

ベイズ最適化・ベイズ推定

強化学習の1つ。

1. 既知データを学習
2. 所望の特性を出す確率の高い記述子候補を推薦
3. 推薦に応じてデータを追加、1へ戻る

ベイズ最適化の手順

- a. カーネル回帰: カーネルは自由に選ぶ
- b. 回帰変数が正規分布 (ベイズ的思考方)
- c. 測定値(誤差)、予測値なども正規分布で連結:
ガウス過程
- d. 予測値の分布関数のパラメータの統計分布
(つまり、予測値の平均と分散が得られる)

統計学

頻度主義 (frequentism) 統計学

確率が結果 (データ) を決める

ベイズ主義 (Bayesianism) 統計学

データは定数、確率分布 (確率分布関数に含まれる定数) が分布をもつ

客観確率 (objective probability)

実験または理論的考察 (思考実験) から求められ、客観的な観測結果と比較できるランダムな事象についての確率

主観確率 (subjective probability)

人間の主観的な信念あるいは信頼の度合

回帰のベイズ的思考方

観測データ (x_i, y_i) が以下のようにあらわされるとする。

$$y_i = w_0 + w_1 x_i$$

頻度主義: 「 w_i は決定していて、
(誤差があつて) y_i に不確定性(確率密度分布)がある」

ベイズ主義: 「データ (x_i, y_i) は決定している。 w_i に確率密度分布がある」

・ w_i の分布が決まれば、 w_i をサンプリングすれば $y_i = w_0 + w_1 x_i$ が計算できる。

例: w_i が正規分布であれば、正規分布の乱数によりサンプリングできる

$$\text{独立な正規分布 (平均 } \mu \text{、分散 } \sigma): p(\mathbf{w}) = \frac{1}{2\pi\sqrt{\sigma_0^2\sigma_1^2}} \exp\left(-\sum \frac{(w_i - \mu_i)^2}{2\sigma_i^2}\right)$$

注: w_i が独立であるとして、共分散を0にしている

1. **最初は w_i に関する情報は無い: 主観確率**

$$\text{適当な事前分布を仮定する: } p(\mathbf{w}) = \frac{1}{2\pi\sigma^2} \exp\left(-\sum \frac{(w_i - \mu_i)^2}{2\sigma^2}\right) = N(\mathbf{w}; \mu, \sigma)$$

w_i の分布と回帰

Launcher2023 - ../config/Launcher tutorial.list - Launc...

ファイル ツール

設定 設定ファイル編集 ja 終了

ランチャ 開発者用 ビュー

外部プログラム
----- Help -----
ヘルプ
インストールマニュアル等
----- Data analysis -----
ベイズ最適化(PHYSBO)
フィッティング
スペクトル解析
2023/1/31チュートリアル (回帰)
2023/2/17チュートリアル (回帰・機械学習)
2023/3/6チュートリアル (非線形最適化)
----- Links -----
リンク

ファイル: 選択

引数:

GUI	?	GUI(exe)	?
	×		×
HP	×	github	×
マニュアル	×	github doc	PDF
引用	×		×
	×		×
	×	tutorial:gp.py	?
edit CLI	×	edit GUI	?

cmd(org): \$(start) \$(file_path)

cmd(conv): start D:\%tkProg%tkProg%tkprog_tutorial%o 実行

message:

wi sampling: configure

係数に分布がある場合の回帰直線の分布を体験できます。
グラフの上でクリックすると、n_add分の係数のサンプリングを行い、
対応する直線を点で描画します

w0c: 1 初期値 w0 の平均値

w1c: 1 初期値 w1 の平均値

sigma: 5.0 初期値 w_i の分布幅

wmin: -15 初期値 分布を描画する w_i の下限

wmax: 15 初期値 分布を描画する w_i の上限

n_add: 100 初期値 クリック1回でのサンプリング数

OK Cancel

w_i の分布と回帰結果: 広い分布

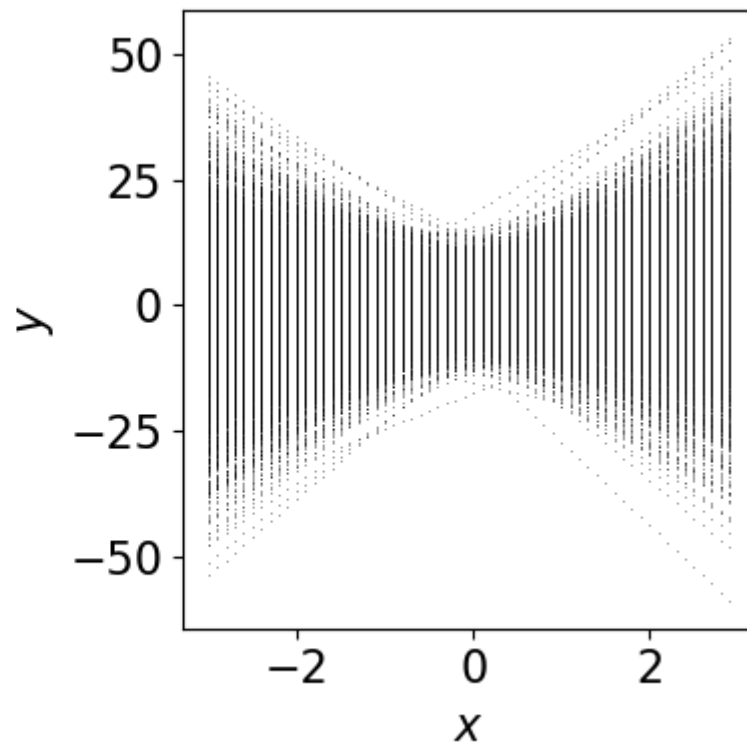
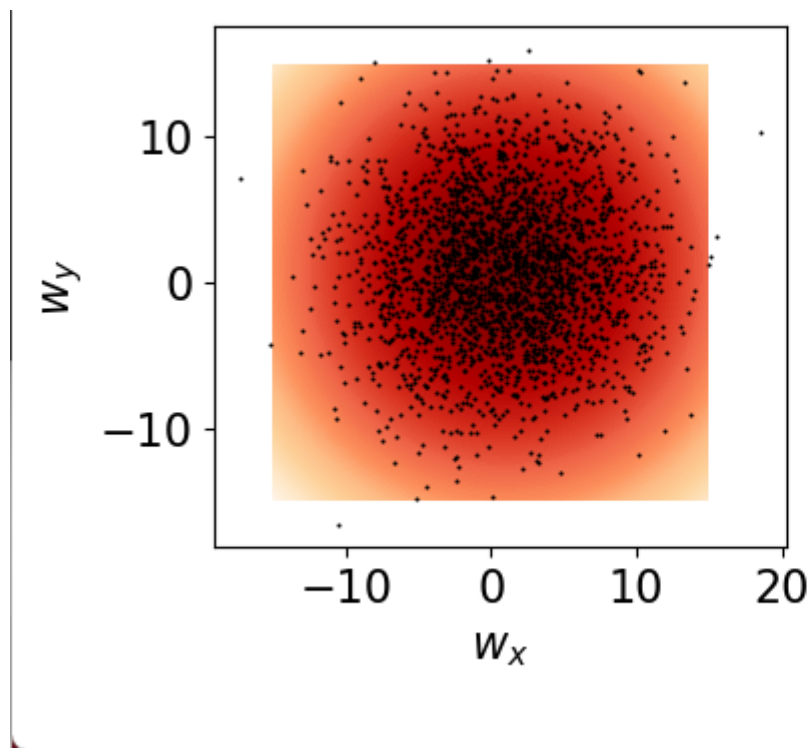
[tkProg]¥tkprog_tutorial¥gp¥gp.py

python gp.py

$$y = w_0 + w_1 * x$$

w_i : 中心 (1, 1)、 $\sigma = 5.0$

マップの強度は対数を取っている



ベイズ的: データを使って w_i の分布を狭める

1. 最初は w_i に関する情報は無い: **事前分布 (主観確率)**

適当な事前分布を仮定する:
$$p(\mathbf{w}) = \frac{1}{2\pi\sqrt{\sigma_0^2\sigma_1^2}} \exp\left(-\sum \frac{(w_i - \mu_i)^2}{2\sigma_i^2}\right) = N(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\sigma})$$

2. データ $D = (x_i, y_i)$ を追加: w_i の範囲を狭めていく

事後分布 (データ D が与えられた場合に w の確率密度分布がどうなるか):

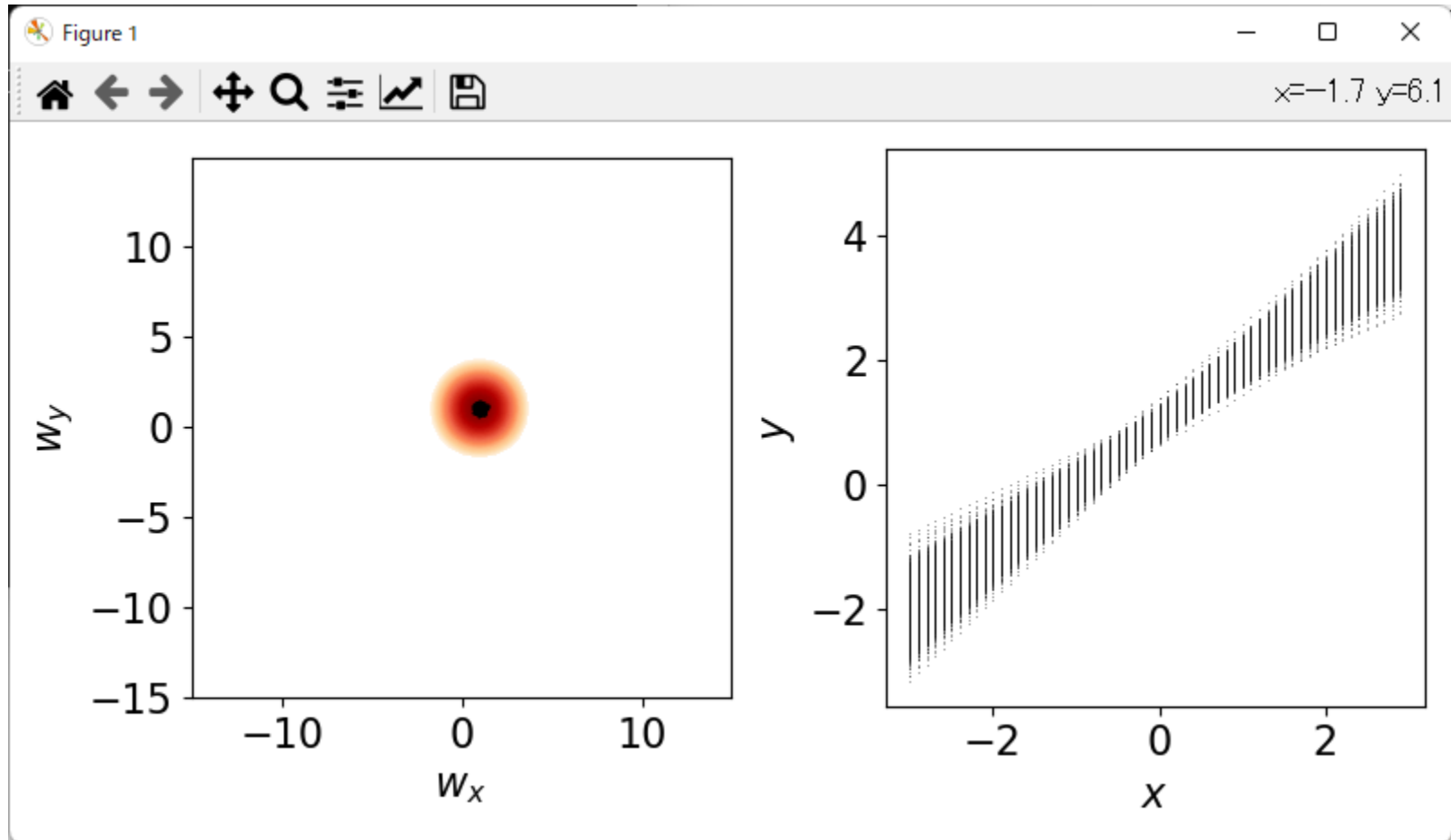
$$p(\mathbf{w}|D) = \frac{1}{2\pi\sqrt{\sigma_0'^2\sigma_1'^2}} \exp\left(-\frac{(\mathbf{w} - \boldsymbol{\mu}')^2}{2\boldsymbol{\sigma}'^2}\right) = N(\mathbf{w}; \boldsymbol{\mu}', \boldsymbol{\sigma}')$$

w_i の分布と回帰結果: 狭い分布

[tkProg]¥tkprog_tutorial¥gp¥gp.py

python gp.py 1 1 0.1

w_i : 中心 (1, 1)、 $\sigma = 0.1$



多変量ガウス分布と線形変換

D 次元の確率変数ベクトル x の正規分布

$$p(x) = \frac{1}{\sqrt{2\pi}^D \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} (x - \langle x \rangle)^T \Sigma^{-1} (x - \langle x \rangle)\right)$$

$\langle x \rangle$: x の平均

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1D} \\ \sigma_{21} & \sigma_{22} & & \\ \vdots & & \ddots & \\ \sigma_{D1} & & & \sigma_{DD} \end{pmatrix} : \text{共分散行列}$$
$$\sigma_{ij} = \frac{\sum (x_i - \langle x \rangle)(x_j - \langle x \rangle)}{n - 1}$$

$|\Sigma|$: Σ の固有値

公式1: 変数変換 $y - \langle y \rangle = A(x - \langle x \rangle)$

$$p(y) \propto \exp\left(-\frac{1}{2} (y - \langle y \rangle)^T A^{-1T} \Sigma^{-1} A^{-1} (y - \langle y \rangle)\right)$$

$$p(y) = \frac{1}{\sqrt{2\pi}^D \sqrt{|\Lambda|}} \exp\left(-\frac{1}{2} (y - \langle y \rangle)^T \Lambda^{-1} (y - \langle y \rangle)\right)$$

$$\Lambda^{-1} = A^{-1T} \Sigma^{-1} A^{-1}$$

x が正規分布をすれば、線形変換しても正規分布になる

w_i の分布をどうやって決めるか

$y = w_0 + w_1 x = \mathbf{w}^T \mathbf{x}$ データ (x_i, y_i)

$$p(\mathbf{w}|\Sigma^{-1}) = \frac{1}{\sqrt{2\pi}^D \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} \mathbf{w}^T \Sigma^{-1} \mathbf{w}\right) \quad (\mathbf{w} \text{ は平均 } 0、\text{共分散 } \Sigma \text{ の正規分布と仮定})$$

誤差ベクトル $\boldsymbol{\varepsilon} = \mathbf{w}^T \mathbf{x}_i - y_i$ も正規分布とする (平均 0、共分散 Λ)

$$p(\boldsymbol{\varepsilon}|\Lambda^{-1}) = \frac{1}{\sqrt{2\pi}^D \sqrt{|\Lambda|}} \exp\left(-\frac{1}{2} \boldsymbol{\varepsilon}^T \Lambda^{-1} \boldsymbol{\varepsilon}\right)$$

最尤推定法: 尤度関数 $p(\boldsymbol{\varepsilon}|\Lambda^{-1})$ を最大化する \mathbf{w} 、 Λ を決めればよいので、

対数尤度 $-\frac{1}{2} \sum (\mathbf{w}^T \mathbf{x}_i - y_i)^T \Lambda^{-1} (\mathbf{w}^T \mathbf{x}_i - y_i)$ を最大化

$\Rightarrow \sum (\mathbf{w}^T \mathbf{x}_i - y_i)^T \Lambda^{-1} (\mathbf{w}^T \mathbf{x}_i - y_i)$ を最小化する

Λ^{-1} を定数と考えて \mathbf{w} を最適化すると

線形最小二乗法と同じになってしまい、意味がない

- \mathbf{w} に正規分布を仮定し**ガウス過程**を利用すると、
正規分布の分散に出てくる \mathbf{w} を**平均化** (周辺化) して
 Λ^{-1} を求める問題にできる

尤度関数

事象 (x_k) が起こる確率を、既知のパラメータ (a_k) の確率密度関数 (PDF)

$$P(X = x_i | a_k) = \prod_i \left\{ \frac{1}{\sqrt{2\pi\sigma_i}} \exp \left[-\frac{\varepsilon_i(x_i | a_k)^2}{2\sigma_i^2} \right] \right\} = \prod_i \left(\frac{1}{\sqrt{2\pi\sigma_i}} \right) \cdot \exp \left[-\sum_i \frac{\varepsilon_i(x_i | a_k)^2}{2\sigma_i^2} \right]$$

などとする。 ($\varepsilon_i(x_i | a_k)$ は誤差。 ($x_i | a_k$) は、 x_i が確率変数で a_k がパラメータであることを示す)

頻度主義の考え方:

データ x_i が得られる確率は $P(X = x_i | a_k)$ 。パラメータ a_k は先に定まっている

ベイズ的な考え方:

データ x_i は確定した事象として考え、

確率分布 (確率変数 a_k) が確率密度関数を持つ と考える

$X = (x_i)$ がわかっているとし、上記の確率密度関数を

パラメータ (a_k) がどれだけ尤もらしいか (尤度) を表す確率密度関数とみなし、

変数 (a_k) の関数として

$P(a_i) = P(x_i | a_i)$ を尤度関数と呼ぶ

(x_i がパラメータ(データ)で a_k が確率変数)

最小二乗法の統計学的基盤: 最尤推定法

最尤推定法

誤差 $\varepsilon_i = f(x_i, a_i) - y_i$ が分散 σ_i の正規分布に従うとする。

データ (x_i, y_i) に対するパラメータ (a_i) の尤度関数は

$$P(a_i) = \prod_i \left(\frac{1}{\sqrt{2\pi\sigma_i}} \right) \cdot \exp \left[- \sum_i \frac{\varepsilon_i(x_i|a_k)^2}{2\sigma_i^2} \right]$$

尤度を最大化するパラメータ a_i を求めるのが「最尤推定法」。

$$\max P(a_i) = \max \ln P(a_i)$$

$$= \min \sum_i \frac{\varepsilon_i^2}{\sigma_i^2}: \text{最小二乗法に一致する}$$

ガウス過程 (Gauss process)

例: 観測データ $\{x_i, y_i\}$ を基底関数 $\phi_k(x_i)$ で展開

観測データに誤差 ε_i があることを考慮すると、

$$t_i = \sum_k w_k \phi_k(x_i)$$

$$y_i = t_i + \varepsilon_i$$

とあらわされるような場合。

このような一連の関数 w_k 、 t_i 、 y_i 、 ε_i が
正規分布に従うとき、これを**ガウス過程**と呼ぶ

Wikipedia:

確率過程 $\{X_t\}_{t \in T}$ は、任意に(有限個の) X_{t_1}, \dots, X_{t_k} を選んで作った
線型結合(あるいはより一般に、 $\{X_t\}_{t \in T}$ を標本関数 X_t 全体からなる
連続濃度の函数空間と見たときの、任意の線型汎関数が正規分布に従うとき、
ガウス過程という。

ガウス過程を利用したベイズ推定の流れ

1. 観測データ $\{x_i, y_i\}$

$$y_i(x_i) = t_i(x_i) + \varepsilon_i \quad \varepsilon_i: \text{誤差}$$

2. 予測値を基底関数 $\varphi(x_i)$ で展開

$$t_i = \sum_{k'} w_{k'} \varphi(x_{i,k'})$$

3. カーネルトリックを用いて計算を簡単化

カーネルは、対称性・正定値性を満たしていればなんでもいい

4. w_k 、 ε_i が正規分布に従うとする: y_i は ε_i と同じ正規分布に従う

分散にでてくる w_k で平均を取り、 w_k を消去。

カーネル回帰ではデータ数と同数の w_k の計算、逆行列の計算が必要だったが、不要になる。

平均、分散を求めることを目的にするため、計算が軽くなる

5. 最尤推定法により、 w_k の平均、分散が求まる

6. ベイズ推定により、 y_i を推定

w_k の具体的な値は使わない

正規分布の表記と公式

平均ベクトル μ 、共分散行列 Σ の正規分布を

$$N(\mu, \Sigma) = N(\text{平均, 共分散行列}) = \frac{1}{\sqrt{2\pi}^D \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}\right)$$

$$N(\mathbf{x}|\mu, \Sigma) = N(\text{確率変数|平均, 共分散行列})$$

と書く（ところどころ $N(\mu, \Sigma^{-1})$ で書いたりもするので、柔軟に対応してください）。

公式1: 線形変換: $y = Ax$ の変化により、

$$p(\mathbf{y}) = N\left(\langle \mathbf{y} \rangle, \left((A^{-1})^T \Sigma^{-1} A^{-1}\right)^{-1}\right)$$

同時分布: ベクトル x が正規分布に従うとき、すべての x_i の同時分布である

$$p(\mathbf{x}) = N(\mu, \Sigma)$$

公式2: 条件付分布: $x = (x_1, x_2)$ のうち x_1 を固定したときに 残りの x_2 が従う分布

$$p(x_2|x_1) = N\left(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

となる (x_1 について積分を取る (周辺化))

基底関数 (特徴ベクトル) で展開、 w の消去

- x の特徴ベクトルを $\varphi(x)$ としてモデル t を展開

$$(t) = (\sum_k w_k \varphi_k(x^{(i)})) = \Phi \mathbf{w}$$

$$\Phi = (\Phi_{ij}) = (\varphi_j(x^{(i)})): \text{計画行列}$$

- \mathbf{w} が正規分布 $N(0, \lambda^2 I)$ に従うとする。
 - 公式1より $t = \Phi \mathbf{w}$ は正規分布 $N(\Phi \mathbf{w}, \mathbf{w}^T \Phi^T \Phi \mathbf{w})$ に従う。
 - $\mathbf{w}^T \Phi^T \Phi \mathbf{w}$ の平均を取ると、

(Φ はデータ $x^{(i)}$ が与えられると定数行列なので)

$$E(\Phi \mathbf{w}) = \Phi E(\mathbf{w}) = 0$$

$$E(\mathbf{w}^T \Phi^T \Phi \mathbf{w}) = \Phi^T \Phi E(\mathbf{w}^T \mathbf{w}) = \lambda^2 \Phi^T \Phi。$$

つまり、 $t = \Phi \mathbf{w}$ は正規分布 $N(0, \lambda^2 \Phi^T \Phi)$ に従う。

※ \mathbf{w} について平均を取った(周辺化)ため、

分布関数から w が消えている

カーネルトリック

- $t = \Phi w$ は正規分布 $N(0, \lambda^2 \Phi^T \Phi)$ に従う。
- カーネルトリック: 分散 $\lambda^2 \Phi^T \Phi$ をカーネル K で置き換える
 $\Rightarrow t = \Phi w$ は正規分布 $N(0, K)$ に従う

$$K = (K_{ij}) = \lambda^2 \Phi^T \Phi$$

$$K_{ij} = \lambda^2 \varphi(x^{(i)})^T \varphi(x^{(j)}): \text{カーネル関数 } k(x^{(i)}, x^{(j)}) \text{ で置き換え}$$

カーネル関数の例:

$$k(x^{(i)}, x^{(j)}) = \theta_1 \exp\left(-\frac{(x^{(i)} - x^{(j)})^2}{\theta_2}\right)$$

動径基底関数 (radial basis function, RBF) / ガウスカーネル

θ_1, θ_2 はハイパーパラメータ

カーネルは任意のものを使える (ガウス関数でなくてもよい)

誤差(ノイズ)とRidge回帰の関係

- 観測値 y_i には誤差(ノイズ) ε_i が含まれる。モデルを t_i とする。

$$y_i = t_i + \varepsilon_i$$

$$(t_i) = \left(\sum_k w_k \varphi_k(\mathbf{x}^{(i)}) \right) = \Phi \mathbf{w}$$

ε_i は正規分布 $N(0, \sigma^2 I)$ に従い、互いに独立(共分散は 0)とする

- y の確率分布は $y_i = t_i + \varepsilon_i$ より (t_i は x が与えらると定数なので、分散は 0)

$$p(\mathbf{y}|\mathbf{t}) = N(\mathbf{t}, \sigma^2 I)$$

- モデル $\mathbf{t} = \Phi(\mathbf{x}^{(i)})\mathbf{w}$ は正規分布 $N(\mathbf{0}, K)$ に従う

- x が与えられた時の y の分布関数は t について平均を取って

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \int p(\mathbf{y}, \mathbf{t}|\mathbf{x}) d\mathbf{t} = \int p(\mathbf{y}|\mathbf{t}) p(\mathbf{t}|\mathbf{x}) d\mathbf{t} \\ &= \int N(\mathbf{y}|\langle \mathbf{t} \rangle, \sigma^2 I) N(\mathbf{t}|\langle \mathbf{t} \rangle, K) d\mathbf{t} \end{aligned}$$

ガウス関数のコンボリューションなので

$$p(\mathbf{y}|\mathbf{x}) = N(\mathbf{y}|\langle \mathbf{t} \rangle, K + \sigma^2 I)$$

Ridge回帰で、正則化係数 $p\alpha$ が σ^2 に対応している
ノイズを考慮することで正則化が自動的に入る
 σ^2 もハイパーパラメータ

- 以下、 $K + \sigma^2 I$ をカーネル $k(\mathbf{x}, \mathbf{x}')$ と書き換える

学習データ (x_i) ($i = 1, 2, \dots, N$) と予測

- N 個の観測値 $D = (x_i)$ ($i = 1, 2, \dots, N$)
 y, t が正規化されているとすると、
 y, t は分散 $k(x, x')$ のガウス過程 $GP(0, k(x, x'))$ から生成されている
- x^* に対する未知の y^* を予測する。観測データ y と合わせた (y, y^*) も y と y^* の同時分布として、同じガウス過程に従う。

$$\begin{pmatrix} y \\ y^* \end{pmatrix} \sim N \left(0, \begin{pmatrix} K & k^* \\ k^{*T} & k^{**} \end{pmatrix} \right)$$

$$k^* = (k(x^*, x_1), k(x^*, x_2), \dots, k(x^*, x_N))$$

$$k^{**} = (k(x^*, x^*))$$

- 条件付分布の公式2より

$$p(x_2 | x_1) = N(\mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (x_1 - \mu_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$$

D と x^* が与えられたときに y^* がとる条件付確率分布は

$$p(y^* | x^*, D) = N(k^{*T} K^{-1} y, k^{**} - k^{*T} K^{-1} k^*)$$

予測値の平均値: $k^{*T} K^{-1} y$

予測値の分散 : $k^{**} - k^{*T} K^{-1} k^*$

$$\text{RBFの場合: } k(x^{(i)}, x^{(j)}) = \theta_1 \exp \left(-\frac{(x^{(i)} - x^{(j)})^2}{\theta_2} \right) + \delta_{ij} \theta_3$$

ハイパーパラメータの推定

学習データの確率:カーネルにハイパーパラメータ θ が含まれる

$$p(\mathbf{y}|\mathbf{x}, \theta) = N(\mathbf{y}|0, K(\theta))$$

最尤推定法:

$$p(\mathbf{y}|\mathbf{x}, \theta) = N(\mathbf{y}|0, K(\theta)) = \frac{1}{\sqrt{2\pi}^D \sqrt{|K(\theta)|}} \exp\left(-\frac{1}{2} \mathbf{y}^T K(\theta)^{-1} \mathbf{y}\right)$$

の対数

$$\log p(\mathbf{y}|\mathbf{x}, \theta) = -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |K(\theta)| - \frac{1}{2} \mathbf{y}^T K(\theta)^{-1} \mathbf{y}$$

が最大になるように、つまり、

$$\log |K(\theta)| + \mathbf{y}^T K(\theta)^{-1} \mathbf{y}$$

が最小になるようにハイパーパラメータを決定する。

=> Cross Validationや非線形最小化 (勾配法など) を使う