

触覚情報に基づく能動的物体認識のための探索行動設計 —相互情報量を規範とする最適制御によるアプローチ—

田中大介 松原崇充 杉本謙二 (奈良先端大)

1. はじめに

人間との共存環境下で稼働するロボットにとって、環境認識の効率化は重要な課題である。本研究では、環境認識問題の一つとして、触覚情報に基づく物体認識問題 (図 1) を取り扱う。触覚情報はロボットが物体に触れることで初めて得られ、触り方に依存して得られる情報が大きく変わる。また、最適な触り方は物体に関して現在までに得られている情報にも依存する。このため、効率的に認識するためには状況に応じた探索行動の設計が重要となる。

これまで様々な探索行動の設計法が提案されているが、行動の計画問題とロボットの制御問題とが別々に取り扱われてきた (e.g. [1-4])。この結果、ロボットにとって追従不可能な行動が計画される、もしくは不必要に大きなトルクが入力されることにより低エネルギー効率で危険な行動となる可能性があった。

我々は、この探索行動設計問題を最適制御問題として取り扱う。最適制御問題は、望む挙動をコスト関数で表し、そのコストが最小となるような制御則を求める問題であり、行動計画と制御問題を同時に解くことが可能である。探索行動設計問題で望む挙動は、対象に対する不確実性を下げることであり、加えてコスト関数の中にエネルギーやトルク制約を考慮することで、省エネルギーでかつ高コンプライアンスな最適探索行動が得られる。

2. 準備

2.1 逐次能動学習による能動的物体認識

我々は、物体認識問題をパラメータ推定問題 [3] として取り扱う。対象物体はそれぞれ固有のパラメータ (以降、物体パラメータと呼称する) が割り当てられており、探索行動により得られた触覚情報から物体パラメータを推定していく。

一般的な能動的物体認識の処理は以下のとおりである。

1. 物体パラメータの初期値をある確率分布として与える (この確率分布を信念 (belief) と呼ぶ)
2. 現時点の信念に基づく最適な探索行動を設計する
3. 設計された探索行動を対象に実行し観測 (触覚情報) を得る
4. 得られた観測から信念を更新する
5. 物体に関して確信が持てるまで (例えば信念の分散が十分小さくなるまで) 2. から繰り返す
6. 信念の平均の最近傍の物体として認識結果を確定する



図 1 触覚情報に基づく物体認識問題

2.2 行動の最適性の定義

本稿では相互情報量を用いて探索行動の最適性を定義する。その時々で最適な探索行動は物体パラメータ θ の確率分布で表される物体の信念 $p(\theta)$ に依存する。相互情報量 $I[\theta, y|x]$ は、ある状態 x にて観測 y が与えられた時の θ の不確実性が減った量、すなわち得られた情報の量を定量的に評価するもので、以下で与えられる。

$$I[\theta, y|x] \triangleq \text{KL}(p(\theta, y|x) \| p(\theta)p(y|x)) \\ = \iint p(\theta, y|x) \log \frac{p(\theta, y|x)}{p(\theta)p(y|x)} dy d\theta \quad (1)$$

この量を最大化する状態 (系列) にシステムを制御することで、パラメータ推定に有効な観測を得る。

3. 提案法: 最適制御に基づく逐次探索行動設計

3.1 対象システムと問題設定

物体認識に用いるロボットとセンサは、それぞれ以下の非線形状態方程式と観測方程式で表される離散時間システムモデルで表されるとする。

$$\mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t) \quad (2)$$

$$\mathbf{y}_t = \mathbf{g}(\mathbf{x}_t, \theta) + \epsilon_t \quad (3)$$

ここで、 $\mathbf{x} \in \mathbb{R}^{d_x}$ は観測可能なロボットの状態、 $\mathbf{u} \in \mathbb{R}^{d_u}$ はロボットへの入力、 $\mathbf{y} \in \mathbb{R}^{d_y}$ はロボットのセンサからの観測、 $\epsilon \in \mathbb{R}^{d_y}$ は平均 0、共分散行列 $\Sigma_\epsilon = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_{d_y}^2\}$ の d_y 次元ガウス分布に従う観測ノイズである。ロボットへの入力は、 $\mathbf{u}_{\min} \leq \mathbf{u} \leq \mathbf{u}_{\max}$ に制限されるとする。また、 $\theta \in \mathbb{R}^{d_\theta}$ は物体パラメータである。

このシステムにより前節で説明した能動的物体認識の処理を行う。処理 2. はロボットのダイナミクス (2) の

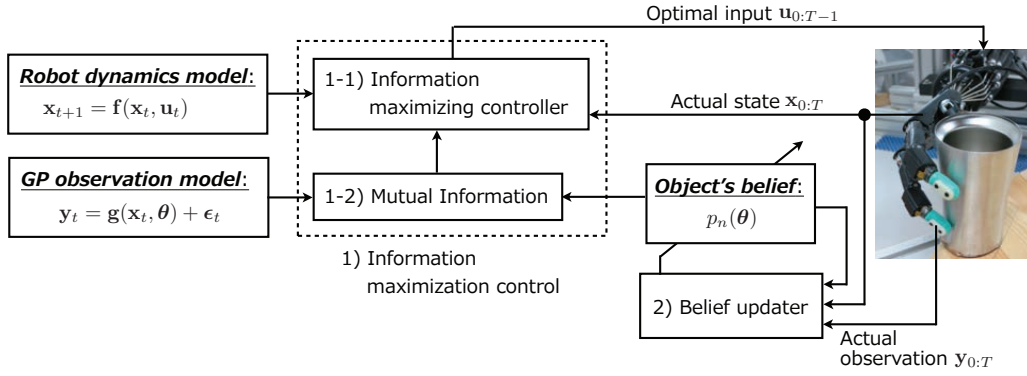


図 2 提案法の概要図．1-1) 得られる情報量を最大にする制御器，1-2) ガウス過程回帰による観測モデルから計算される相互情報量から構成される 1) 情報量最大化制御，2) 得られた観測からの信念の更新の 2 要素で構成される．

制約を考慮した上で，コストの総和が最小となる状態系列に対応する行動系列を探索する問題，つまり相互情報量を規範とする最適制御問題に帰着させる (3.2.1)．このコストには相互情報量 (1) を用いるが，2 重積分を簡単に計算する工夫として，ガウス過程回帰 (Gaussian Process Regression) を用いた観測モデルを利用する (3.2.2)．処理 4. では得られた状態系列から信念を更新する (3.3)．

3.2 情報量最大化制御

3.2.1 局所化に基づく最適制御問題の近似解法

まず，最適制御問題を定式化する．これは，式 (2) のダイナミクスに従うシステムに対して，次式で与えられる評価関数を最小化するような制御則 $u_t = \pi(t, x_t)$ を定める問題である．

$$J_{0:T} = h(x_T) + \sum_{t=0}^{T-1} \ell(t, x, \pi(t, x))$$

ただし， $h(x_T) \geq 0$ は終端コスト， $\ell(t, x, u) \geq 0$ はランニングコストを表す．探索行動設計に対しては，これらコスト関数はエネルギー効率と行動の最適性を用いて与える．実際の与え方については後述する．この問題は，システムが線形でかつ 2 次コストであれば Riccati 方程式に基づく方法で最適制御則 π^* を導出可能であるが，ここではシステムを初期入力列 $\bar{u}_{0:T-1}$ を入力した時の初期状態列 $\bar{x}_{0:T}$ 周りで局所的に線形システムを構築し近似的に解く iterative Linear Quadratic Regulator (iLQR, [5]) を用いる．

この手法を用いる利点は，局所的なフィードバック制御を構築可能なことである．iLQR によりフィードバックゲイン L_t が与えられ，実際に与える入力 u_t は $u_t = \bar{u}_t + L_t(x_t - \bar{x}_t)$ で得られる [5]． \bar{x}_t, \bar{u}_t はそれぞれ iLQR で得られた最適な状態と入力であり， x_t はロボットの実際の状態である．このフィードバックにより，ダイナミクスに物体との相互作用が考慮されていない，もしくはモデル化誤差が存在する場合にも，その状態での最適な入力を生成可能である．この場合にダイナミクスを正確にモデル化することは困難であるため，このようなフィードバックが有効に機能することが期待される．

3.2.2 相互情報量による逐次コストの定義

観測モデル (3) は，ガウス過程回帰 [6] により y の各次元 $y_a = g_a(x, \theta) + \epsilon_a, a \in \{1, 2, \dots, d_y\}$ がモデル化されているとする．訓練データ $\mathcal{D} = \{x^{(i)}, \theta^{(i)}, y^{(i)}\}_{i=1}^N$ が与えられた時，予測分布は x, θ, y_a のそれぞれの訓練データ集合 $\mathbf{X}, \Theta, \mathbf{y}_a$ を用いて以下のガウス分布で与えられる [6]．

$$p(y_a | x, \theta, \mathbf{X}, \Theta, \mathbf{y}_a) = \mathcal{N}(\mu_a(x, \theta; \mathbf{X}, \Theta, \mathbf{y}_a), s_a^2(x, \theta; \mathbf{X}, \Theta, \mathbf{y}_a))$$

本稿では μ_a, s_a^2 の計算に用いられるカーネル関数は以下で表される 2 乗指数カーネル関数 (squared exponential kernel function) を仮定する．

$$k_a(z, z') = \alpha_a^2 \exp\left(-\frac{1}{2}(z - z')^T (\mathbf{H}_a^z)^{-1} (z - z')\right)$$

ここで $z = [x^T, \theta^T]^T \in \mathbb{R}^{d_z}, d_z = d_x + d_\theta$ を用いた．また α_a^2 は g_a の分散で，モデルの簡単化のために x と θ は独立であると仮定し， $\mathbf{H}_a^z = \text{block diag}\{\mathbf{H}_a^x, \mathbf{H}_a^\theta\}$ と定めた． $\mathbf{H}_a^x, \mathbf{H}_a^\theta$ は共に正の要素を対角に持つ対角行列である．ハイパーパラメータは $\gamma_a = \{\alpha_a^2, \sigma_a^2, \mathbf{H}_a^x, \mathbf{H}_a^\theta\}$ であり，訓練データに対する周辺対数尤度関数最大化により得る [6]．

ここからは， θ の信念をガウス分布 $p(\theta) = \mathcal{N}(\mu, \Sigma)$ で表す．ある 1 点の x が与えられた時，相互情報量 $I[\theta, y | x]$ は，GP 観測モデルの訓練データとハイパーパラメータを用いて，以下のように表される．

$$I[\theta, y | x] = -\frac{1}{2} \log\left(\frac{\det \tilde{\Sigma}(x)}{\det \Phi(x, x) \det \Sigma}\right)$$

$$\tilde{\Sigma}(x) = \begin{bmatrix} \Phi(x, x) & \mathbf{C}^T(x) \\ \mathbf{C}(x) & \Sigma \end{bmatrix}$$

ただし， $\Phi(x, x') \in \mathbb{R}^{d_y \times d_y}$ の (a, b) 要素 Φ_{ab} と $\mathbf{C}(x) \in \mathbb{R}^{d_\theta \times d_y}$ は以下で与えられる．

$$\Phi_{ab} = \beta_a^T \Lambda_{ab}(x, x') \beta_b - m_a(x) m_b(x')$$

$$+ \delta_{xx'} \delta_{ab} \left(\alpha_a^2 - \text{Tr}(\mathbf{K}_a^{-1} \Lambda_{aa}(x, x')) \right)$$

$$\mathbf{C} = \Psi(x) - \mu m(x)^T$$

ここで行列 \mathbf{K}_a はカーネル行列であり, その (i, j) 要素 $K_{a,ij}$ は $K_{a,ij} = k_a(\mathbf{z}^{(i)}, \mathbf{z}^{(j)}) + \delta_{ij}\sigma_a^2$ で与えられる. また, $\mathbf{m}(\mathbf{x}) \in \mathbb{R}^{d_y}$ の a 番目の要素は $m_a = \beta_a^T \lambda_a(\mathbf{x})$ であり, $\beta_a = \mathbf{K}_a^{-1} \mathbf{y}_a \in \mathbb{R}^N$ とした. $\Lambda_{ab}(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^{N \times N}$ については付録に示すが, $\lambda_a(\mathbf{x}) \in \mathbb{R}^N$, $\Psi(\mathbf{x}) \in \mathbb{R}^{d_\theta \times d_y}$ については文献 [3] を参照されたい.

3.3 状態系列に基づく信念の更新

現在の信念 $p_n(\theta) = \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ を基にして先に述べた制御法により探索行動が計画後実行され, 時刻 $t = 0$ から $t = T$ までの状態列 $\mathbf{x}_{0:T}^n$ とその時の観測 $\mathbf{y}_{0:T}^n$ が得られたとする. これを用いて信念の更新則は, 文献 [3] で提案されている状態と観測が 1 点ずつ与えられた時の更新則を拡張し, 以下のように与えられる.

$$\begin{aligned} \boldsymbol{\mu}_{n+1} &= \boldsymbol{\mu}_n + \mathbf{C}_n \mathbf{S}_n^{-1} (\mathcal{Y}_n - \mathcal{M}_n) \\ \boldsymbol{\Sigma}_{n+1} &= \boldsymbol{\Sigma}_n - \mathbf{C}_n \mathbf{S}_n^{-1} \mathbf{C}_n^T \end{aligned}$$

ただし, $\mathbf{C}_n \in \mathbb{R}^{d_\theta \times (T+1)d_y}$, $\mathbf{S}_n \in \mathbb{R}^{(T+1)d_y \times (T+1)d_y}$, $\mathcal{Y}_n \in \mathbb{R}^{(T+1)d_y}$, $\mathcal{M}_n \in \mathbb{R}^{(T+1)d_y}$ は以下のように定義した.

$$\begin{aligned} \mathbf{C}_n &= \begin{bmatrix} \mathbf{C}(\mathbf{x}_0^n) & \cdots & \mathbf{C}(\mathbf{x}_T^n) \end{bmatrix} \\ \mathbf{S}_n &= \begin{bmatrix} \Phi(\mathbf{x}_0^n, \mathbf{x}_0^n) & \cdots & \Phi(\mathbf{x}_0^n, \mathbf{x}_T^n) \\ \vdots & \ddots & \vdots \\ \Phi(\mathbf{x}_T^n, \mathbf{x}_0^n) & \cdots & \Phi(\mathbf{x}_T^n, \mathbf{x}_T^n) \end{bmatrix} \\ \mathcal{Y}_n &= \begin{bmatrix} (\mathbf{y}_0^n)^T & \cdots & (\mathbf{y}_T^n)^T \end{bmatrix}^T \\ \mathcal{M}_n &= \begin{bmatrix} (\mathbf{m}(\mathbf{x}_0^n))^T & \cdots & (\mathbf{m}(\mathbf{x}_T^n))^T \end{bmatrix}^T \end{aligned}$$

物体認識は前節で示したように, 行動計画と信念の更新とを繰り返すことで達成される. 例えば終了条件として, 適切な閾値 ϵ を用いて $\|\boldsymbol{\mu}_{n+1} - \boldsymbol{\mu}_n\|^2 < \epsilon$, $\|\boldsymbol{\Sigma}_n\|_F < \epsilon$, もしくは n_{\max} 回繰返すまでと設定することで, 最終的に最近傍の θ に対応する物体との識別結果を得る.

4. 数値シミュレーション

4.1 シミュレーション設定

図 3 に示す 1 リンクアームで, 提案する手法の有効性を検証する. ただし, このアームの可動域は $-\pi/2 \leq q < \pi/2$ とし, 運動方程式は $\Delta t = 0.01[\text{s}]$ でオイラー離散化した. 物体はバネ κ ・ダンパ D を並行に接続した系に基づく押しこみ反力モデルと動摩擦モデルを仮定した. 手先に仮想的に圧力センサが設置されていることとし, 以下の物体押しこみ方向の力 f_1 と擦る方向の力 f_2 を得た.

$$f_1 = -(\kappa \xi_x + D \dot{\xi}_x), \quad f_2 = -\text{sign}(\dot{\xi}_y) \mu' f_1$$

ここで, $\boldsymbol{\xi} = [\xi_x, \xi_y]^T$ は手先のデカルト座標である.

この実験では物体認識問題をダンパ係数 D の推定問題とする. 物体パラメータは $\theta = D$ であり, 探索行動により図 3 の物体の D を推定する. ここからまず観測モデル学習の設定について述べる. バネ定数を

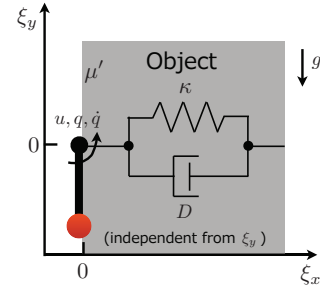


図 3 シミュレーション実験設定. 物体はバネダンパ並列モデル (バネ定数 κ , ダンパ定数 D) で反発力をモデル化し, 表面には動摩擦係数 μ' の摩擦モデルを仮定した. これらは, 物体表面の高さ方向の位置に依存しないと仮定する. アームの先は赤丸で表しており, アームが ξ_x 軸方向と平行のとき $q = 0$ である. また, 図中のアームの姿勢は $\mathbf{x}_0 = [-\pi/2, 0]^T$ に対応している. 入力 u や q, \dot{q} の正方向は反時計回りに定めた.

$\kappa = 1$, 動摩擦係数 $\mu' = 0.5$ と固定し, ダンパ係数 D を $D \in \{1, 3, 5\}$ と変化させた. 観測を $\mathbf{y} = [f_1, f_2]^T$, 状態を $\mathbf{x} = [q, \dot{q}]^T$ として, それぞれの D に対して, $-\pi/2 \leq q \leq \pi/2, -15 \leq \dot{q} \leq 15$ について観測 \mathbf{y} を得て, 訓練データ D を構成した. データ数は各 D に対して 225 個, よって総データ数は $N = 675$ である. 信念の更新回数は $n_{\max} = 10$ とした.

次に最適制御則の導出に関する設定を述べる. 初期状態は $\mathbf{x}_0 = [-\pi/2, 0]^T$ とし, 1 行動の長さ $T = 100$ として初期入力列は $u_t = 3, t = 0, 1, \dots, T-1$ とした. また $p_n(\theta)$ に基づく探索行動についての逐次コスト $\ell_n(t, \mathbf{x}, u)$ は以下のように設定した.

$$\begin{aligned} \ell_n(t, \mathbf{x}, u) &= \exp\left(-\rho^{(n)} \mathbb{I}[\theta, \mathbf{y} | \mathbf{x}]\right) \\ &\quad + 1000 \exp(-(15 - \dot{q})) \\ &\quad + 1000 \exp(-(15 + \dot{q})) \end{aligned}$$

ただし, $\rho^{(n)} > 0$ は最大化問題を最小化問題に変更するために用いた定数であり, 実験的に $\rho^{(n)} = 10 \exp(-0.2n^2)$ と選んだ. また終端コストは $h(\mathbf{x}) = \ell_n(T, \mathbf{x}, 0)$ とした. このコスト関数の第 2 項と第 3 項は, アームの角速度が訓練データの範囲を逸脱しないように入れた罰則項である. さらに, 初期信念として, $p_0(\theta) = \mathcal{N}(3, 5^2)$ を与えた. この平均は, 訓練データの θ の平均に等しい.

4.2 シミュレーション結果

観測モデルから計算される初期信念 $p_0(\theta)$ に対する相互情報量の分布を図 4 に示す. この分布について考察する. ダンパの力は押しこみ方向の速度 $\dot{\xi}_y$ に依存して変化する. $q = 0$ のときには $\dot{\xi}_y$ は 0 になるため, ダンパ係数に関する情報は得られない. 対して, $|q|$ が $\pi/2$ に近い部分では, 大きな $\dot{\xi}_y$ が得られるため, ダンパ係数に関する情報は得られやすい. よって, この相互情報量の分布は妥当だと考える.

次に, $D = 1$ の物体に対する認識結果を図 5 に示す. 横軸は信念の更新回数であり, 各信念の更新には $\mathbf{x}_{0:T}$

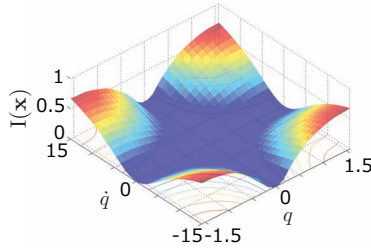


図4 相互情報量の分布．初期信念 $p_0(\theta) = \mathcal{N}(3, 5^2)$ に対して得た．

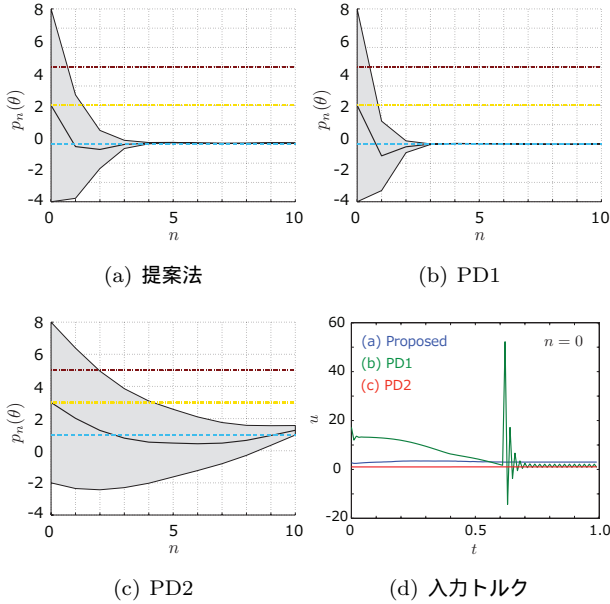


図5 パラメータの推定結果．(a-c) 横軸は信念の更新回数で， $p_0(\theta) = \mathcal{N}(3, 5^2)$ である．真値 $D = 1$ は水色の点線で表した．黒線と灰色の領域は， $p_n(\theta)$ の平均と標準偏差を表す．これらの探索行動は，(a) 提案法，(b) 高いゲインのPDコントローラで，(c) 低いゲインのPDコントローラでそれぞれ得られた．(d) (a-c) のそれぞれの探索行動の入力トルクの関係を表す．

が用いられている．図5(a)は提案法により得た探索行動により信念の更新を行った結果であり， $n = 4$ の時点でほぼ結果が収束していることが確認できる．この結果と比較するため， $\mathbf{x}_d = [\pi/2, 15]^T$ に向かうPDコントローラを設計し，探索行動を行った．この目標値は，図4から設定した．図5(b)は高いゲインに設定した時で，提案法よりも少ない更新回数で収束したように見えるが，図5(d)に示すように高いトルクが入力されており危険である．これに対して低いゲインに設定した場合には入力トルクは小さいが図5(c)に示すように収束は遅い．これらの比較から，提案法の有効性が確認できた．

5. おわりに

本稿では触覚情報に基づく能動的物体識別のための探索行動設計問題に対して，相互情報量を最大化する最適制御によるアプローチを提案した．今後，図6に示す Shadow Dexterous Hand を用いた実機実験シス



図6 実機実験用の設定

テム [4] を用いて提案法の有効性を確認する．

付録A Λ の定義

訓練データとハイパーパラメータを用いて， Λ_{ab} の (i, j) 要素 $\Lambda_{ab,ij}$ は次のように表される．

$$\begin{aligned} \Lambda_{ab,ij} &= \alpha_a^2 \alpha_b^2 \det \left(\left((\mathbf{H}_a^\theta)^{-1} + (\mathbf{H}_b^\theta)^{-1} \right) \Sigma + \mathbf{I} \right)^{-\frac{1}{2}} \\ &\times \exp \left(-\frac{1}{2} (\boldsymbol{\theta}^{(i)} - \boldsymbol{\theta}^{(j)})^\top (\mathbf{H}_{ab}^\theta)^{-1} (\boldsymbol{\theta}^{(i)} - \boldsymbol{\theta}^{(j)}) \right) \\ &\times \exp \left(-\frac{1}{2} (\boldsymbol{\theta}_{ab}^{ij} - \boldsymbol{\mu})^\top \mathbf{R}_{ab}^{-1} (\boldsymbol{\theta}_{ab}^{ij} - \boldsymbol{\mu}) \right) \\ &\times \exp \left(-\frac{1}{2} (\mathbf{x} - \mathbf{x}^{(i)})^\top (\mathbf{H}_a^{\mathbf{x}})^{-1} (\mathbf{x} - \mathbf{x}^{(i)}) \right) \\ &\times \exp \left(-\frac{1}{2} (\mathbf{x}' - \mathbf{x}^{(j)})^\top (\mathbf{H}_b^{\mathbf{x}})^{-1} (\mathbf{x}' - \mathbf{x}^{(j)}) \right) \end{aligned}$$

ただし， $\mathbf{H}_{ab}^\theta = \mathbf{H}_a^\theta + \mathbf{H}_b^\theta$ であり，

$$\boldsymbol{\theta}_{ab}^{ij} = \mathbf{H}_b^\theta (\mathbf{H}_{ab}^\theta)^{-1} \boldsymbol{\theta}^{(i)} + \mathbf{H}_a^\theta (\mathbf{H}_{ab}^\theta)^{-1} \boldsymbol{\theta}^{(j)}$$

$$\mathbf{R}_{ab} = \left((\mathbf{H}_a^\theta)^{-1} + (\mathbf{H}_b^\theta)^{-1} \right)^{-1} + \Sigma$$

を用いた．

参考文献

- [1] J. A. Fishel and G. E. Loeb: “Bayesian exploration for intelligent identification of textures,” *Frontiers in Neurobotics*, vol.6, no.4, 2012.
- [2] J. Sinapov, C. Schenck, K. Staley, V. Sukhoy, and A. Stoytchev: “Grounding semantic categories in behavioral interactions: Experiments with 100 objects,” *Robotics and Autonomous Systems*, vol.62, no.5, pp.632–645, 2014.
- [3] H. Saal, J.-A. Ting, and S. Vijayakumar: “Active sequential learning with tactile feedback,” *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, pp.677–684, 2010.
- [4] D. Tanaka, T. Matsubara, K. Ichien, and K. Sugimoto: “Object manifold learning with action features for active tactile object recognition,” *Proceedings of the 2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2014. to appear.
- [5] W. Li and E. Todorov: “Iterative linear quadratic regulator design for nonlinear biological movement systems,” *Proceedings of International Conference on Informatics in Control, Automation and Robotics*, pp.222–229, 2004.
- [6] C. E. Rasmussen and C. K. I. Williams: *Gaussian processes for machine learning*. MIT Press, 2006.