

MPEG-Based Vision

笹 木 美樹男

Last modified on Jan. 5, 2013.

最終更新日 : 2013 年 1 月 5 日 20 時 40 分

本稿は現在作成中の個人論文であり、今後も改訂されます。

昨今の様々な PC セキュリティ上の懸念に基づき、完成前にハッキング等によって盗用されることのないよう、あえて事前に Web 公開することで筆者のオリジナリティを主張するものです。

(この方針は「神の関数」の HP 全体の著作物に共通します。)

ご理解ください。

1. はじめに

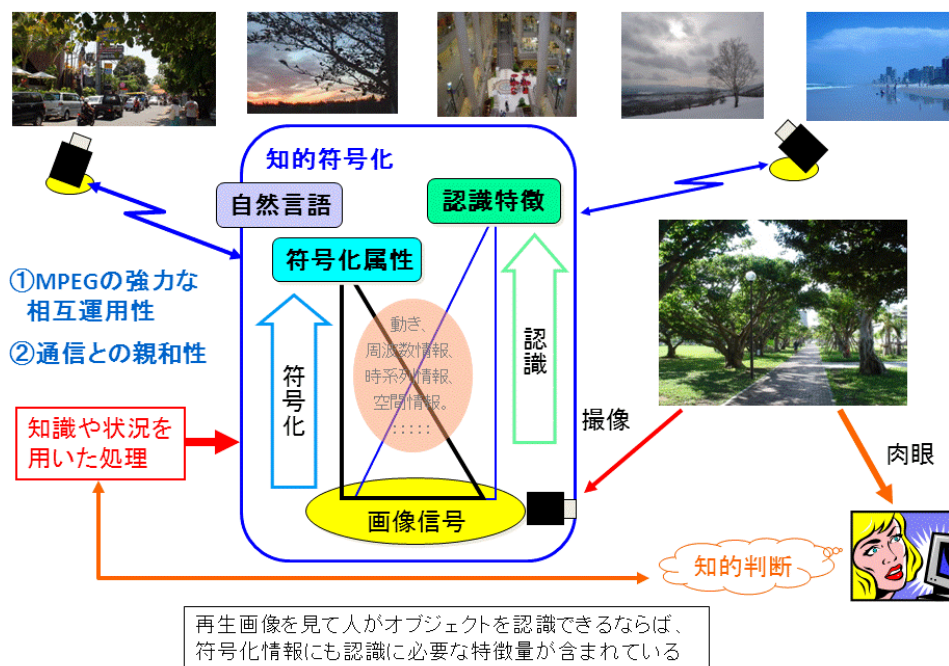
2. 関連研究

3. 理論

3.1 はじめに

本論文で描く理論と実践の概観を図 1 に示す。本論文で提唱する画像認識システムの特徴は以下のとおりである：

- 通信、記録、蓄積、編集、インターネット、クラウドと親和性が高い
- 特に、MPEG の特徴量をターゲットとして上記の親和性を有する
- 認識における記憶・知識処理・予測の重要性に焦点をあて、無記憶条件下での単独の認識器に比べ、複数の認識器の協調や記憶による性能向上を可能にする



プローブ協調と知識処理により符号化画像を認識する

図 3.1 本研究の概観

Fig. 3.1 Targeted overview.

3.2 確信度による認識

3.2.1 確信度ベクトル

各画素ブロックについて、システムは色、テクスチャ、および動き情報をダイレクトに MPEG 符号化データから取り出す。ここで、復号画像の再処理はなくてもよい。そして、ルールベース推論、多変量回帰分析に基づく推論、および簡単な 3D モデルベースの認識が統合され、各ブロック画素を N 個のオブジェクトカテゴリにマッピングする確信

度ベクトルを算出する。

いま、画像中のある画素領域 R （ここでは1つの画素ブロック）に付与するインデックス A_i ($i=1, \dots, N$) の確からしさを確信度 C_i と定義し、

$$\mathbf{C} = (C_1, C_2, \dots, C_i, \dots, C_N)^T \quad (3.1)$$

を確信度ベクトルと呼ぶことにする。

この確信度ベクトルは色やブロック位置、他のオブジェクトとの関係など様々な判定規則で個別に算出した部分的な確信度ベクトルの和として次式で表現する。

$$\mathbf{C} = \mathbf{c}_0 + \sum_{m=1}^M \mathbf{c}_m \quad (3.2)$$

ただし、 \mathbf{c}_0 ：初期条件として与えられる確信度ベクトル（ N 次元列ベクトル）

\mathbf{c}_m ： m 番目（ $m=1, \dots, M$ ）の判定規則に対応する部分的な確信度ベクトル（ N 次元列ベクトル）

である。この確信度ベクトルの最終値において、最大成分に相当するインデックスをマクロブロックに該当するオブジェクトであると判断する。

この手法の利点は従来のオブジェクト認識で用いられてきた特徴空間の分割に基づく手法に比べて、考えられる複数の解釈に対応する確信度という形で認識の状態を保持し、随時ある基準の元で唯一の解釈を動的に選択できるところにある。

3.2.2 確信度ベクトルの推定

確信度ベクトルを推定するための基本的手法としてはルールベースで算出する方法と統計的に推論する方法がある。前者は容易に適用できる反面、様々なシーンやオブジェクトに適応させることが困難であり、応用範囲が限定される。一方、後者は適応性が高い反面、学習操作に手間を要する。したがって、ルールベースと統計的推論の両方をうまく融合することがシステム構築の鍵となる。

3.2.3 確信度の線形回帰モデル

確信度ベクトルは大きく分けて2つのクラスの特徴要因に影響される。1つはシーンの状況に関する特徴要因であり、カメラがシーンを撮影した際の時間や場所、天候などがそれに相当する。これについては後述の予測動作(3.18)において再度詳細に取り上げる。もう1つのクラスは画像に関する特徴要因であり、撮影された画像中のあるブロック画素の色、テクスチャ、動き、2次元位置、推定される3次元位置などの画像特徴に関するものがそれに相当する。いま、特徴要因と確信度ベクトルとの間の因果関係を次のような多変量線形回帰モデルで表現する。

$$\mathbf{C} = \mathbf{c}_0 + \mathbf{W}\mathbf{s} \quad (3.3)$$

ただし、 \mathbf{c}_0 は初期確信度ベクトル（ N 次元列ベクトル）、 \mathbf{W} は特徴要因 \mathbf{s} に対する回帰係数の組を表す N 行 L 列の行列、 \mathbf{s} は特徴要因を表す L 次元列ベクトルである。

3.3 統計的な推論

3.3.1 多変量回帰分析による学習

複数の解釈結果に対する確信度で構成される空間 \mathbf{C} と入力画像や付随情報から得られ

る特徴空間 S との間の回帰関係 W を求める。まず、サンプル映像中の個々のブロック画素の解釈を確信度として教示する。最小 2 乗法を適用すれば、 W^T の推定値 B は次式で計算できる。

$$B = (S^T S)^{-1} S^T Y \quad (3.4)$$

ただし、

$$S = [s(1) \dots s(k) \dots s(K)]^T \quad (3.5)$$

であり、 $s(k)$ は時刻 T_k ($k=1, \dots, K$) における s の値とする。また、 Y は確信度の履歴を表す行列であり、

$$Y = [Y_1 \dots Y_n \dots Y_N] \quad (3.6)$$

である。ここで、 Y_n は n 番目の語彙に対する確信度の時系列を表す K 次元列ベクトルであり、

$$Y_n = (y_n(1), \dots, y_n(k), \dots, y_n(K))^T \quad (3.7)$$

である。更に、 $y(k)$ は時刻 T_k における C の観測値を表す N 次元列ベクトルであり、

$$y(k) = (y_1(k), \dots, y_n(k), \dots, y_N(k))^T \quad (3.8)$$

である。この観測は本稿では人の判定によって行われる。

式(3.2)と式(3.3)により M 種類の状況要因と確信度ベクトルの間の因果関係を表すと次式のようになる：

$$C = c_0 + \sum_{m=1}^M W_m s_m \quad (3.9)$$

ただし、 c_0 は初期確信度を表す N 次元列ベクトルである。 W_m ($m=1, \dots, M$) は要因 s_m に対応する回帰係数行列 ($N \times L_m$) であり、 s_m は L_m 次元列ベクトルを示す。ここで、 $i \neq j$ ならば s_i と s_j の間の相関はないものとした。

以下、式(3.3)～(3.6)と同様にして、

$$B_m \equiv W_m^T = (S_m^T S_m)^{-1} S_m^T Y \quad (m=1, \dots, M) \quad (3.10)$$

ただし、

$$S_m = [s_m(1) \dots s_m(k) \dots s_m(K)]^T \quad (3.11)$$

であり、 $s_m(k)$ は時刻 T_k ($k=1, \dots, K$) における s_m を示す。また、 Y は次式のように、確信度の履歴を表す：

$$Y = [Y_1 \dots Y_n \dots Y_N] \quad (3.12)$$

ただし、 Y_n は K 次元列ベクトルであり、 n 番目の語彙に対応する確信度の履歴を表す。

3.4 シーンの認識

以上により、過去の画像の特徴要因と確信度の履歴を学習する（ここでは MPEG-1 形式のブロック画素単位で）と、新たな画像から抽出した特徴要因から確信度ベクトル C を推定するための回帰関係 W を獲得できる。この C の最大成分に対応するインデックスが各ブロック画素について最適なインデックスであると仮定し、それを第 1 次の認識結果とする。ただし、この処理は後述するターゲットシーンクラスに関して事前に用意した辞書の範囲内で行われる。

3.5 統計的推論手法

時々刻々と移動する車載カメラによる景観認識では次のような性質が望まれる。

- ① 認識率が改善される
- ② 認識できるオブジェクトの種類が増える
- ③ 同じ画素領域に複数の意味を持たせられる

しかしながら、特徴空間をクラスタリングする従来手法では、これらへの対応が困難であった。

3.6 状況ベクトル

時刻 T_k における MPEG 型式の画像中で m_r 行 m_c 列に位置するマクロブロック $MBK(m_r, m_c, k)$ ($m_r=1, \dots, M_r, m_c=1, \dots, M_c$) が何のオブジェクトであることを示す確信度ベクトル $C(m_r, m_c, k)$ はシーンの状況に関する特徴要因と画像信号に関する特徴要因を説明変数として推定できる。即ち、

$$C(m_r, m_c, k) = W_S S_S(m_r, m_c, k) + W_P S_P(m_r, m_c, k) \quad (3.13)$$

と表現できる。

ただし、 $S_S(m_r, m_c, k)$ は $MBK(m_r, m_c, k)$ に投影されるシーンの状況を表す L_S 次元列ベクトルであり、 W_S は S_S に対する回帰係数の組を表す N 行 L_S 列の行列である。

また、 $S_P(m_r, m_c, k)$ は $MBK(m_r, m_c, k)$ に含まれる画像信号の特徴を表す L_P 次元列ベクトルであり、 W_P は S_P に対する回帰係数の組を表す N 行 L_P 列の行列である。

本稿ではこの S_S を状況ベクトルと呼ぶ。いま、状況ベクトルが下記のように構成されるとする。

$$S_S(m_r, m_c, k) = [S_{nenv}^T \ S_{location}^T \ S_{time}^T]^T \quad (3.14)$$

ここで、 $S_{nenv}(m_r, m_c, k)$ は $MBK(m_r, m_c, k)$ に関するシーンの自然環境情報を表す L_{nenv} 次元列ベクトルである。また、 $S_{location}(m_r, m_c, k)$ はシーン中の位置情報を表す $L_{location}$ 次元列ベクトルであり、 $S_{time}(m_r, m_c, k)$ はシーンの時間情報を表す L_{time} 次元列ベクトルである。

状況ベクトルでは一般に代数演算できない特徴量を表現することが求められるため、各状況を言語ラベル（インデックス）で表現して数量化理論 I 類を適用する。例えば、時間帯や天候は次のように表される。

$$S_{time} = (\text{早朝, 朝, } \dots, \text{夕方, 夜})^T \quad (3.15)$$

$$S_{nenv} = (\text{晴, 曇, 雨, } \dots, \text{暴風雨})^T \quad (3.16)$$

ただし、インデックス {早朝, ..., 晴, ...} は対応する事象が真のとき 1, 偽のとき 0 を設定する。

3.7 計測行列の作成

いま、特徴要因行列 S と確信度行列 Y を用いて、

$$\Psi = [S \ Y] \quad (3.17)$$

で表される行列を計測行列と呼ぶことにする。ただし、

$$S = [s(1) \dots s(k) \dots s(K)]^T \quad (3.18)$$

であり、 $s(k)$ は時刻 T_k ($k=1, \dots, K$) における特徴要因ベクトル s の値とする。また、 Y は

式(3.6)と同様に確信度の履歴を表す行列である．ここで， Y に含まれる確信度はサンプル映像中の各ブロック画素に投影されるオブジェクトが何であることを示す正解（Ground Truth）を教示することで獲得される．この正解値は通常、人間の目視判断や高信頼度のセンシング装置で得た判定値、あるいは過去に蓄積された統計的データによって与えられる．このような正解値は、その判定結果であるオブジェクトインデックスに対応する成分の値を1とし、それ以外は0として設定する．この計測行列から式(3.3)の回帰係数行列を算出すれば、新たな入力映像中のブロック画素に対する確信度ベクトルを推定することができる．いま、

$$\mathbf{s} = (\mathbf{S}_P^T \mathbf{S}_S^T)^T \quad (3.19)$$

と定義する．ただし、

$$\mathbf{S}_P = (\mathbf{S}_{color}^T \mathbf{S}_{texture}^T \mathbf{S}_{block_position}^T \mathbf{S}_{motion}^T)^T \quad (3.20)$$

$$\mathbf{S}_{color} = (I_Y, I_U, I_V)^T \quad (3.21)$$

$$\mathbf{S}_{texture} = (A_{hor}, A_{ver}, A_{all})^T \quad (3.22)$$

$$\mathbf{S}_{block_position} = (m_c, m_r)^T \quad (3.23)$$

$$\mathbf{S}_{motion} = (v_x, v_y)^T \quad (3.24)$$

であり、 \mathbf{S}_{color} はMBK(m_r, m_c, k)に関するシーンの色情報を表す L_{color} 次元列ベクトルである．式(3.21)はブロック内平均色であり、2次元DCT係数のDC成分に対応する． $\mathbf{S}_{texture}$ はBLKの2次元DCT係数のAC成分から算出する特徴ベクトルであり、 A_{hor} は水平方向の画素値変化に対応したAC成分の絶対値和、 A_{ver} は垂直方向の画素値変化に対応したAC成分の絶対値和、 A_{all} はブロック全体にわたる画素値変化に対応した成分の絶対値和、即ちAC電力を表す．また、 \mathbf{S}_{motion} はMBK(m_r, m_c, k)における2次元の動きベクトル（画素単位）を表す．

3.8 認識率

上記で得られた確信度ベクトルの最大成分を求めると確信度が最大であるオブジェクトのインデックス $X(m_r, m_c, k)$ を決定できる．これを予め求めておいた正解 $G_{view}(m_r, m_c, k)$ と比較して一致するならMBK(m_r, m_c, k)の認識結果 $R_{view}(m_r, m_c, k)$ を1、一致しないならば0とする．ただし、 $G_{view}(m_r, m_c, k)$ は、MBK(m_r, m_c, k)のGround Truthが視点viewの対象とするオブジェクトインデックス群（以下、オブジェクトリスト） O_{view} に含まれるオブジェクトインデックスのどれかに等しいとき1を取る．それ以外は認識対象としないため0とする．

これに基づき、ある応用上の視点viewに基づいて算出する認識率 Q_{view} (%)を次式で定義することができる．

$$Q_{view} = \frac{100}{K_{view}} \sum_{i=1}^{M_r} \sum_{j=1}^{M_c} R_{view}(i, j, k) \quad (3.25)$$

ただし、

$$K_{view} = \sum_{i=1}^{M_r} \sum_{j=1}^{M_c} G_{view}(i, j, k) \quad (3.26)$$

である。

たとえば安全性という視点 *safe* から見た認識率 Q_{safe} はオブジェクトリスト $O_{safe} = \{\text{先行車, 対向車, 信号, 標識}\}$ に基づいて計算することができる。同様に利便性という観点 *convenience* から見た認識率 $Q_{convenience}$ はオブジェクトリスト $O_{conv} = \{\text{ビル, 緑地, 空, 先行車, 対向車, 信号}\}$ に基づいて計算できる。

一方、観光という視点 *sightseeing* から見た認識率 $Q_{sightseeing}$ はオブジェクトリスト $O_{sightseeing} = \{\text{空, 緑地, ビル, 海, 砂浜}\}$ に基づいて計算できる。ただし、*view* やインデックスの語彙はアプリケーションやユーザプロファイルに依存して決まる。

式(3.25)は画面全体にわたってある視点 *view* に対応するオブジェクトリスト中に含まれるすべてのオブジェクトについての認識率を評価する尺度であるため、シーン認識率と呼ぶことにする。

これに対し、個々のオブジェクトに対する認識性能を評価する際には、情報検索効率の評価で用いる F 尺度（適合率と再現率の調和平均）をブロック画素単位で算出し、その値をオブジェクト認識率と呼ぶことにする [A12]。

3.9 推論特性の合成

多変量線形回帰分析では教示データの追加更新に際して、回帰係数行列を逐次更新式で算出できるという計算上の利点を有する。これは少量の教示データが時々刻々と追加される際の応用に適している。加えて、1台あるいは複数の車両の車載カメラ映像で得た複数の教示データや推論特性を適宜合成してよりよい推論特性を得ることができれば、冒頭で紹介した数々の応用の実現に大きく貢献できる。

この際に教示データの段階で合成して再度学習させることがもっとも平易であるが、一方で行列の計算量が大きくなりすぎるという問題が発生する。そこで、過去の各学習で既に獲得した特性を再利用することを考えた。これは多変量線形回帰分析の線形性に着目して下記のように導出される。

3.9.1 回帰係数行列の合成

ある教示データセット T が確定していれば、式(3.10)に基づいて m 番目の要因に関する回帰係数行列 B_m が計算できる。いま、議論を簡単にするために1つの要因に固定する。そして添え字 m を要因番号から教示データセットの番号 h に読み替えると、

$$B_h = (S_h^T S_h)^{-1} S_h^T Y_h = (A_h)^{-1} S_h^T Y_h \quad (3.27)$$

ただし、

$$A_h = S_h^T S_h \quad (h=1, \dots, H) \quad (3.28)$$

である。このとき、 H 個のデータセット全体を用いて計算される回帰係数行列 $B_{1:H}$ は次式で表される。

$$B_{1:H} = \left\{ \sum_{h=1}^H A_h \right\}^{-1} \left\{ \sum_{h=1}^H A_h B_h \right\} \quad (3.29)$$

したがって、個々のデータセットに関して共分散行列 A_h と回帰係数行列 B_h が各々の学習段階で個別に計算されていれば、それらを用いて上式で $B_{1:H}$ を合成できる。

3.9.2 計算量の評価

いま、式(3.29)の計算量を Q_1 、式(3.10)の計算量を Q_2 とすると、次式で表現される。

$$Q_g = Q_{g_in} + Q_{g_mul} + Q_{g_add}, \quad (g=1, 2) \quad (3.30)$$

ただし、 Q_{g_inv} は次数 $L \times L$ の逆行列計算、 Q_{g_mul} は乗算数、 Q_{g_add} は加算数である。 L 次元正方行列の逆行列計算は式(3.29)、式(3.10)ともに 1 個を要するので、 $Q_{1_inv} = Q_{2_inv}$ となる。そこで、 Q_{g_mul} に着目し、乗算回数 の比 α について下限値を評価すると、

$$\alpha = \frac{Q_{2_mul}}{Q_{1_mul}} > \frac{2K_0}{N} \quad (3.31)$$

となる。ただし、 N は語彙数である。また、簡単化のため、 K_0 を 1 個の教示データに含まれる教示イベント数の平均値と考え、 $K = HK_0$ とした。一例として $N=64$ 、 $K_0=1000$ とすると乗算個数は $1/30$ 以下にできることがわかる。さらに、式(3.29)の方法は計算に必要なデータのサイズが K_0 に影響されないため、データ通信や記憶において式(3.10)よりも有利である。

3.10 主成分分析による多変量回帰分析

まず、特徴ベクトルを MPEG 形式における M_{walk} 個の連続する予測フレームから生成する。ここで、 n 番目の特徴ベクトル ($n=1, \dots, N$) は次式で定義される。

$$\mathbf{x}_n = (x_{n1}, \dots, x_{nP})^T \quad (3.32)$$

ここで、 P は特徴ベクトルの次元数である。7 章の実験では仮に M_{walk} を 5 に設定したので、1 秒に 10 フレームの予測フレームがある場合、時系列の区間は 0.5 秒である。したがって N は $64 \times 5 = 320$ となる。

対応する計測行列は次式で表現される。

$$X = (\mathbf{x}_1 \cdots \mathbf{x}_N)^T \quad (3.33)$$

この行列は、人間が映っている画像の特性を効果的に学習するために、歩行者シーンから抽出したデータをもとにして定義する。したがって、共分散行列 S は次式で定義される：

$$S = \frac{1}{N} X^T X \quad (3.34)$$

S の p 番目の固有ベクトル ($p=1, \dots, P$) は次式で記述される。

$$\mathbf{w}_p = (w_{p1}, \dots, w_{pP})^T \quad (3.35)$$

主成分得点は次式を用いて計算できる：

$$\mathbf{f}_p = X \mathbf{w}_p \quad (3.36)$$

これは容易に正規化行列形式に拡張され、次式で表される：

$$F = (\mathbf{f}_1 \cdots \mathbf{f}_P) = X W \text{diag}(1/\sqrt{\lambda_1} \cdots 1/\sqrt{\lambda_P}) = X W \Delta_S^{-1/2} \quad (3.37)$$

ただし、 $W = (\mathbf{w}_1, \dots, \mathbf{w}_P)$ であり、 λ_p は S の p 番目の固有値を表す。すなわち $\Delta_S \equiv \text{diag}(\lambda_1, \dots, \lambda_P)$ である。

あるシーン中でカテゴリ 'cat' に対応する時区間について PCA を施し、固有値行列 A_{cat} が獲得される場合を考える。いま、人間による教示でシーン中の歩行者 (カテゴリ名 'ped' で表す) を規定すれば、歩行者がシーン中で映っている時区間における主成分得点 (以後、PCS とする) は次式で算出できる：

$$\varphi_{ped} = A_{ped} \mathbf{x} \quad (3.38)$$

同様に、「運転 (車載カメラが搭載された車が動いている) 中で歩行者がカメラ視野内に存在しない」時区間のカテゴリを 'drive' と定義すれば、'drive' に対応する PCS は次

式で表される :

$$\boldsymbol{\varphi}_{drive} = \mathbf{A}_{drive} \mathbf{x} \quad (3.39)$$

このように、PCA の学習段階で学習データをカテゴリごとに集めることにより、そのカテゴリの事象を認識する際の判断基準となる PCS を算出する固有値行列を獲得できる。

3.11 カーネル主成分分析

KPCA (Kernel PCA) は信号処理分野において、従来の PCA よりも高度な品質を達成するために導入されている。KPCA は入力信号の原空間を高次元空間に射影することで線形識別では困難であった識別を容易にする。線形主成分分析 (LPCA) と比較すると、射影後の空間は主軸が湾曲している非線形多様体と考えることができる。以下、文献[G3]を参考にして KPCA の学習フェーズと認識フェーズの算出方法を示す。

3.11.1 学習フェーズ

原空間のベクトル \mathbf{x} で構成される学習データを学習する際に、 \mathbf{x} を高次元空間へ写像する非線形関数を $\boldsymbol{\phi}$ とすれば、カーネル行列 K の i 行 j 列成分は次式で表される。

$$K_{ij} = \boldsymbol{\Phi}(\mathbf{x}_i) \cdot \boldsymbol{\Phi}(\mathbf{x}_j) \quad (3.40)$$

これを正規化したカーネル行列は次式で与えられる。

$$\begin{aligned} \bar{K}_{ij} &= \bar{\boldsymbol{\Phi}}(\mathbf{x}_i) \cdot \bar{\boldsymbol{\Phi}}(\mathbf{x}_j) \\ \bar{K} &= K - I'_N K - K I_N + I'_N K I_N \end{aligned} \quad (3.41)$$

以下、簡単に証明しておく。

$$\begin{aligned} \bar{K}_{ij} &= \bar{\boldsymbol{\Phi}}(\mathbf{x}_i) \cdot \bar{\boldsymbol{\Phi}}(\mathbf{x}_j) = \left\{ \boldsymbol{\Phi}(\mathbf{x}_i) - \frac{1}{N} \sum_{m=1}^N \boldsymbol{\Phi}(\mathbf{x}_m) \right\} \cdot \left\{ \boldsymbol{\Phi}(\mathbf{x}_j) - \frac{1}{N} \sum_{n=1}^N \boldsymbol{\Phi}(\mathbf{x}_n) \right\} \\ &= \boldsymbol{\Phi}(\mathbf{x}_i) \cdot \boldsymbol{\Phi}(\mathbf{x}_j) - \left\{ \frac{1}{N} \sum_{m=1}^N \boldsymbol{\Phi}(\mathbf{x}_m) \right\} \cdot \bar{\boldsymbol{\Phi}}(\mathbf{x}_j) - \bar{\boldsymbol{\Phi}}(\mathbf{x}_i) \cdot \left\{ \frac{1}{N} \sum_{n=1}^N \boldsymbol{\Phi}(\mathbf{x}_n) \right\} + \frac{1}{N^2} \left\{ \sum_{m=1}^N \boldsymbol{\Phi}(\mathbf{x}_m) \right\} \cdot \left\{ \sum_{n=1}^N \boldsymbol{\Phi}(\mathbf{x}_n) \right\} \\ &= K_{ij} - \frac{1}{N} \sum_{m=1}^N \{ \boldsymbol{\Phi}(\mathbf{x}_m) \boldsymbol{\Phi}(\mathbf{x}_j) \} - \frac{1}{N} \sum_{n=1}^N \{ \boldsymbol{\Phi}(\mathbf{x}_i) \boldsymbol{\Phi}(\mathbf{x}_n) \} + \frac{1}{N^2} \sum_{m=1}^N \left\{ \boldsymbol{\Phi}(\mathbf{x}_m) \sum_{n=1}^N \boldsymbol{\Phi}(\mathbf{x}_n) \right\} \\ &= K_{ij} - \frac{1}{N} \sum_{m=1}^N \{ \boldsymbol{\Phi}(\mathbf{x}_m) \boldsymbol{\Phi}(\mathbf{x}_j) \} - \frac{1}{N} \sum_{n=1}^N \{ \boldsymbol{\Phi}(\mathbf{x}_i) \boldsymbol{\Phi}(\mathbf{x}_n) \} + \frac{1}{N^2} \sum_{m=1}^N \left\{ \sum_{n=1}^N \boldsymbol{\Phi}(\mathbf{x}_m) \boldsymbol{\Phi}(\mathbf{x}_n) \right\} \end{aligned} \quad (3.42)$$

この正規化カーネル行列を用いて固有ベクトルとその正規化の計算は以下に定義される。

$$\begin{aligned} \hat{\lambda} \boldsymbol{\alpha} &= \bar{K} \boldsymbol{\alpha} \\ \hat{\boldsymbol{\alpha}}^{(l)} &= \frac{\boldsymbol{\alpha}^{(l)}}{\sqrt{\hat{\lambda}_l}} \quad l = 1, \dots, N \end{aligned} \quad (3.43)$$

3.11.2 認識フェーズ

認識フェーズでは新しい入力ベクトル \mathbf{y} と学習ベクトル \mathbf{x} の間で次式によりカーネル行列を計算する。

$$K_{lj}^{in} = \boldsymbol{\Phi}(\mathbf{y}_l) \cdot \boldsymbol{\Phi}(\mathbf{x}_j) \quad (3.44)$$

式(3.41)と同様に正規化カーネル行列について次式の関係が成立する。

$$\bar{K}^{in} = K^{in} - I'_N K - K^{in} I_N + I'_N K I_N \quad (3.45)$$

高次元空間での主成分を抽出するために固有ベクトル \mathbf{v} 上への写像を考えると、

$$\begin{aligned}
(\mathbf{v}^{(l)} \cdot \bar{\Phi}(\mathbf{y})) &= \sum_{i=1}^N \hat{\alpha}_i^{(l)} (\bar{\Phi}(\mathbf{x}_i) \cdot \bar{\Phi}(\mathbf{y})) \\
&= \sum_{i=1}^N \hat{\alpha}_i^{(l)} \bar{K}^{in}(\mathbf{x}_i \cdot \mathbf{y})
\end{aligned}
\tag{3.46}$$

となるので、式(3.45)により高次元空間での主軸 $\mathbf{v}^{(l)}$ に対する射影値（主成分得点）を各 l について算出できる。この N 個の主成分得点をもとに入力ベクトルと学習ベクトルの間の類似性（例えば学習データが歩行者であれば歩行者らしさ）に関する評価値を算出することができる。

3.12 シーンの構造モデル

一般に、走行車両は局所的には前方に直進することが多いため、車載カメラによる道路画像では図 3.2 のような箱型の立体構造モデルを仮定することが考えられる [A8][A9]。

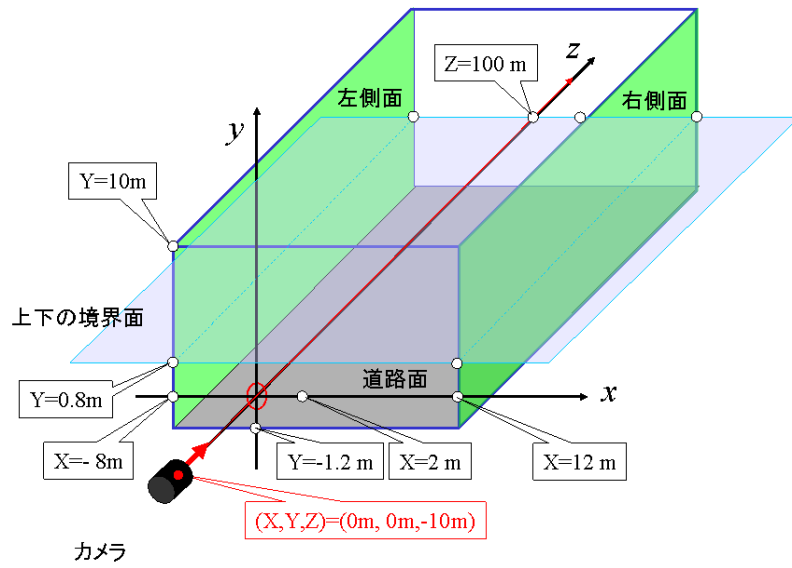


図 3.2 道路シーンの立体構造モデル

Fig. 3.2 Solid structure model of road scene.

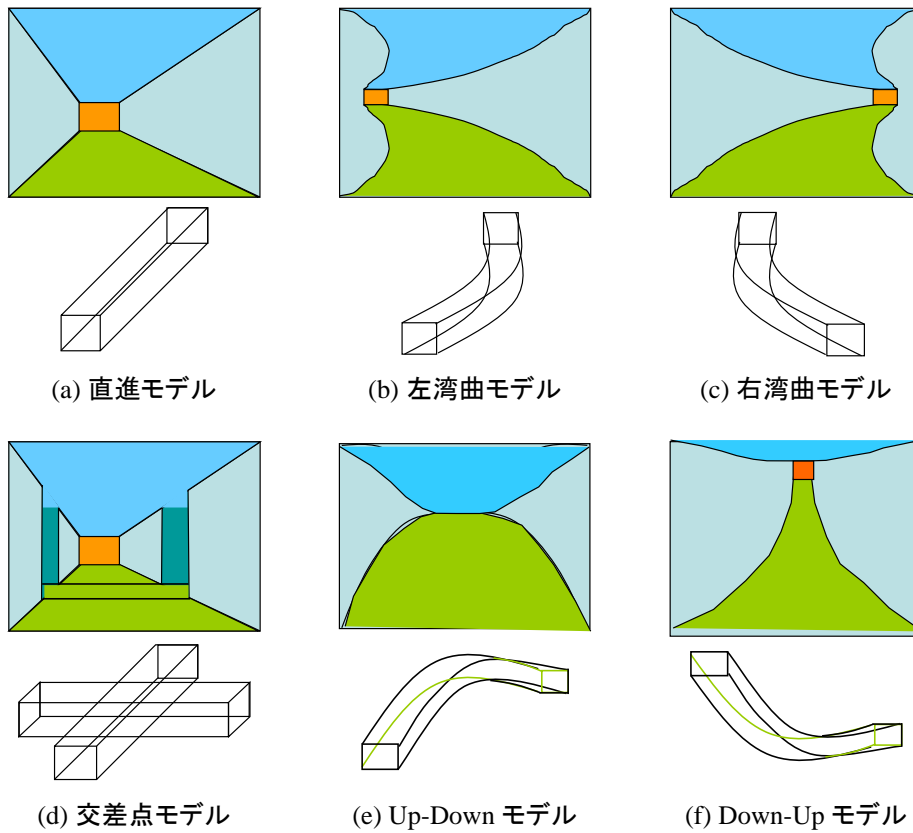


図 3.3 道路シーンの立体構造モデル

Fig. 3.3 Solid structure model of road scene.

3.13 シーンクラス

3.13.1 主観的シーンクラス

主観的シーンクラス(SSC)は人間の定義したシーンの主観的なグルーピングによって決められる。それは画像に含まれる地理的な条件やオブジェクトを用いることで行われる。最初に、筆者は以下の統計的要因に注目することでシーンクラスを考えた。

- (1) シーン中の主要なオブジェクトの配置 (ビル、樹木、緑地、高架、雪、他)
- (2) 道路クラス (一般道、高速道路、他)
- (3) 地理的分類 (山、市街地、海岸、湖畔、他)

これらの要因の経験的考察によって、図 3.4 に示すようないくつかのクラスが考えられる。それらの各々を自然言語表現することは比較的単純である。実際のシーンでは車両や歩行者といった動的な物体が付加的に存在する。さらに、天候や照明の変化、カメラの位置や姿勢の違い、車両からの見え方 (aspect) といった他の要因群も存在する。

3.13.2 計算的シーンクラス

計算的シーンクラス (CSC: Computational Scene Class) は複数のシーンを純粋に機械学習に従って分類しようとするものである。今、

「もし S_1 がシステムに学習されているならば、システムは S_2 を認識できる」

という概念を有向枝で定義できる。これは、より具体的に表現すると、例えば、

「シーン S_1 をシステムの学習データとするとき、システムは S_2 を事前に設定した閾値 α よりも高い認識率で認識する」

と書くことができる。

S_2 から S_1 に向けての有向枝と同時に、もし同様な枝が S_2 から S_1 に向けて定義できるならば、その関係は

「 S_1 と S_2 が同じ計算的シーンクラスに属する」

と表現される。ただし、これはシーン S_1 が一意的な CSC を持つことを意味しているわけではない。すなわち、複数の CSC を S_1 について定義することが可能である。他方、ある S_1 に対して CSC と同時に存在する 1 つ以上の主観的シーンクラス (SSC) が分かっているならば、SSC と CSC の関係は自動的に構築できる。

3.14 意味との相互作用

ここでは計測と意味の両面から自動インデキシングの統一的手法を提案する。シーンの 3次元復元においてよく知られた因子分解法は特異値分解 (SVD) から導出され、次式で表現される。

$$\tilde{W} = MX \quad (3.47)$$

ここで、 \tilde{W} は計測行列 ($2F \times P$) を表し、 F 個のカメラから撮影した対象物体の P 個の 3D サンプル点について平均値分離した 2次元位置を要素として持つ。 M ($2F \times 3$) は正射影の仮定のもとでカメラ特性を表す。そして X ($3 \times P$) は対象物の 3次元形状を表す。

次に、意味の獲得は多変量線形回帰モデルを用いて次式で表現される。

$$C = RS \quad (3.48)$$

ここで、 C は確信度行列 ($N \times P$) であり、 N はオブジェクトカテゴリの数を意味する。いくつかの応用では、各カテゴリは自然言語の語彙で表現できる。また、それは確信度ベクトルの意味的なマッピングの特質により、PCA における主成分に関係づけることがで

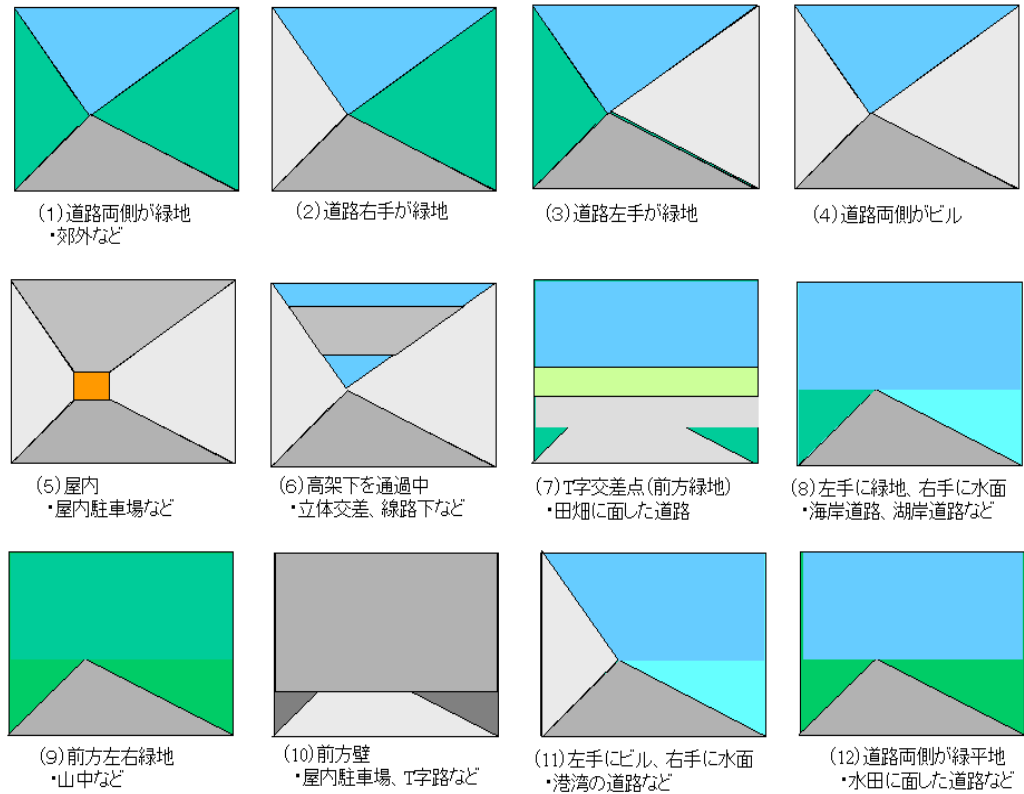


図 3.4 道路シーンの立体構造モデル

Fig. 3.4 Solid structure model of road scene.

きる。それは非線形の場合（カーネル PCA）であっても同様である。 \mathbf{R} ($N \times L$) は回帰係数行列を表す。 \mathbf{S} ($L \times P$) は観測した状態行列であり、 L はターゲットオブジェクト上の各サンプル点における観測状態の数である。基本的にはこの状態はブロック画素に投影された画像情報を通じて獲得される。しかし、容易に周囲の環境や信号やセマンティクスから得られる関連属性に拡張できる。同様に、その拡張により、ユーザプロファイルを含めることも可能である。それは、時間変化特性を有するプロファイルに対しても可能である。分解と同じ形式をとるために、方程式を以下のように変換できる：

$$\mathbf{S} = \mathbf{R}\mathbf{C} = \mathbf{Q}\mathbf{C} \quad (3.49)$$

ここで、 \mathbf{Q} ($L \times N$) と \mathbf{R} は \mathbf{R} の一般化疑似逆行列を表す。したがって、 \mathbf{S} の学習データと \mathbf{C} の計測データから、次式の最小 2 乗法により \mathbf{Q} を推定することができる：

$$\mathbf{Q} = \mathbf{S}\mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1} \quad (3.50)$$

式(3.49)はカメラ番号 f ($f=1, \dots, F$) とは独立であることに注意されたい。したがって、2 つの式は次のように統合できる：

$$\bar{\mathbf{W}} \equiv \begin{bmatrix} \tilde{\mathbf{W}} \\ \tilde{\mathbf{S}} \end{bmatrix}, \quad \bar{\mathbf{W}} = \mathbf{M}\bar{\mathbf{X}} = \begin{bmatrix} \mathbf{M} & \mathbf{A} \\ \mathbf{B} & \tilde{\mathbf{Q}} \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \mathbf{C} \end{bmatrix} \quad (3.51)$$

ただし、 $\bar{\mathbf{W}}$ ($(2+L)F \times P$) は拡張した計測行列を示す。そして、当面は $\mathbf{A}=\mathbf{0}$ かつ $\mathbf{B}=\mathbf{0}$ と仮定する。 $\tilde{\mathbf{S}}$ ($LF \times P$) と $\tilde{\mathbf{Q}}$ ($LF \times N$) はそれぞれ F \mathbf{S} s と F \mathbf{Q} s, から容易に作成できる。この方式によりシステムは、膨大な 2 次元画像メディアデータベースから、対象物とそれを含むシーン全体に関して、その 3 次元形状と動きのみならずその意味までも同時に獲得

することが可能になる。

しかしながら、式(3.51)にかかわる基本的問題が残されている。一つ目はSVDと多変量回帰 (MVR) の違いである。SVDは形状とカメラ運動を同時に導出し、その解は一意的と考えられる。他方、MVRの関係は基本的には射影でありQとCについては複数の解が存在する。二つめの違いはSVDが学習プロセスを要さず、MVRは学習が必要である点である。もし式(3.51)のMVR部分のQがSVD部分のM (光学的カメラ特性) と同様の確固たる性質を持っていれば、QとCについてもSVDで同時に獲得できる。

第2のアイデアは式(3.51)を単純に線形射影問題として考えることである。この場合、たとえカメラ特性Mが透視変換特性を持っているとしても、最小2乗法を使って説くことができる。Mが透視変換特性の場合はSVDでは解けない。これは因子分解法が一般的に正射影を仮定しているためである。このアイデアは式(3.51)が統計的学習問題として考えられているということである。しかしながら、それはカメラ番号とカメラ配置の自由度という問題を引き起こすかもしれない。もしF個のカメラが固定されているかあるいはかなり制約のある軌道の上に位置しているとき、この統計的学習は可能になる。

知識処理の観点でいうと、式(3.51)における行列Aと行列Bは特殊な役割を持っている。AはCからWへの写像を意味する。すなわち、サンプル点の2次元位置情報は各点が所属するオブジェクトカテゴリの確信度に影響されるかもしれない。第2に、BはXからSへの写像を意味しており、これによって、形状情報が観測状態に影響を与える可能性があることを示している。

さらに、これらのAとBは時空間的な発展過程を持つと同時に、人間のフィーリングや情動にかかわるセンシング特性に関するいくつかの数量的表現への非線形写像を有するであろう。この期待は式(3.51)における基本的な数学的メカニズムが線形システムを志

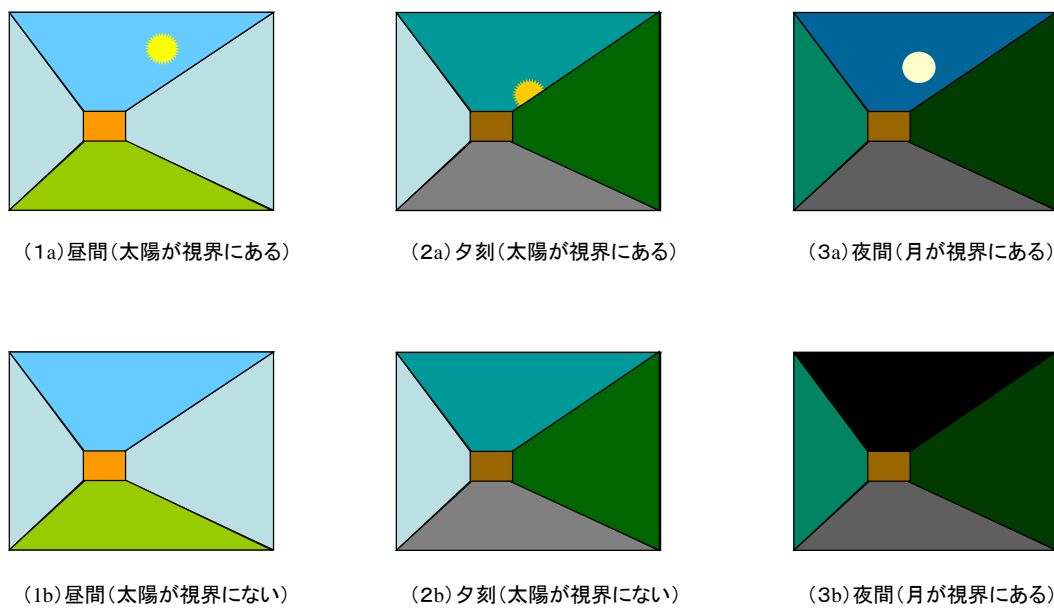


図 3.5 道路シーンの立体構造モデル

Fig. 3.5 Solid structure model of road scene.

向することを意味する。そしてシステムにおける他の非線形機能は主として非線形マッピング、論理操作、および状態遷移構造のような巡回アーキテクチャで実現される。

観点のカテゴリでいえば、式(3.51)はさらに拡張できる。拡張した観点が U (user), であるとする、上記ですでに議論した2つの観点は X (3D シーン構造)と C (コンテンツの意味)である。これより式(3.51)は次式のように拡張できる：

$$W^* = M^* X^* \quad (3.52)$$

ただし、

$$W^* \equiv \begin{bmatrix} W \\ S \\ P \end{bmatrix}, \quad M^* \equiv \begin{bmatrix} M_X & A_C & A_U \\ B_X & M_C & B_U \\ C_X & C_C & M_U \end{bmatrix}, \quad X^* \equiv \begin{bmatrix} X \\ C \\ U \end{bmatrix} \quad (3.53)$$

である。計測行列 W^* における P は、実際上はたとえば、システムのシーン理解に影響を与えるユーザ属性を数量化理論で表現したユーザプロフィールに相当する。

3.15 アプリケーションと語彙

カメラが車載カメラのように移動するときは、特にそれが前に進み通り抜けるとき、ラベリング問題は突然難しくなる。スナップショットはラベル付けすることができ、シーンクラス (SC) を生成することはできるが、ビデオシーンは絶えず変化し新しい SC を作成する。この結果、たとえ観測したビデオシーンがトータルとして一つのラベルしか持たない場合でも、それは局所的には多くの異なるラベルを含むであろう。

語彙の定義はきわめてターゲットアプリに依存する。シーンクラス獲得で考えた目的は検索、分類、エンタテインメント、ロケーションウェアサービス、環境認識など多岐にわたる。あるビデオシーケンスについて SC の時間スパンを考えねばならない。これについての簡単なアイデアはあるシーンについての SC ラベルを SC ラベルの時系列の作成を通じて決定していくことである。ここで、シーン中の各単独のフレームについての SC ラベルの各々は次のように定義できる：

$$\alpha_k = \text{Index}\{p(X_k^* | W_{1:k}^*, X_{1:k-1}^*)\} \quad (3.54)$$

$$\alpha = \text{SceneClass}\{\alpha_{1:K}\}, \quad \alpha_{1:K} = \{\alpha_1, \dots, \alpha_K\}$$

ただし、 K はあるシーンに含まれる全フレーム数であり、 $p()$ は確率分布を意味する。 Index はオブジェクトベースの確信度からフレーム時刻 k における SC ラベルへの写像である。 SceneClass は一連のインデックスからシーン全体に割り当てた単一のラベルへの写像である。この式の背後にある基本的なアイデアはパーティクルフィルタリングベースの状態トラッキング [E-7] である。シーンデータベースの定義と集め方はターゲットアプリによって定められ、語彙と調整に用いる論理操作が(3.54)式の2つの非線形写像に含まれる。さらに、この CSC の自動生成機構は非線形写像により時系列表現された特徴のラベリングを可能にするが、カーネル主成分分析 (KPCA) や多様体学習 [E-12] に関する研究とも関連が深い。

2つ目の問題は視点 (viewpoint) の影響である。これはターゲットアプリの制約のもとの評価側の変化である。後述するように、CSC の有向グラフ表現は視点の影響を受ける。自然言語理解の類推に基づき、単語間の意味ネットワークの有向グラフ構造は、視覚的概念形成に向けた出発点のひとつとして CSC 生成に適用できる。

3.16 複数手法の統合化

3.16.1 3S 認識モデル

この考えは形状、スペクトル、状況という3つの大きなトリガーで有機的に活性化される。まず1番目の形状ベースのオブジェクト認識手法は既に数多く存在し、通常はこれをもっともポピュラーなアプローチである。しかしそれらはしばしば、ぼけや隠れ、動的な変化、あいまいさ、個々の差異といった深刻な問題を呈する。

2番目にスペクトルベースの認識を考える。これは空間スペクトルと時間スペクトルの両方に適用される。このスペクトルベースのアプローチは形状ベースの手法が苦しんできた問題に対してロバストであるが、認識タスク全体を単独でカバーできない。

3番目の状況は、対象領域を信号処理によりパターン分類した中間結果を、知識ベースで判定する際に必要となる。この状況要因は観測者とアプリケーションに依存するが、パーティクルフィルタのような仮説の生成棄却を伴うトラッキングにおいては良好な仮定を提供する。しかし、その際に対象領域の形状とスペクトルが十分に解析されていなければ、必要以上に多くの確率的解釈を生成し、最終解の判定が困難になる恐れがある。

そこで、筆者はこれら3つのトリガーの動的な活性化手法（以下、3Sモデル）を提案する。これによってシステムはメインの戦略を随時スイッチングできる。そして確信度は一連のサブタスクを統合し、3Sモデルにおけるその動的スイッチング挙動を通じて認識に向かう。

3.17 映像内容の予測と推定

これから撮影しようとするシーンのシーンクラスがシステムに与えられている場合、

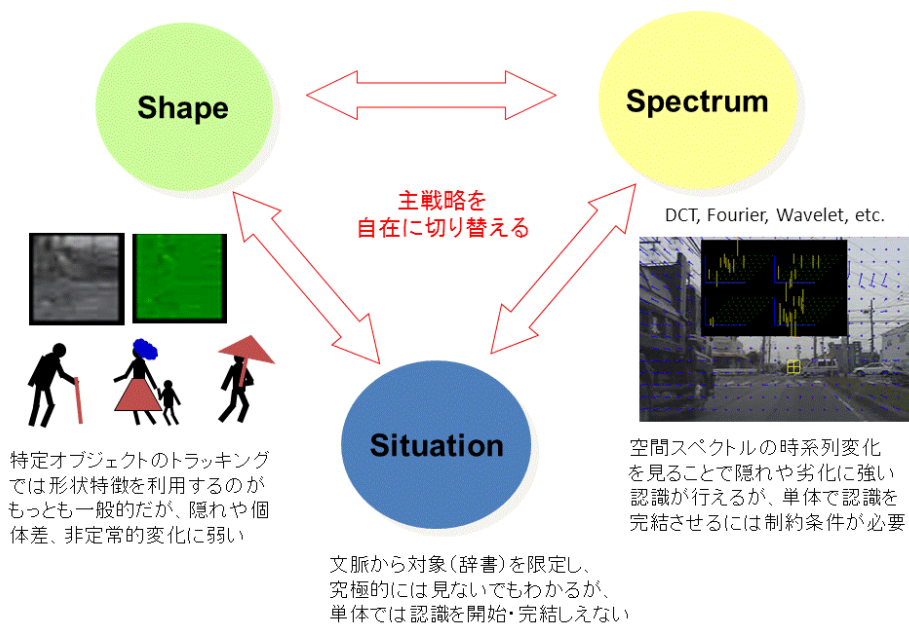


図 3.6 3Sモデルの概念

Fig. 3.6 Concept of 3S model.

システムはシーン内のオブジェクトやシーン構造の認識を行うための語彙と認識特性を

事前に用意することができ、解空間の不要な探索はなくなる。これは近年 GPS 機能を有するカメラが増大し、撮影した内容に位置情報を付与できることからきわめて現実的なことである。また、異なる撮影者が Web 上で同一の場所のシーンをアップロードして共有することにより、さまざまな時間条件のもとで過去のシーンが蓄積され、スタティックな事象（建物や地形の形状など）あるいは周期的変化を伴う現象（波が打ち寄せるシーン、木々が揺れるシーン、天候の変化に伴うシーン、四季折々の変化、人の流れなど）については十分予測可能になってきている。その蓄積される画像の付与される情報には次のような場合が考えられる：

- (1) 場所の情報が与えられている場合
 - a) カメラの位置姿勢およびその時間変化（軌道と光軸がわかる場合
 - b) 視野内のオブジェクトやシーン構造が既知または推定できる場合
 - c) web や地図などにある説明情報（description）が使える場合
- (2) 時間はわかるが場所の情報が不明の場合
- (3) 時間も場所も不明の場合

他にも、カメラの各種設定パラメータ（ズームや照度など）が影響することが考えられる。映像の予測と推定にはこのような付帯記述や検索インデックスを有する過去のデータベースが極めて重要と考えられるが、必ずしもそれらは現段階で整備されているわけではない。そこでまず、未整備の過去データを用いる統計的手法を考察する。

いま、次のような時系列モデルを仮定する：

$$\mathbf{y}_n = \mathbf{t}_n + \mathbf{s}_n + \mathbf{p}_n + \mathbf{d}_n + \mathbf{c}_n + \mathbf{w}_n \quad (3.55)$$

ここで、 \mathbf{y}_n は K 次元列ベクトルで、 k 番目の成分は時刻 n に k 番目の状態節点（particle）（ $k = 1, \dots, K$ ）で観測される歩行者の数であり、その状態とは地理的スポットとオブジェクトの組み合わせとして定義される。同様に、 $\mathbf{t}_n, \mathbf{s}_n, \mathbf{p}_n, \mathbf{y}_n, \mathbf{c}_n, \mathbf{w}_n$ はそれぞれトレンド成分、季節成分、定常 AR（自己回帰）成分、曜日成分、時刻の成分、そして白色雑音に対応する K -次元ベクトルである。もし、マクロスケールの状態ベクトル \mathbf{x}_n を図 3.7 のように定義すれば、 \mathbf{x}_n は次式で表現される。

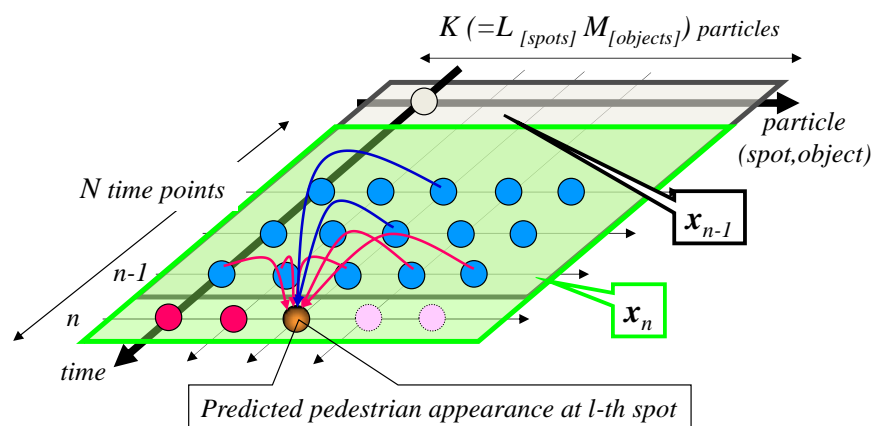


図 3.7 歩行者出現予測のグラフィカルモデル

Fig. 3.7 Graphical model for prediction of pedestrian appearance..

$$\mathbf{x}_n = (y_{n,1}, \dots, y_{n,K}, \dots, y_{n-(N-1),1}, \dots, y_{n-(N-1),K})^T \quad (3.56)$$

したがって、式(3.55)はよく知られた状態空間モデルに変換され、次式で表現できる：

$$\begin{cases} \mathbf{x}_n = F\mathbf{x}_{n-1} + G\mathbf{v}_n \\ \mathbf{y}_n = H\mathbf{x}_n + Q\mathbf{w}_n \end{cases} \quad (3.57)$$

ただし F, G, H , および Q はそれぞれ、 $KN \times KN$ 行列、 $KN \times M_v$ 行列、 $K \times KN$ 行列、 $K \times M_w$ 行列である。 \mathbf{v}_n と \mathbf{w}_n はそれぞれ、 M_v -次元のシステム雑音、 M_w -次元の観測雑音である。式(3.55)の右辺における各項を抽出するためには \mathbf{x}_n に関していくつかの仮定を追加する必要がある。また、各係数行列の推定に際しては、歩行者数に対応する莫大な量の正解値を用意する必要がある。

3.18 予測動作における状況ベクトル

上記の特徴ベクトル（目的変数）が発生する際に、その説明変数となる状況ベクトルは表 3.2 の {場所, 日時, 天候} で構成される。ここで、場所は 3 次元（各次元が区間に相当）、日時は {年, 月, 曜日, 時間帯} の合計 25 次元、天候は {晴, 曇, 雨} の 3 次元とした。

状況ベクトルは更に拡張でき、日種効果, 季節効果, 時間帯効果, 年中行事の影響なども考慮することができる。また、道路環境や地理要因も同様に記述できる。本稿では以下の 3 つのベクトルで構成する。

\mathbf{x}_i : 各成分は番号 i ($i=1, \dots, I$) を割り振られた地点を地理的に特定する語彙（座標, 地名に関する言語表現）、すなわち非数値情報に対応し、これらを数量化理論により数値化した値をとる。具体的には、各語彙を予測条件として指定するときに 1, そうでなければ 0 とする。

\mathbf{t}_j : 各成分は番号 j ($j=1, \dots, J$) を割り振られた離散時刻を特定する語彙（季節, 月, 曜日, 時間帯）と対応しており、予測条件として指定するときに 1, そうでなければ 0 とする。

\mathbf{s}_k : 各成分は番号 k ($k=1, \dots, K$) を割り振られた付帯状況要因を特定する語彙と対応しており、 \mathbf{t}_j における \mathbf{x}_i の付帯状況（天候や周囲環境など）を示す。各要因は予測条件として指定するときに 1, そうでなければ 0 とする。

状況ベクトルは原理的には本来、代数演算可能な数値属性（気温, 照度, 湿度, 道路幅, 集客施設までの距離 [12] など）を含むことができる。10 章では記述の容易さを重視し、語彙の数値化情報のみを用いた場合の実験例を示す。

3.19 主成分回帰分析

上記の特徴ベクトルが状況ベクトルに回帰すると仮定して多変量回帰分析を行うと、状況行列が特異となり、共分散行列に対する良好な逆行列を算出できない場合がある。このようなとき、原空間に主成分分析を施して次元を削減すれば各オブジェクトの確信度が予測可能になる。以下、その手順を示す。

まず、上記の \mathbf{x}_i と \mathbf{s}_j を成分として構成される状況ベクトルを次式で表す。

$$\begin{aligned} \mathbf{S} &= (S_1, \dots, S_L)^T = [S_{place}^T, S_{time}^T, S_{env}^T]^T \\ &= [\mathbf{x}_i^T, \mathbf{t}_j^T, \mathbf{s}_k^T]^T \end{aligned} \quad (3.58)$$

これにより， $\{i, j, k\}$ の3項組がクラス分けされた状況イベントを構成することがわかる．これを H 個のイベントについて収集すると，次式の計測行列 S ($H \times L$) が構成できる．

$$\mathbf{S} = (\mathbf{S}_{[1]} \dots \mathbf{S}_{[H]})^T \quad (3.59)$$

これをもとにして，共分散行列 V ($L \times L$) は次式で計算される．

$$\mathbf{V} = \frac{1}{H} \mathbf{S}^T \mathbf{S} \quad (3.60)$$

この固有値問題を解いて得られる V の p 番目の固有ベクトル ($p=1, \dots, P$) は次式で表現される．

$$\mathbf{w}_p = (w_{p1}, \dots, w_{pL})^T \quad (3.61)$$

これにより， p 番目の主成分に対応する主成分得点を H 個の状況ベクトルについて計算した主成分得点ベクトルは次式で表現される．

$$\mathbf{q}_p = \mathbf{S} \mathbf{w}_p \quad (3.62)$$

これを正規化した行列形式 Q ($H \times P$) に拡張して，次式で表す．

$$\begin{aligned} \mathbf{Q} &= (\mathbf{q}_1 \dots \mathbf{q}_p) \\ &= \mathbf{S} \mathbf{W} \text{diag}(1/\sqrt{\lambda_1}, \dots, 1/\sqrt{\lambda_p}) = \mathbf{S} \mathbf{W} \Delta_p^{1/2} \end{aligned} \quad (3.63)$$

ただし， $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_p)$ であり λ_p は V の p 番目の固有値を表す．すなわち，

$$\Delta_p \equiv \text{diag}(\lambda_1, \dots, \lambda_p). \quad (3.64)$$

である．

3.20 予測器

ここでは移動物体が出現することを過去の撮影結果から予測することを考える．

3.20.1 基本動作

過去の統計に基づく予測頻度の行列 F ($H \times M$) は多変量回帰予測によれば，

$$\begin{aligned} \mathbf{F} &= \mathbf{Q} \mathbf{B} \\ \mathbf{Q} &= [\mathbf{Q}^{(1)} \dots \mathbf{Q}^{(H)}]^T \end{aligned} \quad (3.65)$$

ただし， \mathbf{B} は回帰係数行列 ($P \times M$)， $\mathbf{Q}^{(h)}$ は h 番目のイベントに対応する主成分空間の状況ベクトル ($P \times 1$)，そして \mathbf{Q} は主成分空間の状況行列 ($H \times P$) を表す．これより \mathbf{B} は最小2乗法を用いて次式で算出できる．

$$\mathbf{B} = (\mathbf{Q}^T \mathbf{Q})^{-1} \mathbf{Q}^T \mathbf{Y} \quad (3.66)$$

ただし， \mathbf{Y} ($H \times M$) は H 個のイベントにおける M 種類のオブジェクト (表2では $M=12$) の出現頻度を表すデータ行列である． \mathbf{Y} は H 個の各イベントにおいて画像認識器が出力した確信度行列 \mathbf{C} ($N_B \times M_O$) から歩行者 (表2では $M_O = M/2 = 6$) と判定された画素領域に注目し， M 個の出現頻度値に変換することで生成される．ここで， N_B は画像中の認識単位である画素領域の数， M_O は画像認識が判別する歩行者オブジェクトの種類の数とする．いま，この変換操作を次式で表現する．

$$\mathbf{Y} = \Psi(\mathbf{C}) \quad (3.67)$$

この C の信頼度は画像認識の認識率 R によって定められるため、 C から生成した Y を学習データとする回帰係数行列 B は R の影響を受ける。いま、この不完全な認識器が推定した真理値を Estimated Truth (以下, ET) と呼ぶことにする。一方、人間や完全な認識器が与える正解値は一般に Ground Truth (以下, GT) と呼ばれ、その信頼度はほぼ 100% となる。しかし、ET が機械により自動獲得できるのに対し、GT の獲得は一般に多大なる労力を要する。そこで、ET を学習データとして用いることを想定し、GT も含めて次式で表現する。

$$\hat{Y} = [\hat{f}_{hm}] \quad (h=1, \dots, H)(m=1, \dots, M) \quad (3.68)$$

ただし、 f_{hm} は h 番目のイベントに対応する画像中の m 種類目の歩行者オブジェクトの出現頻度を示す。

3.20.2 ET の変化を見積もるモデル

次に、この学習データ Y^{\wedge} によって獲得できる予測性能を見積もるため、 f_{hm} を次式でモデル化する。

$$\hat{f}_{hm} = A(r_{hm}, g_{hm}) \quad (3.69)$$

ただし、

r_{hm} : 画像認識器の信頼度 (平均認識率)。

g_{hm} : 出現頻度に関する GT 値。

$A(r, g)$: 画像認識で得られる出現頻度の ET 値を認識率 r と GT 値 g から見積もるモデル関数。

である。いま、画像認識率 R を再現率 α と適合率 β の調和平均で定義すると、

$$R = \frac{2\alpha\beta}{\alpha + \beta}, \quad \alpha = \frac{g - u}{g}, \quad \beta = \frac{g - u}{g - u + d} \quad (3.70)$$

となる。ただし、 g は GT 値、 u は未検出個数、 d は誤検出個数を表す。このとき、認識率 R で計測される出現頻度 f^{\wedge} の最大値 f_{max}^{\wedge} と最小値 f_{min}^{\wedge} はそれぞれ、

$$\hat{f}_{max} = g \frac{2-R}{R} \geq \hat{f} \geq g \frac{R}{2-R} = \hat{f}_{min} \quad (3.71)$$

と表現される。今、 f^{\wedge} が上記の区間内で同じ出現確率で分布すると仮定すると、その平均値は

$$\hat{f}_{avr} = \frac{\hat{f}_{max} + \hat{f}_{min}}{2} = A_R g \quad (3.72)$$

$$A_R = \frac{1}{R} + \frac{1}{2-R} - 1$$

となる。したがって、

$$A(r, g) = \left(\frac{1}{r} + \frac{1}{2-r} - 1 \right) g \quad (3.73)$$

を式(3.69)に適用すれば Y^{\wedge} を算出できる。この Y^{\wedge} を過去の統計データとして予測頻度の行列 F の近似値 F^{\wedge} を算出すると、

$$\begin{aligned}\hat{F} &= Q\hat{B} \\ \hat{B} &= (Q^T Q)^{-1} Q^T \hat{Y}\end{aligned}\quad (3.74)$$

となる。

3.20.3 認識器の性能改善

このような予測とそれに対応する観測（画像認識）が過去に H_F 回行われたとすると、予測対象となった地点で車両が実際に観測した出現頻度の行列 F_1 と予測頻度の行列 F_0 との間の誤差 ΔF は次式で表される。

$$\Delta F = F_1 - F_0 \quad (3.75)$$

ただし、

$$F_a = [\mathbf{f}_{a,1}, \mathbf{f}_{a,2}, \dots, \mathbf{f}_{a,H_F}]^T \quad (a=0,1) \quad (3.76)$$

$$F_0 = Q\hat{B} \quad (3.77)$$

である。

車載器，センター，携帯，PC，定点カメラなどによる ET や GT は予測器の学習データとなる。このとき式 (3.75) の予測誤差から次式の評価尺度を算出する。

$$E_F = \sum_{i=1}^{H_F} \sum_{j=1}^M |\Delta F_{ij}| \quad (3.78)$$

この E_F がある閾値 E_{FTH} を超えるとき、予測器と認識器を修正すれば半自動的に性能を改善できるであろう。いま、予測対象とする地域に応じて経験的に 2 つの閾値 f_{th} と Φ_{th} を設定し、次の制約式を満たすまで学習データを増やす。

$$\{\|\mathbf{f}(\mathbf{S})\|_1 \geq f_{th}\} \text{ AND } \{\|\Phi(\mathbf{S}, N)\|_1 \geq \Phi_{th}\} \quad (3.79)$$

ただし、

$$\begin{cases} \Phi(\mathbf{S}, N) = \left\{ \sum_{n=1}^N \mathbf{r}_n \mathbf{f}_n(\mathbf{S}) + \mathbf{f}_0 \right\} \\ \mathbf{r}_n = \text{diag}(r_{n1}, \dots, r_{nm}, \dots, r_{nM}) \\ \mathbf{f}(\mathbf{S}) = \frac{1}{(N+1)} \Phi(\mathbf{S}, N) \\ \mathbf{f}_n = (f_{n1}, \dots, f_{nM})^T \end{cases} \quad (3.80)$$

である。ここで、 n は認識器の識別番号であり、 $n=1, \dots, N$ 。 \mathbf{f}_0 は初期の歩行者頻度ベクトル、 \mathbf{f}_n は認識器 n が出力する歩行者頻度ベクトル、そして、 \mathbf{f} は N 個の認識器の出力に関する平均値ベクトルである。また、 \mathbf{S} は認識時の状況ベクトルであり、 r_{nm} は認識器 n のオブジェクト m に関する認識率を表す。

この更新動作はセンターから信頼度付きの ET 値または事前に算出した認識特性を配信することで達成される。その後、上述の予測，認識，差分評価のプロセスを繰り返し、 E_F が E_{FTH} 以下ならば修正を終了する。 $E_F > E_{FTH}$ であれば、認識パラメータを改善し、認識，予測， E_F の最小化を繰り返す。

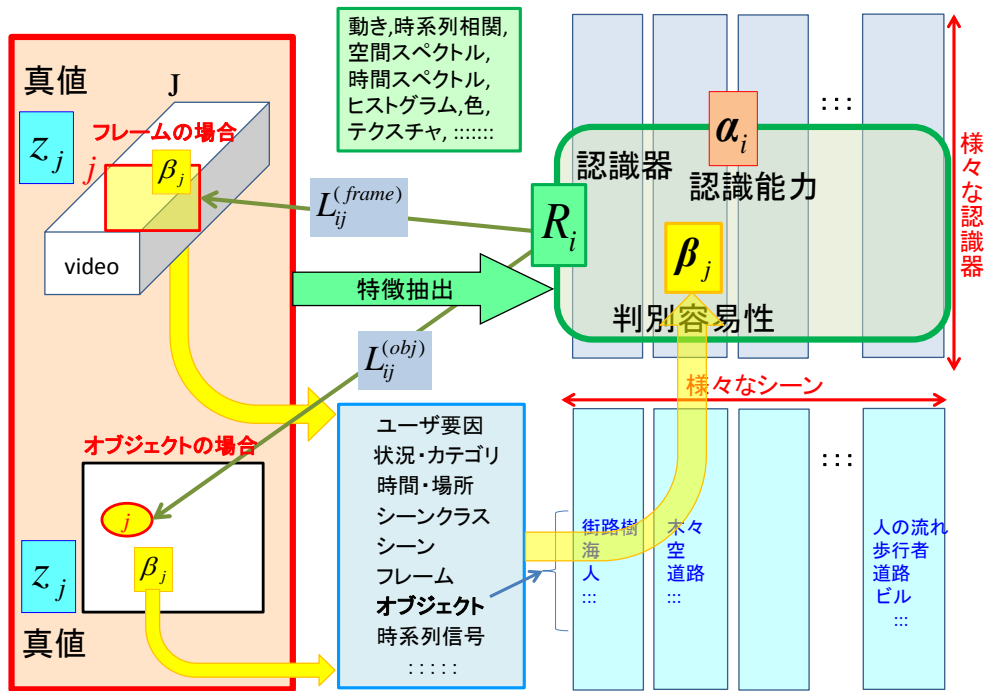


図 3.8 映像と認識器の関係

Fig. 3.8 Relation between images and recognizers.

3.20.4 EM アルゴリズムの応用

複数の認識器を巡回的に協調させるという手法をより高度なものにした例として、EM アルゴリズムを導入したGLAD (Generative model of Labels, Abilities, and Difficulties)手法[H1]がある。この手法を参考にして認識に関連する属性要因を拡張し、MBVに適用する。

いま、 K 個のブロック画像 $BLK(k), (k=1, \dots, K)$ の集合を考える。各ブロック画像は N 個の対象カテゴリの 1 つに属するとする。

このとき、ブロック画像 $BLK(k)$ のオブジェクトインデックス z_k を決定したい。観測したインデックスは主としてブロック画像の判別容易性、真のオブジェクトインデックス、認識器の認識能力の 3 つに依存する。

いま、ブロック画像 $BLK(k)$ の判別容易性を次式で表現する：

$$\beta_k \in [0, +\infty) \quad (3.81)$$

また、画像 $BLK(k)$ の真のインデックスを z_k とする。また、認識器 L_i の認識能力のパラメータを次式で表現する：

$$\alpha_i \in (-\infty, +\infty) \quad (3.82)$$

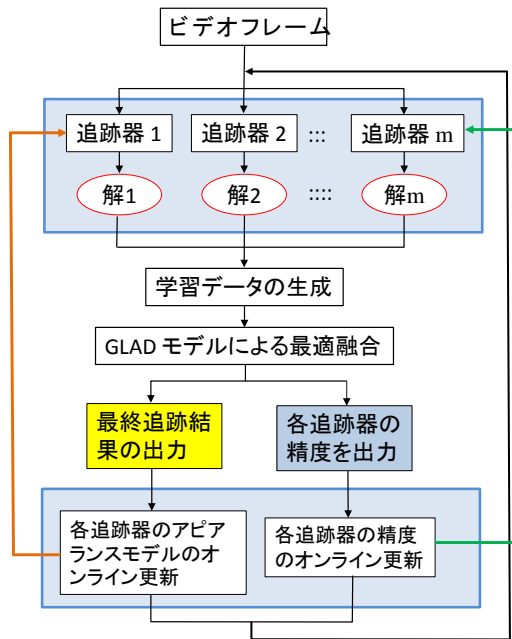


図 3.9 追跡アルゴリズム

Fig. 3.9 Tracking algorithm..

α_i は認識対象とする画像 $BLK(k)$ に関連する以下の要因の各々について定義できる。

- Level-7) ユーザ要因 : ユーザが求める認識タスクおよびユーザの嗜好性
- Level-6) 状況・カテゴリ : シーン撮影時の状況と分類カテゴリ
- Level-5) 時間・場所 : シーンが撮影された時間と場所
- Level-4) シーンクラス : シーンクラス (3.13 節参照)
- Level_3) シーン
- Level_2) フレーム
- Level_1) オブジェクト
- Level_0) 時系列信号 : $BLK(k)$ の特徴量およびそれを含む時系列

L_i が $BLK(k)$ に与えたインデックスを l_{ik} とし、それが真のインデックスに等しい確率を次式のモデルで表現する。

$$p(l_{ik} = z_k | \alpha_i, \beta_k) = \frac{1}{1 + e^{-\alpha_i \beta_k}} \quad (3.83)$$

このモデルのもとで獲得したインデックスで正しいものに関する対数尤度 (log odds) は認識能力と判別容易性の双一次関数として次式で表現される。

$$\log \frac{p(l_{ik} = z_k)}{1 - p(l_{ik} = z_k)} = \alpha_i \beta_k \quad (3.84)$$

図 3.12 はモデルの因果関係を示すグラフィカルモデルである。 z_k, α_i, β_k は既知の先験的

分布でサンプルされる。これらは式(3.83)によって観測インデックスを決定する。観測イ

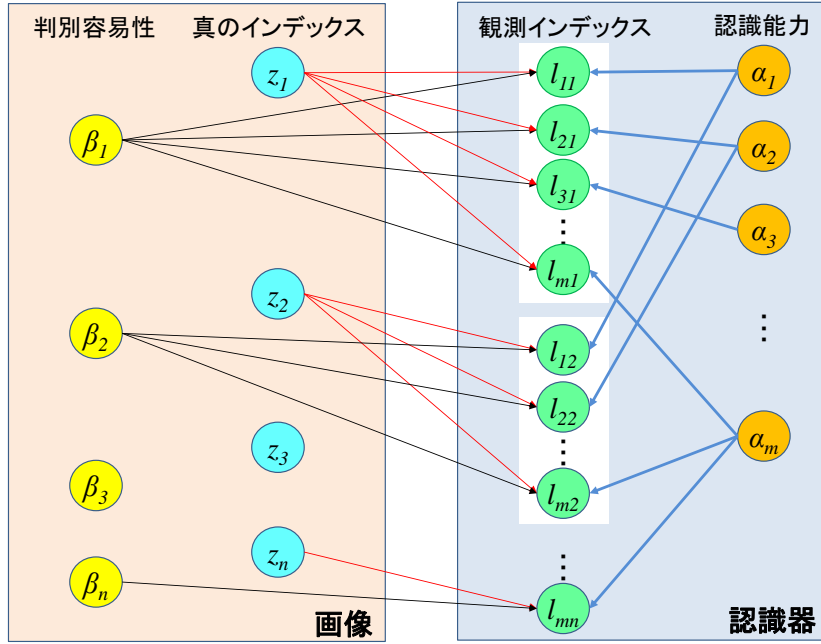


図 3.10 グラフィカルモデル

Fig. 3.10 Graphical Model..

ンデックスの集合 $L = \{l_{ik}\}$ が与えられたとき、求められるタスクは $Z = \{z_k\}$, α_i, β_k について最も尤度の高い値を同時に推定することになる。そこで、これらの最尤推定を行うために期待値最大化 (EM) 手法を適用する：

【E ステップ】

画像 $BLK(k)$ について与えられたすべてのインデックスを

$$l_k = \{l_{ik} | \hat{k} = k\} \quad (3.85)$$

と表す。ただし、すべての認識器が各ブロック画像をインデキシングするわけではない。この場合、 l_{ik} 中の認識器に関する識別番号 i は画像 $BLK(k)$ のインデックスを付与した認識器のみ参照する。すなわち、

$$i \in \{i_L(k,1), \dots, i_L(k, I_k)\} \quad (3.86)$$

であり、参照した認識器群全体の認識能力を表すベクトル $\alpha(k)$ は次式で表される：

$$\alpha(k) = (\alpha_{i_L(k,1)}, \dots, \alpha_{i_L(k, I_k)}) \quad (3.87)$$

また、 K 個のブロック画像全体の判別容易性を表すベクトル β は次式で表される：

$$\beta = (\beta(1), \dots, \beta(K)) \quad (3.88)$$

最新の M ステップにおいて得た $\alpha(k)$, β が与えられたときのすべての $z_k \in \{0,1\}$ と観測したインデックスについて事後確率は次式で計算する：

$$\begin{aligned} p(z_k | l, \alpha(k), \beta) &= p(z_k | l_k, \alpha(k), \beta_k) \propto p(z_k | \alpha(k), \beta_k) p(l_k | z_k, \alpha(k), \beta_k) \\ &\propto p(z_k) \prod_{i=i_L(k,1)}^{i_L(k, I_k)} p(l_{ik} | z_k, \alpha_i, \beta_k) \end{aligned} \quad (3.89)$$

ここで、グラフィカルモデルから得られる条件独立仮定により、 $p(z_k | \alpha, \beta_k) = p(z_k)$ と

した。

【M ステップ】

パラメータ $(\alpha(k), \beta)$ が与えられたときに観測変数と隠れ変数の組 (\mathbf{l}, \mathbf{Z}) の結合対数尤度の期待値として定義される関数 Q を最大化する。ここで前段の E ステップで計算した \mathbf{Z} の事後確率を用いる：

$$\begin{aligned} Q(\alpha, \beta) &= E[\ln \{p(\mathbf{l}, \mathbf{Z} | \alpha, \beta)\}] \\ &= E[\ln \{\prod_k (p(z_k) \prod_i p(l_{ik} | z_k, \alpha_i, \beta_k))\}] \\ &= \sum_j E[\ln \{p(z_k)\}] + \sum_{ik} E[\ln \{p(l_{ik} | z_k, \alpha_i, \beta_k)\}] \end{aligned} \quad (3.90)$$

ただし、 $\mathbf{Z}, \alpha, \beta$ があたえられたとき、 l_{ik} は条件付き独立である。

ここで、前段の E-ステップで推定した $\alpha_{old}, \beta_{old}$ を用いて \mathbf{Z} の期待値を計算する。勾配上昇法 (gradient ascent) を用いて Q を局所的に最大化する α, β の値を見つける。

ここまではフレーム内における巡回的な最適化計算である。

3.20.5 トラッキングの性能改善

上記は複数の認識または計測結果をもとに予測した結果と新たに認識した結果との差分を評価するというプロセスを繰り返すことで認識器の性能をフレーム内で改善させようとするものであった。したがって、フレーム間の認識性能については依然として各認識器の個別の性能に頼らざるを得ない。

これに対し、近年、ラベリングそのものをトラッキングすなわちフレーム間認識のレベルで改善する試みが行われている。その手法は以下の3つに大別される：

- (1) マルコフ確率場 (MRF) を用いる手法 (東大上条研など)
- (2) モンテカルロフィルタ、パーティクルフィルタによる手法 (INRIA 他)
- (3) 上述の GLAD 手法を時間方向に拡張した手法 [H2] (ハルビン工科大他)

ここでは複数の認識器を協調させるという点で前節の発展形に位置づけられるとの観点から、(3) に着目してその概要を示す。

いま、 m 個の不完全な認識器が動画像の入力列 $\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_T$ を観測すると考える。ただし、

$$\mathbf{s}_t = \{x_t^1, x_t^2, \dots, x_t^{N_t}\} \quad (3.91)$$

はフレーム t で得られたサンプル集合であり、 x_t^k は \mathbf{s}_t における k 番目のブロック画像 $BLK(k)$ から構成される特徴ベクトルに相当する。 x_t^k に対するクラスインデックスの候補集合を次式で記述する。

$$\mathcal{L}_t^k = \{i^k | \hat{k} = k\} \quad (3.92)$$

ただし、 $l_t^{ik} \in \{0, 1\}$ は i 番目の認識器によって与えられる。弱い学習によるトラッカーの設定は以下のように行う。

以下の手順を $t=1,2,\dots$ について繰り返す。

Step-1) サンプル集合 $\mathbf{s}_t = \{\mathbf{x}_t^1, \mathbf{x}_t^2, \dots, \mathbf{x}_t^{N_t}\}$ を受け取る

Step-2) この段階ではまだ \mathbf{s}_{t-1} までしか学習していない不完全認識器群を用いて、各 \mathbf{x}_t^k について $\mathbf{L}_t^k = \{l_t^{ik} | \hat{k} = k\}$ を予測する

Step-3) \mathbf{s}_t の正解インデックスと上記予測結果を比較し、各不完全認識器の認識精度 (accuracy) を算出する。

Step-4) この認識精度と正解インデックスを \mathbf{s}_t とともに学習する。

Step-5) ロバストなトラッキング結果であれば不完全認識器群を更新する。

学習データは *heuristic* に選ばれるが、学習完了後はあらたな学習データがなくても、入力に対して最も妥当な bounding box の位置ベクトルと各トラッカーのトラッキング精度を同時に推定する。

m 個の不完全認識器が時刻 t のフレームにおける現在の探索窓においてブロック画像集合から得られる $\mathbf{s}_t = \{\mathbf{x}_t^1, \mathbf{x}_t^2, \dots, \mathbf{x}_t^{N_t}\}$ を受け取るとする。ここで、 \mathbf{x}_t^k は $BLK(k)$ から構成される特徴ベクトルである。 $\mathbf{P}_t = \{\mathbf{p}_t^1, \mathbf{p}_t^2, \dots, \mathbf{p}_t^m\}$ を時刻 t のフレームについて不完全トラッカー群 $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ が出力したトラッキング結果 (bounding box で表示) の集合とする。また、 $D(\mathbf{a}, \mathbf{b})$ を2つのブロック画像 \mathbf{a}, \mathbf{b} 間の中心位置の距離 (画素) とする。そしてトラッカー T_i が出力した \mathbf{x}_t^k のラベルを $l_t^{ik} \in \{0,1\}$ と記述する。各トラッカー T_i についてフレーム t の正例に属する上位 N_{BEST} 個の確信度 (トラッカー T_i で推定される) を有するブロック画像の集合を Top_t^i で表す。 \mathbf{s}_t 中の各ブロック画像についてトレーニング集合 (学習データ) に入れるかどうかの選別を以下の手順で行う：

```

For  $k = 1, \dots, N_t$  {
  If  $(\mathbf{p}(\mathbf{x}_t^k) \in \mathbf{P}_t = \{\mathbf{p}_t^1, \mathbf{p}_t^2, \dots, \mathbf{p}_t^m\})$ 
  {
    For  $i = 1, \dots, m$  {
      If  $(D(\mathbf{p}(\mathbf{x}_t^k), \mathbf{p}_t^i) \leq D_{th} \parallel \mathbf{x}_t^k \in \text{Top}_t^i)$   $\{l_t^{ik} = 1\}$ 
      else  $\{l_t^{ik} = 0\}$ 
      endif
    }
  } else { neglect the image patch  $\mathbf{x}_t^k$  }
}

```

ここで、 $\mathbf{p}(\mathbf{x})$ は特徴ベクトル \mathbf{x} を構成するブロック画像の位置情報ベクトルを表す。文献[H2]では $N_{\text{BEST}}=10$ 、 $D_{th}=5$ と設定されているが、採用した認識器や問題に応じて変わると考えられる。

さらに、ブロック画像 \mathbf{x}_t^k の判別容易性 β_t^k を m 個の不完全認識器による多数決法で初期化すると次式のようなになる：

$$\beta_t^k = \exp\left(\kappa_1 \times \left|\frac{2}{m} \sum_{i=1}^m l_t^{ik} - 1\right|\right) \quad (3.93)$$

ここで、 κ_1 は正規化係数である。

m 個のトラッカーの精度 α_i^t , ($i=1, \dots, m$) について時間方向の更新は次式で行われる。

$$\alpha_i^t = (1 - w_i^t) \alpha_i^{t-1} + w_i^t M_i^t \quad (3.94)$$

ただし、 w_i^t は学習率であり、また、 M_i^t は正解との比較により成功したと判断されるトラッカーについて 1、それ以外のモデルでは 0 となる。誤った更新を避けてロバスト性を最大化するために、学習率の設定は次式のように適応的に行う：

$$w_i^t = \begin{cases} \frac{\kappa_2}{1 + e^{-\alpha_i}}, & \text{if } M_i^t = 1 \\ \frac{\kappa_2}{1 + e^{\alpha_i}}, & \text{if } M_i^t = 0 \end{cases} \quad (3.95)$$

ただし、 κ_2 は正規化係数であり、文献[H2]では 0.1 に設定されている。トラッカー i が成功して精度 α_i が大きいことがわかれば、学習率 w_i^t を大きい値に設定して迅速に適応させる。成功しない場合は w_i^t を小さくする。

以上により学習段階が完了し、以後、各トラッカーにおいて予測が実行できる。

3.21 複数の認識エージェント

時系列ステレオの類推から、過去のデータに基づく映像予測は複数の視点による認識と同等に考えることができる。ここで、個々の認識器の動作を単なる画像センシングのみならず、他の認識器との協調、さらには情報交換によるマクロ概念の形成という動作としてとらえれば、個々の単体は認識器というよりも認識エージェントと呼ぶほうがふさわしい。複数の認識エージェントの協調には次のような効果がある。

- 1) 認識率を向上させる
- 2) 時空間に伴う変化を予測する
- 3) 新しい認識属性を獲得する（複数の 2 次元画像が立体構造を再構成するように）

もちろん、認識が完結する時間に制約を設けなければ、単体のエージェントが観測を繰り返すことで上記と同等の効果を得ることは可能である。しかし、リアルタイムに複数の場所に存在することができなければならないようなケースには対処できない。

認識エージェントの協調には様々な応用例がある。ロボット、プローブカー、Web 上の人工知能、人間社会における知の形成の説明などがそれにあたる。本論文ではこれらの説明は割愛し、MBV における動画像認識に焦点をあてる。

参考文献

A. 筆者文献

- [1] M. Sasaki, "A Proposal for High-Efficiency Coding of the Definite Image Sequence", *PCSJ88*.
- [2] 笹木, "モデルベース顔領域抽出と領域別量子化による顔画質向上", 1997 年映像情報メディア学会年次大会 (ITE'97), p.381-382, 1997 年 7 月.
- [3] M.Sasaki, A.Toyoda, "3D Motion Estimation of a Moving Object using Model-Based Techniques", ITS95 pp.1075-1081, November 1995.
- [4] 笹木, 難波, 粉川, 小川, "車載情報通信におけるメディア適応化手法の提案 - 制約充足問題を用いた適応化手法の検討 -", 信学技報 ITS2004-22, Vol.104 No.325, pp.35-42, 2004 年 9 月
- [5] N. Day, S. Sekiguchi and M. Sasaki, "21 Mobile Applications" in *Introduction to MPEG-7* edited by B. S. Manjunath et al. , pp. 353-361, WILEY, 2002.
- [6] 笹木美樹男, 難波秀彰: "ITSにおける知的情報支援とそれを支える視覚通信複合技術 - 知識共有型の環境認識をめざして -", 2006 年度人工知能学会全国大会 (JSAI 2006) 1E2-2, 2006 年 6 月.
- [7] 笹木美樹男: "統計的推論に基づく景観認識の検討", 情報処理学会研究報告, 2006 AVM-54-8 (2006).
- [8] 笹木美樹男: "色情報と知識処理に基づく車載カメラ映像のインデキシング", デンソーテクニカルレビュー Vol.12 No.1 2007, p.58-68.
- [9] Mikio Sasaki, Kenji Muto, Kunihiko Sasaki: "Estimation of Driving Safety from On-Board Captured Scene by using Probe Image Database", *ITST2007*, June 2007.
- [10] Mikio Sasaki, Hideaki Nanba: "Pedestrian Tracking by using MPEG-based Video Signal Processing", *IEICE Technical Report*, vol. 108, no. 344, IE2008-106, pp. 17-22, Dec. 2008.
- [11] Mikio Sasaki: "MPEG-based Scene Class Recognition", *IEICE Technical Report*, vol.109, no.148, IE2009-55, July 2009.
- [12] 笹木美樹男: "MPEG 映像に適した景観認識手法の提案", 画像電子学会誌, vol.38, no.6, pp.890-899 (2009-11).
- [13] Mikio Sasaki: "MPEG-based Scene Class Recognition", http://homepage3.nifty.com/MikioSasaki/scene/090724_IEICE_SceneClass.pdf

B. 画像符号化・映像メディア処理

- [1] "Information technology - Coding of moving pictures and associated audio for digital storage media at up to about 1,5 Mbit/s - Part 1: Systems", ISO/IEC11172-1 First edition 1993-08-01.
- [2] "Information technology - Coding of moving pictures and associated audio for digital storage media at up to about 1,5 Mbit/s - Part 2: Video", ISO/IEC11172-2 First edition.
- [3] 石川, 相澤, 羽鳥, "簡易モデルのあてはめによる 3 次元仮想空間の構築手法", 画像符号化・映像メディア処理シンポジウム (PCSJ・IMPS98) I-1.4, pp.9-10, Oct.26-28, 1998.
- [4] ISO/IEC 15938-3 Information technology - Multi-media Content Description Interface Part 3: Visual.
- [5] ISO/IEC 15938-5 Information technology -- Multimedia Content Description Interface Part 5: Multimedia Description Schemes.
- [6] B.S.Manjunath *et al.*, "Introduction to MPEG-7", WILEY, 2002.
- [7] Youti Horry, Ken-ich Anjyo and Kiyoshi Arai, "Tour into the picture", Proc. SIGGRAPH'97, Los Angeles, California, pp.225-232, August 3-8, 1997.
- [8] 岩崎敏紀, 横山貴紀, 渡辺俊典, 古賀久志, "MPEG ビデオデータの動きベクトルを用いた圧縮領域における移動物体の検出と追跡", 電子情報通信学会論文誌 D, Vol. J91-D, No. 6. (June 2008), pp. 1592-1603.

C. シーン理解

- [1] Hoiem, A. A. Efros and M. Hebert, "Recovering Surface Layout from an Image", *International Journal of Computer Vision* 75(1), 151-172, 2007.
- [2] G. Pirirou, P. Bouthemy, and J.-F. Yao: "Recognition of Dynamic Video Contents With Global Probabilistic Models of Visual Motion", *IEEE Trans. Image Process.*, vol. 15, no.11, pp.3417-3430, Nov. 2006.
- [3] "DIGITAL VIDEO RETRIEVAL at NIST - TREC Video Retrieval Evaluation", [http://www-nlpir.nist.gov/projects/trecvid/](http://www.nlpir.nist.gov/projects/trecvid/)

- [4] J. Vogel, B. Schiele, "Semantic Modeling of Natural Scenes for Content-Based Image Retrieval", *International Journal of Computer Vision* 72(2), 133-157, 2007.
- [5] Anna Bosch, Andrew Zisserman, and Xavier Munoz, "Scene Classification Using a Hybrid Generative/Discriminative Approach", *IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI)*, vol.30, no.4, pp.712-727, Apr. 2008.
- [6] Joakim Jitén Söderberg: "Multidimensional Hidden Markov Model Applied to Image and Video Analysis", *A Thesis for the degree of DOCTOR OF PHILOSOPHY*, Institute Eurecom Sophia Antipolis, February 2007.
- [7] Jianbo Shi and Carlo Tomasi. Good features to track. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 593–600, 1994.
- [8] Ashutosh Saxena, Min Sun, and Andrew Y. Ng, "Make3D: Learning 3D Scene Structure from a Single Still Image", *IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI)*, vol.31, no.5, pp.824-840, May 2009.
- [9] D. Hoiem, A. A. Efros and M. Hebert: "Recovering Surface Layout from an Image", *International Journal of Computer Vision* 75(1), 151-172, 2007.
- [10] "DIGITAL VIDEO RETRIEVAL at NIST - TREC Video Retrieval Evaluation", <http://www-nlpir.nist.gov/projects/trecvid/>
- [11] 浅田, 白井, "マルチセンサ情報を動的に統合することによる道路シーンの解釈とモデリング", *情報処理学会論文誌*, Vol.31, No.12, pp.1743-1754 (Dec.1990).
- [12] 杉山, 目黒, 金子, "少数の透視投影画像に基づいて生成された仮想空間における3次元移動表現", *信学技報 MVE* 2001-148 (2002-3).
- [13] Vittorio Ferrari, Andrew Zisserman, "Learning Visual Attributes", <http://eprints.pascal-network.org/archive/00003146/01/ferrari07.pdf>
- [14] Bangpeng Yao, Gary Bradski, Li Fei-Fei, "A Codebook-Free and Annotation-Free Approach for Fine-Grained Image Categorization", http://ai.stanford.edu/~bangpeng/document/cvpr12_0654.pdf

D. 因子分解法

- [1] 金出武雄, コンラッド・ポールマン, 森田俊彦, "因子分解法による物体形状とカメラ運動の復元", *電子情報通信学会論文誌 D-II*, J74-D-II-8, pp. 1497-1505, Aug.1993.
- [2] J. P. Costeira, "A Multibody Factorization for Independently Moving Objects", *International Journal of Computer Vision* 29(3), 159-179(1998).

E. 歩行者検出・トラッキング

- [1] Jianbo Shi and Carlo Tomasi. Good features to track. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 593–600, 1994.
- [2] P. Viola, M. J. Jones, "Robust Real-time Object Detection", *Cambridge Research Laboratory Technical Report Series CRL2001/01*, Feb. 2001.
- [3] G. Antonini, "A discrete choice modeling framework for pedestrian walking behavior with application to human tracking in video sequences", *Thèse EPFL*, no 3382 (2005).
- [4] [http://www.gavrila.net/Research/Pedestrian Detection/Pedestrian Projects/SAVE-U Final Report.pdf](http://www.gavrila.net/Research/Pedestrian%20Detection/Pedestrian%20Projects/SAVE-U_Final_Report.pdf)
- [5] Xiang, Tao; Gong, Shaogang; "Beyond Tracking: Modeling Activity and Understanding Behaviour.", *International Journal of Computer Vision*, Vol. 67, No. 1, April 2006 , pp. 21-51(31).
- [6] J.-M. Odobez, D. Gatica-Perez, and S. O. Ba, "Embedding Motion in Model-Based Stochastic Tracking", *IEEE Trans. Image Processing*, vol.15, no.11, pp.3514-3530, Nov. 2006.
- [7] E. Arnaud, E. Mémin, "Partial Linear Gaussian Models for Tracking in Image Sequences Using Sequential Monte Carlo Methods", *International Journal of Computer Vision* 74(1), 75-102, 2007.
- [8] Bo Wu, Ram Nevatia, "Detection and Tracking of Multiple, Partially Occluded Humans by Bayesian Combination of Edgelet based Part Detectors", *International Journal of Computer Vision* Volume 75, Issue 2(November 2007), 247 – 266.
- [9] M. Lyell, R. Flo, M. Mejia-Tellez, "Simulation of Pedestrian Agent Crowds, with Crisis", <http://neCSI.org/events/iccs6/viewpaper.php?id=367>.
- [10] D. M. Gavrila, S. Munder, "Multi-cue Pedestrian Detection and Tracking from a Moving Vehicle", *International Journal of Computer Vision* 73(1), 41-59,2007.
- [11] J. Tilp, et al., "Pedestrian Protection based on Combined Sensor Systems", Paper Number 05-0156, <http://www-nrd.nhtsa.dot.gov/pdf/nrd-01/esv/esv19/05-0156-O.pdf>
- [12] Ahmed Elgammal, Chan-Su Lee, "Tracking People on a Torus", *IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI)*, vol.31, no.3, pp.520-538, Mar. 2009.

F. 一般物体認識

- [1] 神谷保徳, 矢野良和, 大熊繁, 高橋友和, 井手一郎, 村瀬洋: "一般物体認識のためのタイプの異なる局所特徴の統合利用", 信学論 D, Vol. J92-D, No.5, pp628-638 (2009)
- [2] 三浦純, 森田英夫, ヒルドミヒヤエル, 白井良明: "SVM による物体と位置の視覚学習に基づく屋外移動ロボットの位置推定. 日本ロボット学会誌, 25(5), pp. 792-798, 2007.
- [3] Wan-Lai Zhao, Chong-Wah Ngo, "lip-vireo Manual Version 1.01", <http://www.cs.cityu.edu.hk/~wzhao2/lip-vireo.htm>

G. 予測・推定・多変量解析・時系列分析

- [1] J.-S. R. Jang: "ANFIS: Adaptive-Network-Based Fuzzy Inference System", *IEEE Trans. Syst., Man, Cybern.*, vol. 23, no.3, pp.665-684, May/June 1993.
- [2] Y. Takane, Michael A. Hunter, "Constrained Principal Component Analysis: A Comprehensive Theory", *AAECC 12*, 391-419 (2001).
- [3] 滝口哲也, 有木康雄, "Kernel PCA を用いた残響下におけるロバスト特徴量抽出の検討," 情報処理学会論文誌, Vol.47, No.6, pp. 1767-1773, 2006. <http://www.me.cs.scitec.kobe-u.ac.jp/~takigu/pdf/2005/kpca.pdf>

H. GLAD 法

- [1] Jacob Whitehill, Paul Ruvolo, Tingfan Wu, Jacob Bergsma, and Javier Movellan, "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise", <http://mplab.ucsd.edu/~jake/OptimalLabeling.pdf>
- [2] B. Zhong, et al., "Visual Tracking via Weakly Supervised Learning from Multiple Imperfect Oracles", *CVPR2010*.

I. ITS 関連

- [1] K. Muto, K. Matsugatani, M. Egawa, and H. Nanba, "Omni-Directional Picture Database Using Probe Cars", *12th World Congress on ITS*, 6-10 Nov. 2005.
- [2] Thaned SATIENAM, "A Study on Pedestrian Accidents and Investigation of Pedestrian's Unsafe Conditions in Khon Kaen Municipality, Thailand", *Journal of the Eastern Asia Society for Transportation Studies*, Vol.5, October, 2003.
- [3] 関本義秀, 菊地英一, 佐藤圭一, 秋山祐樹: パーソントリップデータを活用した人の流れの時空間的な詳細化, 第 28 回交通工学研究発表会論文集, pp.197-200, 2008.
- [4] <http://www.dft.gov.uk/pgr/roads/vehicles/intelligentspeedadaptation/>
- [5] 熊谷, 伏木, 横田, 君田, "特徴空間射影によるプローブカーデータのリアルタイム補完", 情報処理学会論文誌, Vol.47, No.7, pp.2133-2140 (2006).
- [6] T. Fushiki, et al, "Traffic Condition Prediction by Use of Floating Cars", *IPSJ Journal* Vol. 43 No. 12, pp. 3801-3808, Dec. 2002.
- [7] M. Kumagai, et al, "Efficient Forecast Process for Nationwide Traffic Information Services", *Technical Report Of IEICE ITS2004-20(2004-09)*
- [8] T. Fushiki, et al, "Traffic Condition Prediction by Use of Floating Cars", *IPSJ Journal* Vol. 43 No. 12, pp. 3801-3808, Dec. 2002.

J. その他

- [1] 北, 津田, 獅々堀, "情報検索アルゴリズム", 共立出版, 2002.