

*Pattern
Recognition
and
Machine
Learning*



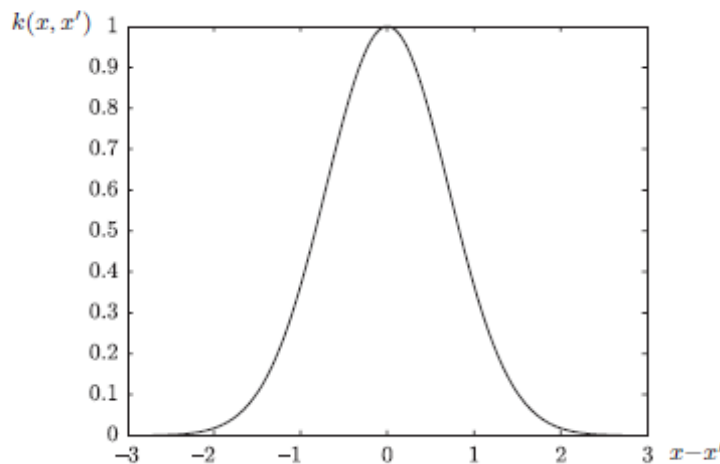
第6章 カーネル法

修士2年
藤井 敬士

カーネル法とは

■ カーネル関数を用いたデータ解析手法

- カーネル関数とは、二つの入力 $\mathbf{x}=(x_1, \dots, x_d)$, $\mathbf{x}'=(x'_1, \dots, x'_d)$ から計算される関数 $k(\mathbf{x}, \mathbf{x}')$.
- 直観的には、 $k(\mathbf{x}, \mathbf{x}')$ は \mathbf{x} と \mathbf{x}' の近さのようなものである
- 不変カーネル $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ と 均一カーネル (RBF) $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$



よく使われるカーネル関数の例

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\beta \|\mathbf{x} - \mathbf{x}'\|^2)$$

$\|\mathbf{z}\|$ の意味は $\mathbf{z}=(z_1, \dots, z_d)$

$$\|\mathbf{z}\|^2 = \sum_{m=1}^d z_m^2$$

また、 β は適当に決めるパラメータである

カーネル法とは

- x に対して関数 $y = \sum_{j=1}^n \alpha_j k(\mathbf{x}^{(j)}, \mathbf{x})$ を当てはめ、二乗誤差
 $r_k(y, \mathbf{x}; \boldsymbol{\alpha}) = (y - \sum_{j=1}^n \alpha_j k(\mathbf{x}^{(j)}, \mathbf{x}))^2$ を最小化する α を求める.

- 二乗誤差の総和は $R_k(\boldsymbol{\alpha}) = \sum_{i=1}^n r_k(y^{(i)}, \mathbf{x}^{(i)}; \boldsymbol{\alpha}) = (\mathbf{y} - \mathbf{K}\boldsymbol{\alpha})^T (\mathbf{y} - \mathbf{K}\boldsymbol{\alpha})$

ここで,

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & k(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(n)}, \mathbf{x}^{(1)}) \\ k(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & k(\mathbf{x}^{(2)}, \mathbf{x}^{(2)}) & \dots & k(\mathbf{x}^{(n)}, \mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(1)}, \mathbf{x}^{(n)}) & k(\mathbf{x}^{(2)}, \mathbf{x}^{(n)}) & \dots & k(\mathbf{x}^{(n)}, \mathbf{x}^{(n)}) \end{pmatrix}$$

この解は、 \mathbf{K} が正則ならば $\boldsymbol{\alpha} = (\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T \mathbf{y}$

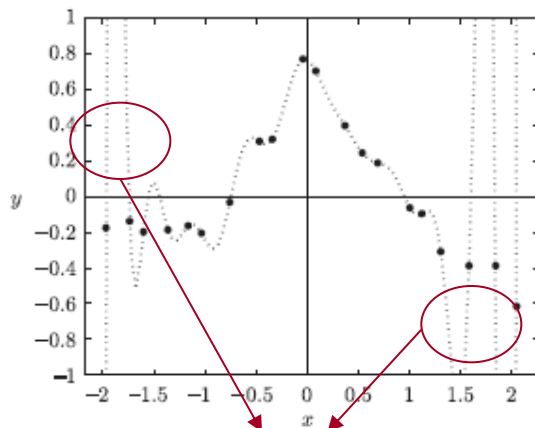
任意の \mathbf{x}, \mathbf{x}' について $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$ が成り立つので \mathbf{K} は対称
行列

$$\boldsymbol{\alpha} = \mathbf{K}^{-1} \mathbf{y}$$

すなわち、 $\mathbf{K}^T = \mathbf{K}$ だから、 $(\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T = (\mathbf{K}^2)^{-1} \mathbf{K} = \mathbf{K}^{-1}$ よって、

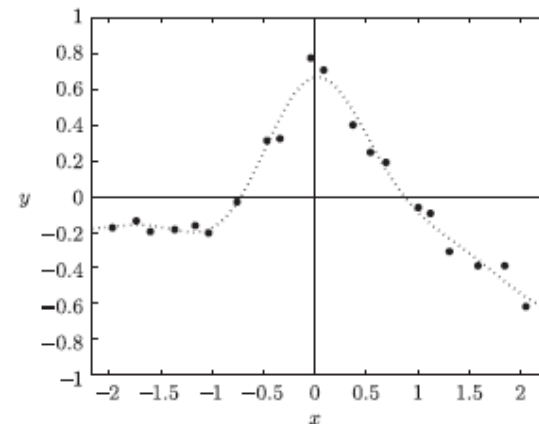
カーネル法とは

■ 推定結果



過学習による誤推定

正則化



$\lambda=0.01$ の場合の関数

$R_{k,\lambda}(\alpha) = (\mathbf{y} - \mathbf{K}\alpha)^T(\mathbf{y} - \mathbf{K}\alpha) + \lambda \alpha^T \mathbf{K} \alpha, \quad \lambda > 0$ を最小化する

正則化項

$$-K(\mathbf{y} - \mathbf{K}\alpha) + \lambda \mathbf{K} \alpha = 0 \quad \text{より,} \quad \alpha = (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

※正則化パラメータ λ の取り方には任意性が残る

λ が小さいと不安定な解, λ が大きいと $\alpha=0$ に近づく

カーネル法を用いることの利点

1. サンプルの増加と共にどんどん複雑に出来る
(正則化パラメータを適当にとると)複雑な関数を表現することが出来る.
2. 線形性と非線形性を両方持つ
 $y = \sum_{j=1}^n \alpha_j k(\mathbf{x}^{(j)}, \mathbf{x})$ より, 決めるべきパラメータ α については線形性を持つが, 入力データについては非線形な関数を表現できる.
3. 高次元・非数値データへの適用
カーネル関数の中身は1次元の実数に限らず, 高次元, 文字列, グラフ構造などについても同様に扱うことが出来る.
4. カーネル関数のモジュール化
最適解 α は行列 \mathbf{K} のみに依存し, カーネル関数がどんなものであるかは関係ない. つまり, カーネル関数を計算する部分とそれ以降の処理を分割できる.

双対表現

- パラメータベクトル \mathbf{w} を直接扱う代わりに、最小二乗法のアルゴリズムをパラメータベクトル \mathbf{a} で表現しなおすこと.
- この表現によって、カーネル関数が見える形になる.

□ $J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}.$ のような正則化された二乗
和誤差の最小

化を考える($\lambda \geq 0$ とする).

□ $J(\mathbf{w})$ $\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}.$

すなわち、係数が \mathbf{w} の関数であるような $\phi(\mathbf{x}_n)$ の線形結合となる. ここで
 $a_n = -\frac{1}{\lambda} \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}$ $\phi(\mathbf{x}_n)^T$ で与えられるような計画行列. また,
として $\mathbf{a}=(a_1, \dots, a_N)$ とする.

双対表現

- $\mathbf{w} = \Phi^T \mathbf{a}$ を $J(\mathbf{w})$ に代入すると, $J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$
ここで, $\mathbf{t} = (t_1, \dots, t_n)$ とする.

- 次に, $N \times N$ の対称行列で, その要素が $K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$
で表されるグラム行列 $\mathbf{K} = \Phi \Phi^T$ を定義す. $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$.

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a} \\ \mathbf{a} &= (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}. \end{aligned}$$

\mathbf{w} を消去して \mathbf{a} について解くと,

- これを線形回帰モデルに代入し直すことによって, 新たな \mathbf{x} に対する予測 $y(\mathbf{x})$ は以下のようになります
$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \underline{k(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}}.$$

カーネル関数のみで表現可能

双対: $\phi(\mathbf{x})$ の要素の線形結合によって \mathbf{a} が表現できることから, パラメータベクトル \mathbf{w} を用いたもとの定式化を復元できる.

特徴ベクトル $\phi(\mathbf{x})$ を明示的に考えなくても, カーネル関数で表現できる

双対表現

- 双対表現 においては, $N \times N$ 行列の逆行列を求めることでパラメータ \mathbf{a} が得られ $\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$.
- もともとの表現においては, $M \times M$ 行列の逆行列を求めればよかった.
- 通常は $N \gg M$. しかし, 双対表現を用いることで特徴ベクトル $\phi(\mathbf{x})$ を明示的に考えずに, 高次元や無限次元の特徴空間を間接的に扱うことが出来る.
- 回帰のための確率的な線形モデルとガウス過程の双対性
- サポートベクトルマシンとの関連性(7章)

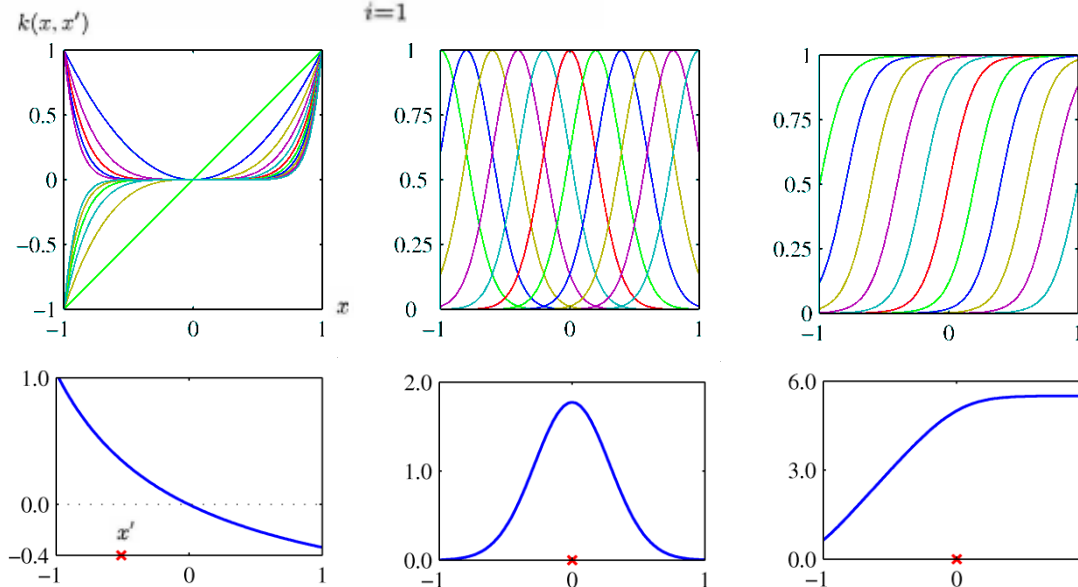
カーネル関数の構成

- カーネル置換を行う＝有効なカーネル関数を構成する必要
 1. 特徴空間への写像 $\phi(x)$ を考え、対応するカーネルを構成する.
 2. カーネル関数を直接定義する.
- 1. カーネル関数は以下のように定義されている.

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x').$$

$\phi_i(x)$ は基底関

数



カーネル関数の構成

2. 次の例で考える. $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$

2次元の入力空間 $\mathbf{x} = (x_1, x_2)$ を考えて, 上式を展開

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\ &= \phi(\mathbf{x})^T \phi(\mathbf{z}). \end{aligned}$$

特徴空間への写像は $\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T$ の形を持ち,
すべての2次の項を含む.

関数 $k(\mathbf{x}, \mathbf{x}')$ が有効なカーネル

\Leftrightarrow 任意の $\{\mathbf{x}_n\}$ に対して, 要素が $k(\mathbf{x}_n, \mathbf{x}_m)$ で与えられるグラム行列 \mathbf{K} が半正定値であること.

新たなカーネル関数の構成法

- $k_1(\mathbf{x}, \mathbf{x}')$ と $k_2(\mathbf{x}, \mathbf{x}')$ が有効なカーネルであるとき, 下の関数もカーネル関数として有効である.

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}')$$

• $c > 0$ は定数

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$$

• $f(\cdot)$ は任意の関数

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$$

• $q(\cdot)$ は非負の係数を持つ多項式

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$$

• $\phi(\mathbf{x})$ は \mathbf{x} から \mathbf{R}^M への関数

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

• $k_3(\cdot, \cdot)$ は \mathbf{R}^M で定義された有効なカーネル

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$$

• A は対称な半正定値行列

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}'))$$

• \mathbf{x}_a と \mathbf{x}_b は $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ であるような変数

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T A \mathbf{x}'$$

• k_a と k_b はそれぞれの特徴空間において有効なカーネル関数

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b).$$

カーネル関数の例

- ガウスカーネル

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$$

- 生成モデルに基づくカーネル

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}').$$

- 配列 \mathbf{X} と \mathbf{X}' の類似度を測るカーネル

$$k(\mathbf{X}, \mathbf{X}') = \sum_{\mathbf{Z}} p(\mathbf{X}|\mathbf{Z})p(\mathbf{X}'|\mathbf{Z})p(\mathbf{Z}).$$

- フィッシャーカーネル

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\theta, \mathbf{x})^T \mathbf{F}^{-1} \mathbf{g}(\theta, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\theta, \mathbf{x})^T \mathbf{g}(\theta, \mathbf{x}').$$

- シグモイドカーネル

$$k(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^T \mathbf{x}' + b)$$

RBFネットワーク

- 線形基底関数モデル(3章)では, 基底関数の形を考えていなかった.
→RBF(動径基底関数:radial basis function)が良く使われる.

$$\phi_j(\mathbf{x}) = h(\|\mathbf{x} - \mu_j\|)$$

- もともとは, 関数補間 (目的変数の値を正確に表現できる関数を求めること) のために導入された.

$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|).$$

- 入力変数にノイズが含まれる場合の補間にも使われる

入力変数 \mathbf{x} に含まれるノイズが, 確率分布 $\nu(\xi)$ に従う確率変数 ξ によって表されるとき

$$E = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n + \xi) - t_n\}^2 \nu(\xi) d\xi.$$

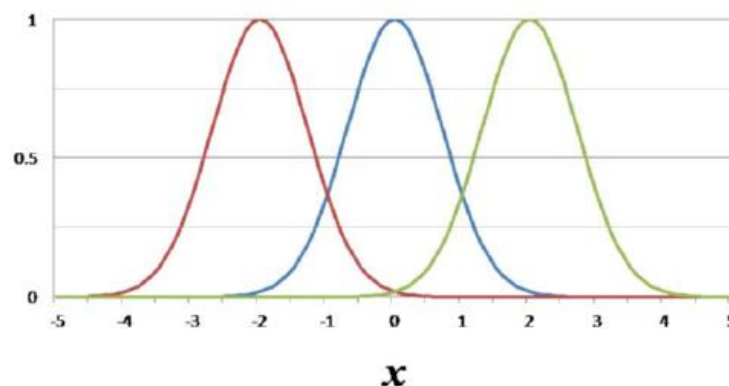
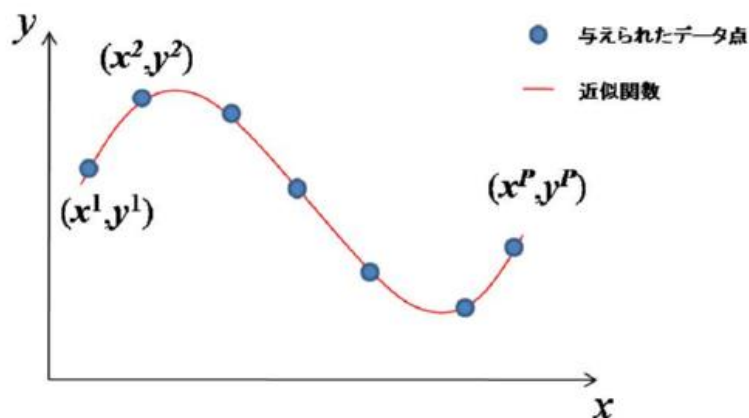
変分法を用いて以下のように最適化でき, RBFも求められる

$$y(\mathbf{x}) = \sum_{n=1}^N t_n h(\mathbf{x} - \mathbf{x}_n), \quad h(\mathbf{x} - \mathbf{x}_n) = \frac{\nu(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N \nu(\mathbf{x} - \mathbf{x}_n)}.$$

Nadaraya-Watsonモデル

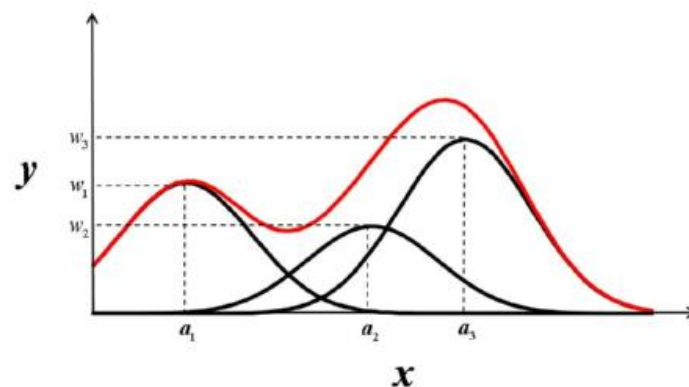
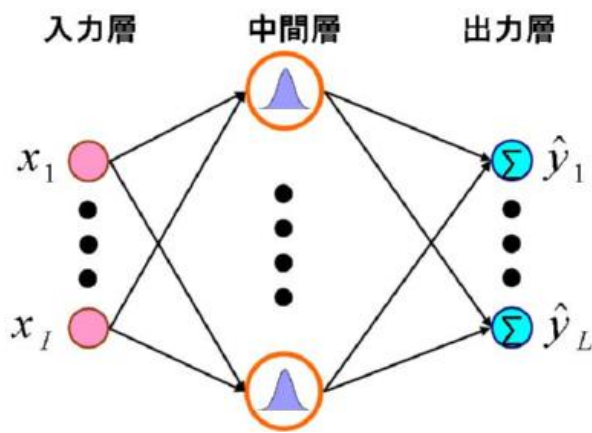
RBFネットワーク

- 3層から構成されるニューラルネットワーク
- 最小二乗法によって関数の最良近似法を導くことができる
=安定した学習が可能
- ガウス関数を基底関数として用いることが多い



RBFネットワーク

- ネットワーク構造の中間層にさまざまなRBFを使用
- 出力層は中間層出力の多重和
→複数のRBFに重み付けをして足し合わせることで任意の関数を実現
- RBFの中心値と重みを調整することで学習データと出力の誤差を小さくする.



RBFネットワーク

デモ



Nadaraya-Watsonモデル

- 訓練集合を $\{\mathbf{x}_n, t_n\}$ として, 同時分布 $p(\mathbf{x}, t)$ を推定するために, Parzen推定法を用いる. $f(\mathbf{x}, t)$ は密度関数の要素.

$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n).$$

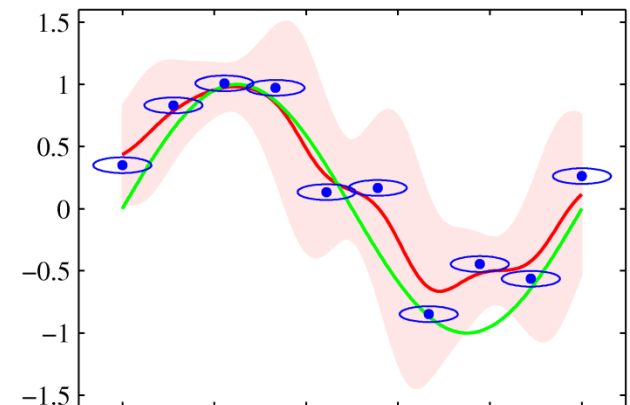
回帰関数 $y(\mathbf{x})$ は, 目標変数の条件付き期待値

$$\begin{aligned} y(\mathbf{x}) &= \mathbb{E}[t|\mathbf{x}] = \int_{-\infty}^{\infty} t p(t|\mathbf{x}) dt \\ &= \frac{\int t p(\mathbf{x}, t) dt}{\int p(\mathbf{x}, t) dt} \\ &= \frac{\sum_n \int t f(\mathbf{x} - \mathbf{x}_n, t - t_n) dt}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt}. \end{aligned}$$

$$\begin{aligned} y(\mathbf{x}) &= \frac{\sum_n g(\mathbf{x} - \mathbf{x}_n) t_n}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} \\ &= \sum_n k(\mathbf{x}, \mathbf{x}_n) t_n. \\ k(\mathbf{x}, \mathbf{x}_n) &= \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)}. \\ g(\mathbf{x}) &= \int_{-\infty}^{\infty} f(\mathbf{x}, t) dt \end{aligned}$$

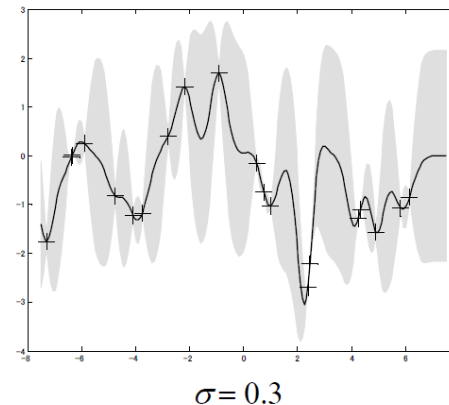
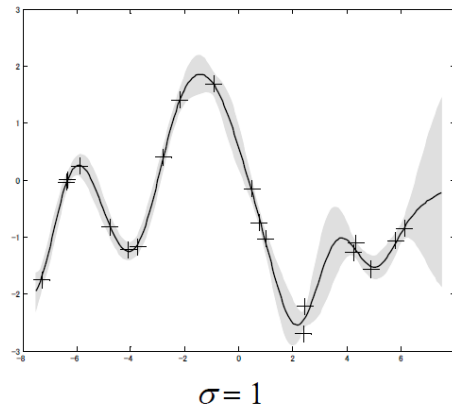
Nadaraya-Watsonモデル(カーネル回帰)

- データ点 \mathbf{x} に近いデータ点 \mathbf{x}_n ほど大きな重みを与えることができる



ガウス過程

- 関数 $y(\mathbf{x})$ の上の確率分布として定義され、任意の点集合 x_1, \dots, x_N に対する $y(\mathbf{x})$ の値の同時分布がガウス分布に従うもの
- 線形回帰モデルだけでは訓練データを増やせば増やすほど予測がほぼ期待値に一致してしまう。しかし、一般化されたガウス過程のモデルを導入すると、訓練データに近いところでは分散が小さく、離れるごとに分散が大きくなるモデルとなる。



$$\text{共分散関数 } R(s, t) = \exp\left(-\frac{1}{2\sigma^2}(s-t)^2\right)$$