = u_p
データ X_j の第 p 主成分 = $u_p^T X_j$



PCAの非線形化

$$\max_{\|a\|=1} \text{Var}[a^T X] = \frac{1}{N} \sum_{i=1}^N \left(a^T X_i - \frac{1}{N} \sum_j a^T X_j \right)^2$$

$$\max_{\|f\|=1} \text{Var}[f(X)] = \frac{1}{N} \sum_{i=1}^N \left(f(X_i) - \frac{1}{N} \sum_j f(X_j) \right)^2$$

Fを探す空間として
RKHS H_k をとると

$$\max_{f \in H_k, \|f\|=1} \text{Var}[f, \Phi(X)] = \frac{1}{N} \sum_{i=1}^N \left(\langle f, \Phi(X_i) \rangle - \frac{1}{N} \sum_j \langle f, \Phi(X_j) \rangle \right)^2$$

$f, \Phi(X_i)$ は H_k のベクトル $\Rightarrow H_k$ における線形な問題

特徴ベクトル $\Phi(X_1), \dots, \Phi(X_N)$ に対する H_k 内のPCA

カーネルPCA

カーネル k を設定

特徴ベクトル $\Phi(X_1), \dots, \Phi(X_N)$ に対する H_k 内のPCA

$$\max_{f \in H_k, \|f\|=1} \text{Var}[f, \Phi(X)] = \frac{1}{N} \sum_{i=1}^N \left(\langle f, \tilde{\phi}_i \rangle \right)^2 \quad \tilde{\phi}_i = \Phi(X_i) - \frac{1}{N} \sum_j \Phi(X_j)$$

$$f = \sum_{i=1}^N \alpha_i \tilde{\phi}_i \quad \text{としてよい}$$

$$\text{そうすると、分散} = \frac{1}{N} \sum_{i=1}^N \left(\sum_{a=1}^N \alpha_a \langle \tilde{\phi}_i, \tilde{\phi}_a \rangle \right)^2 = \alpha^T \tilde{K}^2 \alpha \quad \tilde{K}_{i,j} = \langle \tilde{\phi}_i, \tilde{\phi}_j \rangle$$

主成分は

$$\max_{\alpha} \alpha^T \tilde{K}^2 \alpha$$

$$\text{subject to} \quad \alpha^T \tilde{K} \alpha = 1 \quad \leftarrow \quad \|f\|_{\tilde{H}_k}^2 = \left\langle \sum_i \alpha_i \tilde{\phi}_i, \sum_j \alpha_j \tilde{\phi}_j \right\rangle = \alpha^T \tilde{K} \alpha$$

ちょっとしたまとめ

- カーネルによる非線形化
 - カーネル化
 - 線形関数 $w^T X \Rightarrow$ 非線形関数 $f(X) \Rightarrow$ 再生性 $f(X_i) = \langle f, \Phi(X_i) \rangle_{H_k}$
 - RKHS内の特徴ベクトルに対する線形手法とも考えられる
 - データとパラメータの内積を使って表される線形手法(射影、相関、分散共分散 etc.)なら同様に適用可能
 - 例: SVM, カーネルPCA, カーネルCCA (正準相関分析)、スプライン平滑化、カーネルFisher判別分析等
- 非線形化することの特徴
 - 線形では捕らえることのできない特性の把握
 - ただし、捕らえることのできる非線形性はカーネルに依存