

MPEG-Based Vision

2015 年 月 日

笹 木 美樹男

目次

1.	序論	4
1.1	本論文の背景	4
1.2	本論文の目的	6
1.3	本論文の構成	7
2.	関連研究	10
2.1	景観認識	10
2.2	一般物体認識	10
2.3	MPEG 画像の解析	10
2.4	歩行者認識	11
2.5	ヒューマントラッキング	11
2.6	構造復元	11
2.7	シーンの予測	11
2.8	アクティブセンシング	12
3.	理論	13
3.1	はじめに	13
3.2	確信度による認識	14
3.2.1	確信度ベクトル	14
3.2.2	確信度ベクトルの推定	15
3.2.3	確信度の線形回帰モデル	15
3.3	統計的な推論	15
3.3.1	多変量回帰分析による学習	15
3.4	シーンの認識	16
3.5	統計的推論手法	16
3.6	状況ベクトル	17
3.7	計測行列の作成	17
3.8	認識率	18
3.9	推論特性の合成	19
3.9.1	回帰係数行列の合成	19
3.9.2	計算量の評価	19
3.10	主成分分析による多変量回帰分析	20
3.11	カーネル主成分分析	22
3.11.1	学習フェーズ	22
3.11.2	認識フェーズ	22
3.12	シーンの構造モデル	24
3.13	シーンクラスの生成条件	24
3.13.1	主観的シーンクラス	24
3.13.2	計算的シーンクラス	24
3.14	意味との相互作用	25

3.15	アプリケーションと語彙	27
3.16	複数手法の統合化	28
3.16.1	3S 認識モデル	28
3.17	映像内容の予測と推定	29
3.18	予測動作における状況ベクトル	30
3.19	主成分回帰分析	31
3.20	予測器	31
3.20.1	基本動作	31
3.20.2	ET の変化を見積もるモデル	32
3.20.3	認識器の性能改善	33
3.20.4	EM アルゴリズムの応用	34
3.20.5	トラッキングの性能改善..... エラー! ブックマークが定義されていません。	
3.21	複数の認識エージェント	37

1. 序論

1.1 本論文の背景

インターネットやマルチメディア、デジタル放送の分野で普及している MPEG 系の画像圧縮技術との親和性が高い映像処理技術は、相互運用性の観点から強く望まれている。MPEG-7 や MPEG-21 などのマルチメディア記述に関する標準化動向もそれに呼応するが、内容の自動記述に関してはまだ研究開発途上にある。ここで、映像を復号してから認識するのではなく MPEG データとして符号化された特徴量を用いて認識できる手法があれば、復号後にビットマップデータを再処理する必要がなくなる。その結果、演算量やメモリを削減でき、通信帯域も有効利用できるという点で大変魅力的である。

1980 年代後半、マルチメディア技術の国際標準規格の根幹となる動画像符号化方式を策定するための大きな動きとして、MPEG の国際標準化活動は開始された。MPEG はまず、CD に動画を蓄積することを手掛けた。当初 ITU-T H.261 などの TV 会議用画像伝送・圧縮アルゴリズムを母体として検討が開始されたが、蓄積系の画質要求やランダムアクセス性を考慮した結果、現在の MPEG 標準の第 1 弾として 1993 年に MPEG-1 が国際規格化された。続く第 2 弾として現在の DVD や地デジの基本となる MPEG-2 標準が MPEG-1 をベースにして標準化された。当時は符号化には専用プロセッサが必要であり、復号も PC 上のソフトウェアのみで行うにはまだ十分ではなかった。

MPEG の台頭により動画像圧縮の基本方式がほぼ収束していったこともあり、画像符号化分野における次世代の研究テーマの一つとして知的符号化が提唱されるようになった。これは従来の符号の概念を言語や知識のレベルに高めようというもので、符号化対象物の 3 次元モデルを送受信側で共通の知識として、移動・回転情報や変形情報のみを符号化伝送する（モデルベース符号化）ことにより飛躍的な低ビットレート化をはかるという研究がなされた。一方で、これを人の顔に適用した研究成果は表情合成や超低ビットレート TV 電話などの基本アルゴリズムとして検討された。また、一方では車両や建物などの三次元位置姿勢推定手法としても研究がなされた。

モデルベース符号化の最大の難点は精密な三次元モデルの自動獲得手法が確立されていないことと画質の不自然さにあった。前者に関しては移動ステレオや因子分解法、非線形最適化手法などが有力な候補として検討された。特に因子分解法は 1990 年代から 2000 年代までに精力的に研究がなされたが、いまだなお完全なる実用化には至っていない。

この三次元形状獲得の流れに並行してシーンの認識対象を概念レベルで構造化し、記述するための記述手法に関する研究も行われた。これは主として Web や放送編集分野における画像検索応用を目的とするため、自動記述まで必須とはされていないが、近年、その人による記述作業の低減のため、自動化技術の開発が望まれている。そしてこの自動記述は画像理解とクロスする。

画素単位の処理によるシーン認識については多くのコンピュータビジョン関連の研究がある。だが最近ではさまざまな応用において、莫大な分量の圧縮画像がベースバンド画像にとってかわろうとしている。ビデオ圧縮においてもっとも影響力のある形式のひ

とつは MPEG であり，これはブロック画素を主として MC-DCT (Motion Compensation - Discrete Cosine Transform) で処理する．また，MPEG の ISO 標準の恩恵により，対応する符号化デバイスも低コストで広く用いられている．過去 20 年間に於いて，疑いなく MPEG は大量の洗練された技術，ソフトウェア，ハードウェア，システムによりバックアップされてきた．

コンピュータの発展とあいまって，1980 年代後半から映画の CG 技術は飛躍的に発展を遂げ，作成，表示において素晴らしい成果を収めてきた．だが，一方で内容に関する記述の生成や意味理解に関しては現時点においても，思ったほど自動化はできていない．しかし，MPEG による動画像メディアの普及により，それにつながる技術として，画像理解を 1 枚のフレームで完結させるのではなく，時間方向の連続フレームのひとつままとまりに対するインデキシングが考えられた．これは認識対象を概念レベルで XML などの構造化記述言語を用いて定義しておくことで達成される．このインデキシングの標準化作業もメディア記述へのアプローチとして ISO/MPEG-7 や MPEG-21 などで行われた．

それらの動きと並行してインターネットコンテンツとしての MPEG の普及はますます進み，さらに，デジタル放送としての MPEG や H.264 がそれに拍車をかけた．その結果，いまではデジタルハイビジョンが一般家庭で当たり前になり，携帯電話やスマホでも自在に閲覧できるようになった．そして，次の世代に向けてスーパーハイビジョンへの歩みがすでに開始されている．このように，符号化・放送・通信・メディア応用は急速に進歩を遂げた．

以上により，放送，編集，インターネットなどにおける動画像処理の応用で最も普及



図 1.1 “beach”のシーン例

Fig. 1.1 Scene examples of “beach”.

している符号化形式である MPEG を対象として認識理解を考えることは、従来コンピュータビジョンで行われてきたビットマップベースの画像処理に比べて極めて現実的である。これは、ユーザエンドで用いる画像符号化形式が BMP や AVI であることは実はあまり多くないからである。

従来、画像理解手法は主として 3 次元形状の獲得や領域分割に注力し、意味へのマッピングは応用上の課題として切り分ける傾向にあった。一方で近年、シーンに出現するオブジェクトをあらかじめ認識辞書中の有限個の語彙として登録し、画像をブロック画素単位で色解析することで辞書中の語彙に対応付け、景観を概略記述する方法が提案されている[B5][B6]。必ずしもベースバンドビデオを処理の出発点とするのではなく、映像メディアの普及・伝搬・加工に伴う表現形態を考慮した認識手法は近年の新しい方向性を示している。本論文では MPEG の特徴量の認識と知識処理によって景観を認識するシステムを MPEG-Based Vision と呼び、その実際に即した技術検討を行う。

1.2 本論文の目的

MPEG 映像データは生成された時点で DCT (Discrete Cosine Transform) 係数と動きベクトルに関する特徴抽出を完了している。本論文ではこの MPEG 映像の特徴量を利用し景観を認識する手法を提案する。ここでは、一般のコンピュータビジョンのようなリッチな特徴量は用いず、人間の知識処理に相当するメカニズムを念頭に置いて実現をはかる。具体的には、過去の知識データベースの効率的利用、異なる複数の認識エージェントによる情報共有と知識交換などがその本質に相当する。

例えばたくさんのシーン（静止画，動画ともに）の中から“beach”という単語を付与できるシーンを分類することを考える。人が選出した結果の例を図 1.1 に示す。だが、そ

シーンの自動分類ができる...

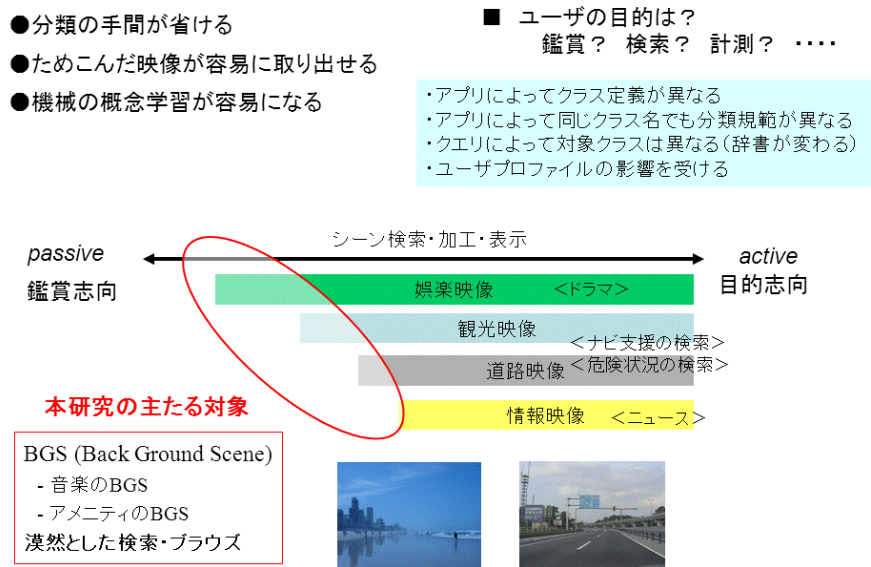


図 1.2 本研究の概観 → 要修正

Fig. 1.2 Targeted overview.

の分類作業の際、対象が静止画であってもその静止画から想起される人間の視野での時系列事象（個人の記憶により異なる）が処理対象となっている。これはすなわち、画像から獲得される信号情報のみが分類に使われているのではなく、想起される知識と記憶が脳内で評価された結果として分類が達成されていると考えられる。つまり、分類の本質は信号処理に知識と記憶による処理が加わって初めて人間に近いものになる。同様に、分類の単位を静止画、動画というレベルから 1 フレーム内のオブジェクトのレベルに下げても、その判別規範は単なる信号処理だけではなく、知識と記憶に基づく処理が大きく影響すると考えられる。

このような知的処理の応用にはまず、シーンの自動分類などがある。シーンの自動分類ができれば

- 1) 分類の手間が省ける
- 2) ためこんだ映像が容易に取り出せる
- 3) 検索や鑑賞を容易にする
- 4) 機械のレコメンド
- 5) 機械の概念学習が容易になる

といった恩恵がある。

第 2 の応用としては、環境理解や特定の事象（交通事象、自然風景など）の認識、移動体検出、歩行者検出といったものが考えられる。これらは知識処理や処理対象を工夫することでセキュリティシステム、遠隔監視、行動理解などにも応用できると考えられる。

第 3 の応用としては、複数の認識器の協調により概念を獲得するプロセスの実現が考えられる。認識器を搭載した複数エージェントによる協調的理解についてはインターネットや携帯、無線通信を介しての幅広い応用が期待されており、ヒューマンプローブやプローブカー、ロボットなどでの展開が見込まれる。

近年のニューロサイエンスの進歩により、将来的には脳内イメージの視覚化[ATR]やニューロコンピュータによる生体視覚の模倣[IBM]が可能になるとみられるが、今日のデジタル民生機器と同様に一般に普及するのは 2030 年以降のことと考えられる。したがって本論文では、現行のマルチメディア技術の延長上で景観認識や画像理解のソリューションを模索する。

1.3 本論文の構成

本論文では主として MPEG 動画像を学習し、ブロック画素単位でのインデキシングをベースとして最終的にはシーンを理解することを目指す。ただし、画像処理にすべてを委ねるのではなく、絶えず知識処理に重点をおいたメカニズムを想定して認識システムを構成する。

「2 章 関連研究」では、画像符号化、静止画理解、動画像理解、複数認識器の協調に関連する研究動向を概観する。

「3 章 理論」では本研究でめざす MPEG Based Vision が画像符号化に立脚した画像認識システムであり、蓄積系や通信系の情報システムとの連携により信号処理のみならず知

識処理も併用した画像理解が可能になることを理論的に示す。

「4章 ルールベースによる景観認識」では MPEG 映像の特徴量を利用して景観を認識する手法の第 1 段階として、ルールベースによる手法を提案する。まず、符号化特徴量と状況情報からブロック画素単位で特徴ベクトルを形成する。この特徴ベクトルから語彙の確信度ベクトルを推定し、最大確信度の語彙を認識結果とする。

「5章 景観認識を用いた歩行者認識の改善」では MPEG 映像の特徴量を利用して抽出した景観情報に基づき、従来の静止画パターン認識ベースの歩行者認識に含まれる誤検出や未検出をオブジェクトの位置関係を用いて削減する手法を述べる。ここで、3章で述べた多変量回帰分析による学習を景観認識に導入する。

「6章 統計的推論に基づく景観認識」では、5章の静止画ベースの多変量回帰分析をさらに発展させ、MPEG 動画像の特徴量である動き補償差分に注目した動画像の学習とそれに基づく景観認識を行う。また、各オブジェクト毎の認識率に関する時間変化についても考察する。

「7章 移動体検出」では、MPEG 動画像の特徴量を時系列分析することにより、歩行者や車両などの移動体を検出する。特に、動き補償差分の DCT 係数で構成する時系列ベクトルについて主成分分析やカーネル主成分分析を行う。また、時空間画像についても考察する。

「8章 シーンクラスと景観認識」では、6章で行った景観認識手法をもとに複数のシーン間の交差学習を行い、認識率の評価に基づき、有向グラフを形成してシーンのクラスタリングを試みる。また、主成分分析結果に FFT を施すことで樹木の揺れや白波、歩行者特徴などを景観から抽出する。

「9章 予測と認識の相互作用」では、移動体検出を物理的な対象物の出現予測プロセスに位置づけ、単なる画像処理ではなく予測と認識の相互作用と考えた。具体的には車載カメラや定点カメラ、携帯電話などから得る出現頻度情報をセンサーや車載装置に集積して得られる時空間分布から、想定する時間帯に所定の区間内で出現する歩行者の頻度を予測することを試みる。

「10章 カメラ移動に基づく認識」では、2次元処理や予測で得た推定解を意味のレベルから評価し、改善することでシーンクラスを判別し、その結果として推定解を正解に近づけることを考える。ここでは因子分解法を導入し、カメラ軌道に応じた三次元復元の度合いを考察する。さらに、意味を扱う次元に拡張した一般化因子分解法を試行し、誤認識低減をめざす。

「11章 複数の認識器の協調」ではひとつの視点では認識困難なシーンについて、複数の

認識器を強調させることで認識能力を向上させることを考える。さらには、異なる認識能力を持つ複数の認識器を集合知演算で統合することにより、正解を教示しなくても推定解が正解に近づいていくことを示す。

12章では全体のまとめと今後の展望について述べる。

2. 関連研究

2.1 景観認識

J. Vogelら [C4]は720×480 pixelsの自然風景の静止画を72×48 pixelsのブロック画素単位に区切り，個々のブロック画素について色，エッジ，濃度共起行列からなる180次元の特徴ベクトルを算出した．さらに，それらをSVMで9個のカテゴリに分類し，個々のカテゴリのシーン全体における割合を9次元ベクトルで表現して，**Concept-Occurrence Vector**という特徴ベクトルをシーン一つに割り当てた．これをさらにSVMで6分類して**Coasts, Rivers / Lakes, Forests, Plains, Mountains, Sky / Clouds**のインデックスを割りあてた⁵⁾．類似する手法を三浦らは屋外移動ロボットのカメラ映像を対象として考案し，フレーム内のブロック画素単位でSVMを適用して4種類の物体クラスの判別を行った [F2]．

また，J. JitenらはHMMを用いて，静止画にブロック画素単位でインデックスを割り振る実験を行い，動画中の複数の人物トラッキングにも応用した [C6]．

また，G. Pirirouらは動きの特徴を用いてビデオデータの意味的なクラスタリングを行った [C2]．同様の試みはビデオ検索に関する世界的なコンソーシアムTRECVID [C3]においても精力的に行われている．ここで扱うインデックスにはITS関連の語彙もいくつか含まれている．

2.2 一般物体認識

他方，一般物体認識という観点から異なる複数のタイプの特徴量を統合化し，静止画中に識別対象が存在するかしないかを判別する2クラス識別問題において従来法を凌いだという報告もある [F1]．

しかしながら，以上の研究は比較的限られた内容の映像を対象としており，車載カメラで撮影する走行映像のように撮像条件や構成オブジェクトの変化が激しい景観には適用困難である．一方で，走行映像を対象とした研究では車両や歩行者，白線，標識などに特化した応用特定の画像処理が志向され，IT分野で幅広い相互運用性を有するMPEGとの互換性はほとんど考慮されなかった．これは車を「走るIT環境」と考えれば大いなる技術的損失である．本稿では通信や記録に適したMPEGとの互換性に力点を置き，景観認識手法を提案する．

2.3 MPEG 画像の解析

岩崎ら [B8]はMPEGビデオデータの動きベクトルから直接得られる特徴画像に注目した．この特徴画像は移動領域を記録し，移動物体が停止していると考えられる領域を次フレーム以降も蓄積する．この履歴特性をもつ領域の特徴画像を用いることにより，移動物体が静止した場合やピクチャタイプなどによって動きベクトルが出現しない場合でも，移動物体の検出と追跡が可能であることを示した．

2.4 歩行者認識

歩行者認識の分野では多くのアプローチがあり，大域的にはいくつかのセンサーヒュージョンシステムが存在する [E10], [E11], [E4]. センシングの観点ではレーザレーダ，ミリ波レーダ，超音波デバイス，ビデオカメラ，赤外線カメラ，ステレオカメラなど様々なデバイスが存在する．本稿では単眼の視覚システムにおける画像処理に焦点をあてる．

画像処理に基づく歩行認識では，SVM (Support Vector Machine)と Ada Boost の組み合わせ[E2]，ニューラルネットワークベースの手法[E13]，モデルベース手法，テンプレートマッチング，見え方 (appearance) ベースの手法などが存在する．これらは静止画ベースの手法に分類され，たいていは画像シーケンスの動きや時間的性質を用いない．だが，隠れ或不鮮明または雑音のある信号環境で人間を検出する際は，この動きや時間的性質が極めて重要となる．

画像通信やプローブカーを想定した歩行者認識[A9]なども注目される．また，集合知や学習効率化という点では，誤りを含む学習データから正しいラベルを推定する手法などが画像認識や機械学習の分野でも取り入れられようとしている [H1][H2].

2.5 ヒューマントラッキング

近年，さまざまなヒューマントラッキング手法が開発されている [E5], [E6], [E8]. その中のいくつかは重要なセキュリティ要求に端を発し，それに加えて，ときには行動理解や非定常行動や異常性の検出[E5]，そして人の移動軌跡の解析[E3]にも及ぶ．さらには，群衆 (crowd) の認識も重要となり，基本的な研究が開始された[E9].

他の観点では，誤検出の削減 [A9]や，隠れの問題の解決[E8]，そしてロバストネスの獲得[E6]，などステップバイステップの研究もある．確信度ベースのシーン認識とパターンベースの歩行者検出 (ニューラルネットなど) の統合化手法の提案もなされている．これはパターンベースの認識手法が単独では擬人形状を判別できないことに加え，隠れや形状変化，常識的判断に関しても様々な欠点を有するためである．このような背景に従って，筆者はビデオシーケンスベースの歩行者トラッキングに焦点をあてる．

2.6 構造復元

構造復元に関しては，有名な因子分解法を用いた手法[D1][D2]以外では，A. Saxena *et al.* は MRF を用いて単一の静止画から 3D シーン構造を推論しようとしている [C8]. しかしながら，それらはシーンの 3次元構造の意味解釈にまでは言及していない．

MPEG ベースの認識に関しては，3D 動き抽出，顔トラッキング，シーン認識，ヒューマントラッキングなどの研究がある．それらの研究は主としてメディア記述を志向したのものであり，メディア記述と意味との間の自律的な相互作用には達していない．

2.7 シーンの予測

シーンを物理的事象の観測結果と考える場合，シーンの予測は事象の学習と予測に基づ

いて行われることが多い。たとえばシーン予測の応用を交通安全における危険予知に絞るならば、車両を含む交通事象や人の流れ、人の出現に着目して行うのが妥当と思われる。2003年にタイにおいて歩行者の事故発生に関する時空間分布の調査結果が報告された[12]。これによると時空間の特徴から事故発生を確率的に予測することが可能になる。一方で、地域の実情と過去の多くの事例を踏まえて動物や通学団などの出現を標識で警告することも行われている。これらは認識器を時空間に応じて最適に設定することを可能にする。さらには時系列変化としての交通事象（たとえば車両の動き）そのものを定点観測で獲得し、collapsed Gibbs Samplingによってシステムに学習させ、所定の時空間条件と画像特徴に基づいてシーン予測するという試みもなされている[C19]。同様の認識・予測実験が集団の動きについても適用され、異常行動の検出などが試みられている[C20]。しかしこれらはいずれも定点観測が基盤であり、時々刻々と変化する歩行者事情や道路環境の不連続的变化に適応できるものではない。

一方で、過去の人の流れに関する大規模調査データを用いて人の出現を予測する試みも行われている[A14][A15]。

EPFL（ローザンヌ工科大学）では特定の人の歩行（特に方向決定）に関するモデルを作成し、特定の地域（地下鉄駅前など）を対象として人の流れを認識する実験を行った[E3]。ただし、これは空間や時間の変化、交通事情、周囲の集客施設の影響までは考慮していない。

2.8 アクティブセンシング

KINECTのようにアクティブセンシングによるソリューションもあるが、すべての映像メディアがKINECTのセンシングを具備しているわけではないため、応用対象となる映像が極めて限定される。過去に蓄積された映像の解釈にも適用できない。このことは近年、車載環境認識の実用化において期待されているLIDARについても同様である。LIDARは静的環境の3次元センシングについては高信頼度であるが、人や動物その他非剛体的特徴を示す物体や動的環境に関しては未解決の課題が多い。

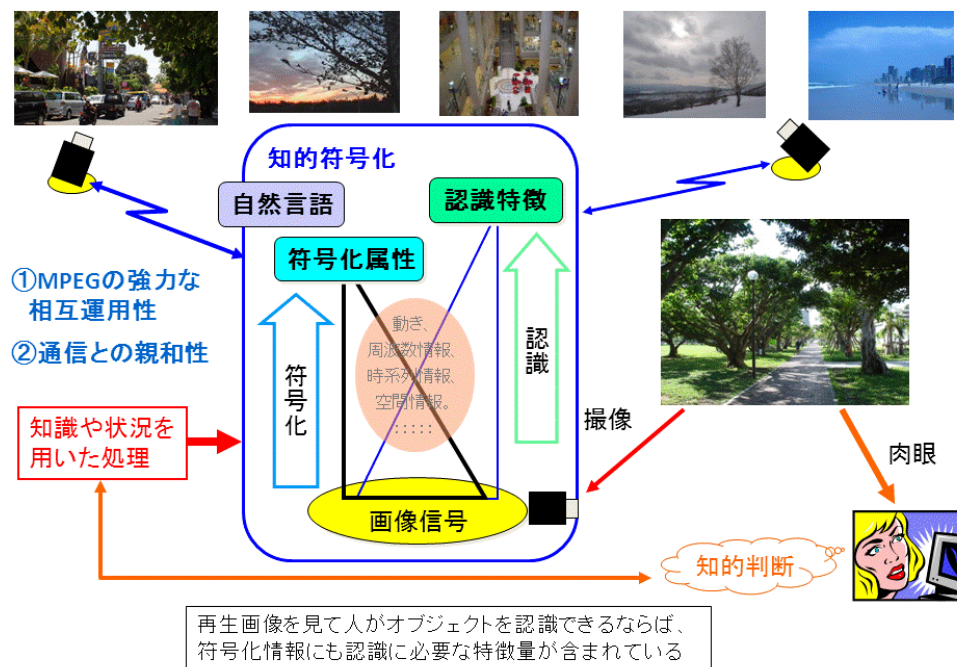
3. 理論

3.1 はじめに

本研究では画像符号化に含まれる特徴抽出と情報量の削減，記号生成のプロセスが極めて人間の認識プロセスに近いことに着目する（図 3.1）. 符号化器と復号化器の双方にオブジェクトに関する知識を共有させ，フレーム間の変化情報をオブジェクトの判別結果やオブジェクトの動情報として符号化する枠組みは，記号生成，自然言語との対応づけなども含めて知的符号化と呼ばれ，画像認識と密接な関係がある分野として研究が行われてきた. これを MPEG ベースで行うことは，MPEG の強力な相互運用性，および通信との親和性から考えて，応用面でのメリットが極めて大きい. さらに，時空間上に分散して存在する複数の認識器（符号化器）においてプローブ協調と，description を用いた知識処理を可能にするため，性能面でもブレークスルーを生み出す可能性がある.

本研究でめざす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
- 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
- 2次元の認識と3次元の認識の統合を可能にする
- シーンに関する知識は description で記述され，カメラを含む観測者の特性は profile



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 3.3 本研究の概観

Fig. 3.3 Targeted overview.

に記述される．いずれもクラス生成によって効率的に管理される

● 認識器の内部および複数の認識器間の学習，推定，予測動作は確信度で統一され，確信度の伝搬により常識やルールの生成，システムの最適化が行われる

一般に人間は，視覚にしても聴覚にしても，対象とする事象のすべてを 100%の精度でセンシングできているわけではない．視野の空間分解能でいえば周辺部は中心部よりも明らかに劣化している．一方，時間方向においても，まばたきや眼球を動かしている場合，あるいは移動体を認知する場合，などは決して一定の分解能ではない．それでも機械に比べて柔軟かつ安定した認知動作を人間が行える理由の一つには生体的な視覚機構そのものの素晴らしさ（これさえも視力における個人差や年齢による違いが存在する）もさることながら，記憶と経験に裏打ちされた外界事象に対する予測能力と推論能力が，時空間的なセンシングの欠落や劣化を補うためと考えられる．

すなわち，人間がシーンを認識する際の予測能力と推論能力に相当する情報処理能力を機械に持たせることができれば，画像処理や特徴抽出が劣化した場合でもシーンを適切に認識することができるはずである．たとえば，歩行者の挙動モデルを知識ベースとするシーンの予測や推定，道路や建物の構造を既知とする場合の移動体の判別などがその典型的な例である．

したがって本研究では，画像情報のソースを MPEG 映像に限定しても，事象の連続性や常識に基づいて不完全情報からでも認識を可能にする知識処理能力を具備することで，システムは従来の画像処理偏重のシーン認識機構と同等あるいはそれ以上の認識率を達成できる，という仮説に立つ．この中には，個々のオブジェクトやシーン全体に関する時間変化特性の獲得，適切な出現予測，非観測部分の推定，などが含まれる．それが MPEG 特徴量の範囲でどの程度まで可能かはまだ現時点では見極めてはいない．だが，少なくとも MPEG 再生画像を人間が見てターゲットオブジェクトを判別できるならば，上記の仮説を裏付けるに十分な特徴量が符号化画像にも存在するはずである．

本章の次節ではまず，MPEG-Based Vision の中の景観認識手法として，有限個の語彙に対応した確信度ベクトルをブロック画素単位で定義し，状況要因と画像特徴を観測量とする多変量線形回帰分析で推定する方法を定式化する．

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，

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

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。

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

ただし、 \mathbf{c}_0 は初期確信度ベクトル（ N 次元列ベクトル）、 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 \mathbf{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 \mathbf{B} は次式で計算できる。

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

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。

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

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

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，

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

$$\mathbf{S}_{color} = (I_r, I_u, I_v)^T \quad (3.21)$$

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

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

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

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 $O_{safe} = \{\text{先}\}$

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.

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

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 次元正方形の逆行列計算は式(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：

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

同様にして、「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：

$$\varphi_{drive} = A_{drive} \mathbf{x} \quad (3.39)$$

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

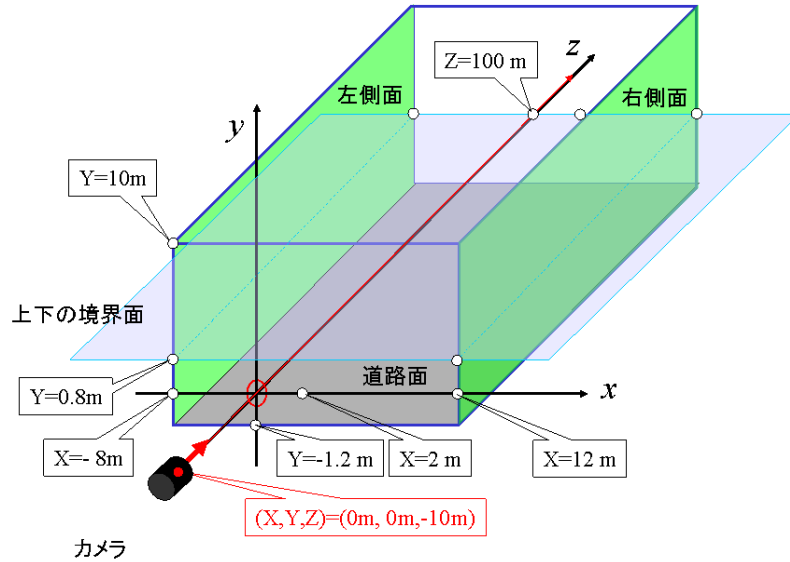


図 3.4 道路シーンの立体構造モデル

Fig. 3.4 Solid structure model of road scene.

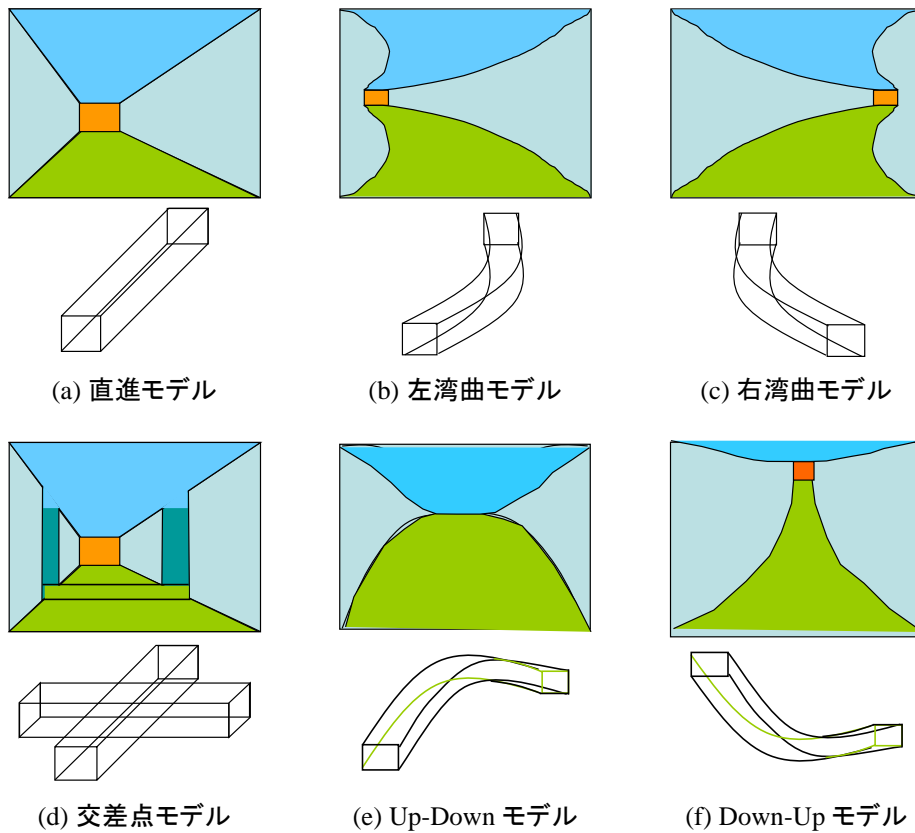


図 3.5 道路シーンの立体構造モデル

Fig. 3.5 Solid structure model of road scene.

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 Φ 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。

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

これを正規化したカーネル行列は次式で与えられる。

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

以下、簡単に証明しておく。

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

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。

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

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。

$$K_{ij}^{in} = \Phi(\mathbf{y}_i) \cdot \Phi(\mathbf{x}_j) \quad (3.44)$$

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、

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

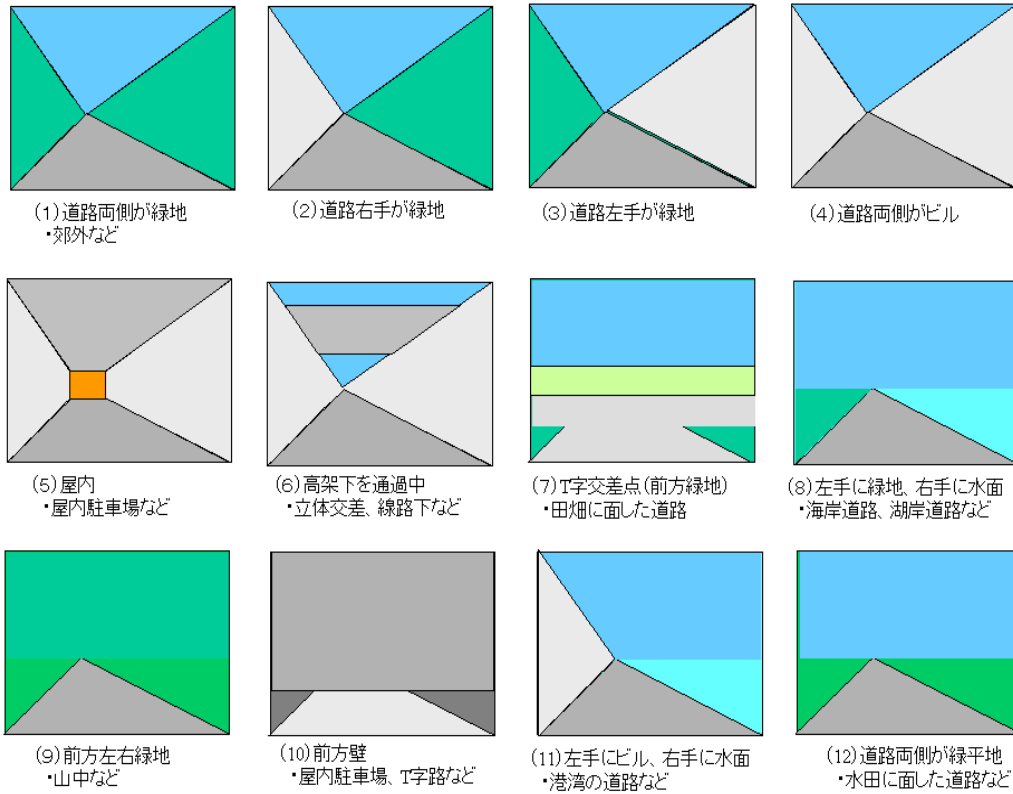


図 3.6 道路シーンの立体構造モデル

Fig. 3.6 Solid structure model of road scene.

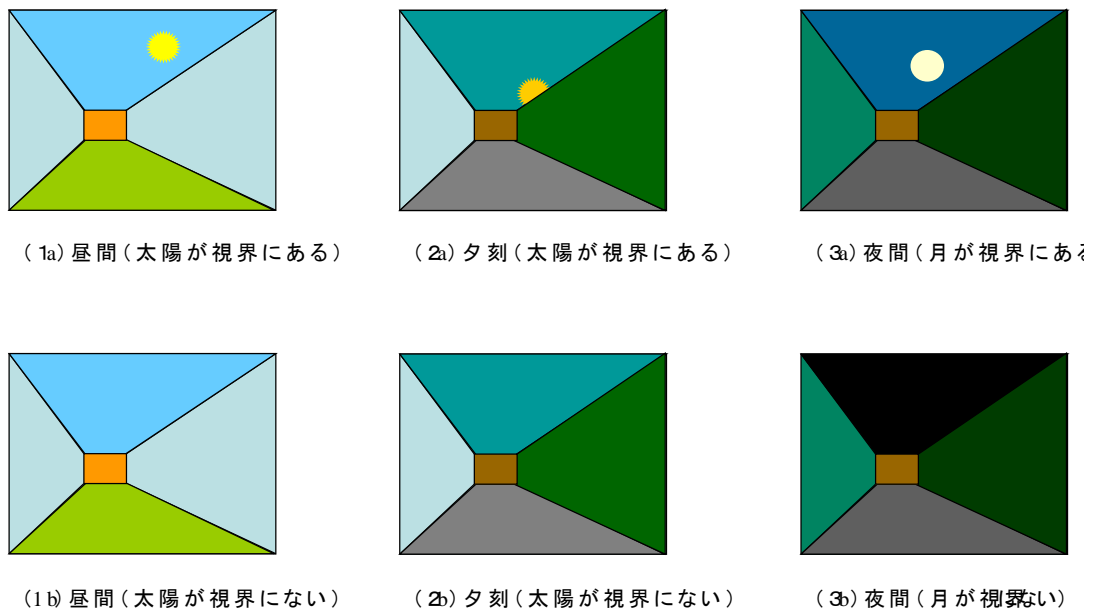


図 3.5 道路シーンの立体構造モデル

Fig. 3.7 Solid structure model of road scene.

となるので、式(3.45)により高次元空間での主軸 $v^{(l)}$ 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]。より長い距離または広い時間範囲を想定すると、箱型のみでは対応が困難となり、実際には少なくとも図 3.3 に示すようなバリエーションが必要となる。

3.13 シーンクラスの生成条件

シーンの分類によってクラスを生成する際には、次のようなことを考慮せねばならない：

- 1) 適用するアプリケーションによって、クラス定義が異なる
- 2) 適用するアプリケーションによって、同じクラス名でも分類規範が異なる
- 3) 検索を行う際、クエリによって検索対象とする対象クラスが異なってくる（すなわち、辞書が変わる）
- 4) シーンのカテゴリ規範はユーザプロファイルの影響を受ける

本章では、車載カメラによる道路交通シーンを分類する際のシーンクラスを例にとり、以下の議論を進める。

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（図 3.5）、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基づい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」

3次元情報と意味の統合表現

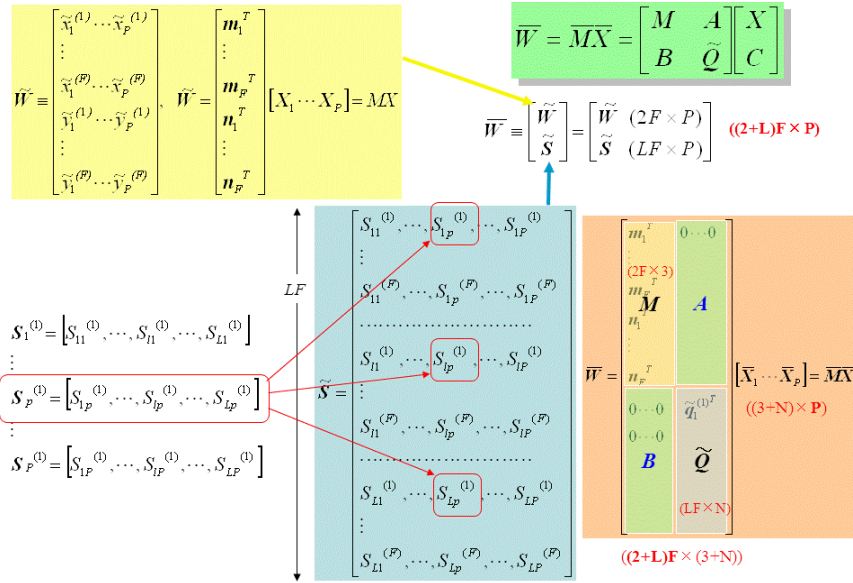


図 3.6 因子分解法の属性拡張

Fig. 3.6 Extension of factorization regarding attributes..

と表現される．ただし，これはシーン S_i 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_i について定義することが可能である．他方，ある S_i 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

(1) 因子分解法の拡張

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 $(2F+3) \times P$ 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． R $(N \times L)$ は回帰係数行列を表す． 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で観測される状態の数を示す．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、ユーザプロフィー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ィー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。したがって、2つの式は次のように統合できる：

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

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 個の \mathbf{S} と F 個の \mathbf{Q} から容易に作成できる(図3.6)。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)には基本的な問題が残されている。1つめの問題はSVDと多変量回帰(MVR)が異なる計算メカニズムであるという点である。SVDは計測行列を2つの行列に分解し、形状とカメラ運動を同時に導出する。その解は計量拘束に基づき一意的に決定される。ところがMVRは基本的には射影であり、 \mathbf{Q} と \mathbf{C} については複数の解が存在するため、SVDの計量拘束に相当する制約をいれない限り一意的には決定されない。

2つめの問題は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 \mathbf{Q} がSVD部分の \mathbf{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、 \mathbf{Q} と \mathbf{C} についてもSVDで同時に獲得できる。

(2) 統計的学習問題としての解法

実際的な解法のひとつは、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、カメラ特性 \mathbf{M} が透視変換特性を持っているとしても、最小2乗法で説くことができる。一方、因子分解法は一般的に正射影を仮定しているため、 \mathbf{M} が透視変換特性の場合にはSVDでは解けない。これは式(3.51)を統計的学習問題として考えていることを意味する。しかしカメラ番号や各カメラの配置に関して自由度が残されている場合、この統計的学習は困難となる。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上に位置しているならば、この統計的学習は可能になる。

(3) 意味と形状の相互作用

知識処理の観点でいうと、式(3.51)における行列 \mathbf{A} と行列 \mathbf{B} は特殊な役割を持っている。 \mathbf{A} は \mathbf{C} から \mathbf{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ーの確信度に影響されるかもしれない。第2に、 \mathbf{B} は \mathbf{X} から \mathbf{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 の行列要素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数量化を担う非線形写像があると考えられる。この仮定に基づけば、式(3.51)における認識の基本構造は行列による線形システムのままでよく、動画像理解システム全体として有する他の非線形機能は主として行列要素と外部変数を関係づける非線形マッピング（確信度からインデックスへの変換およびその逆変換、シーンクラスから状況変数や制御入力への変換、過去の時系列に対する Bayes 的な変換など）、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 を創成する。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で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、通常はこれより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いければ、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。

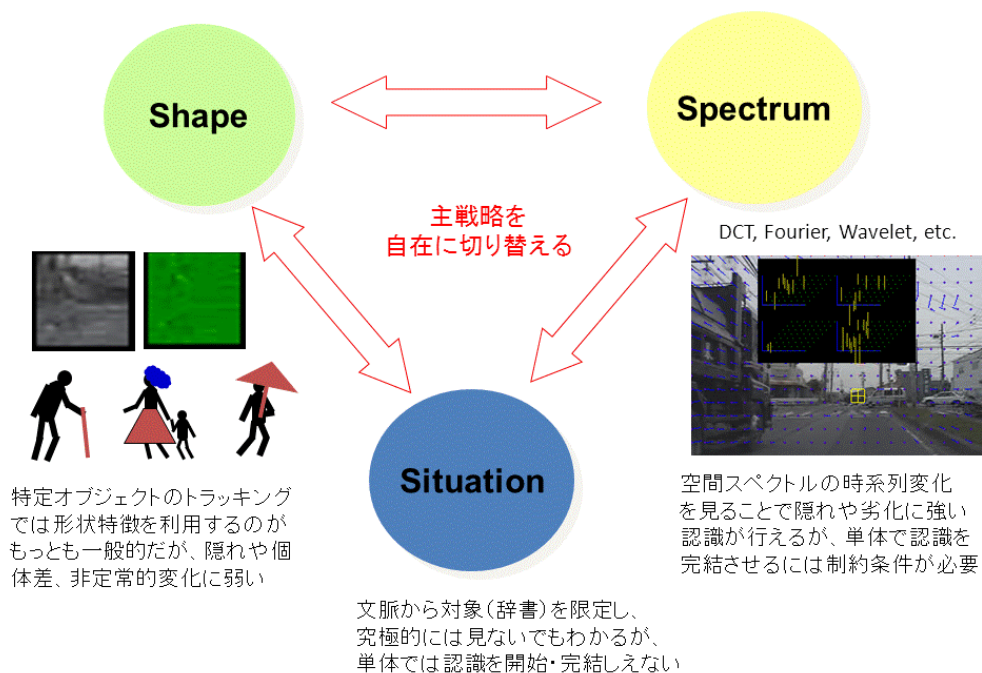


図 3.7 3S モデルの概念

Fig. 3.7 Concept of 3S model.

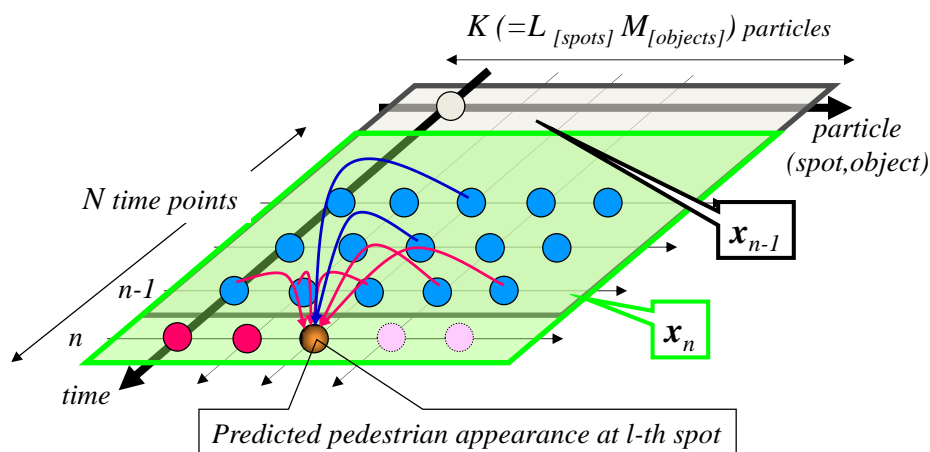


図 3.8 歩行者出現予測のグラフィカルモデル

Fig. 3.8 Graphical model for prediction of pedestrian appearance..

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（図3.7）。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、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 3.8 のように定義すれば， \mathbf{x}_n は次式で表現される．

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

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：

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

ただし $\mathbf{F}, \mathbf{G}, \mathbf{H}$ および \mathbf{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 9.2 の {場所，日時，天候} で構成される．ここで，場所は 3 次元 (各次元が図 9.3 の各区間に相当)，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．9 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。

$$\begin{aligned} \mathbf{S} &= (\mathbf{S}_1, \dots, \mathbf{S}_L)^T = [\mathbf{S}_{place}^T, \mathbf{S}_{time}^T, \mathbf{S}_{env}^T]^T \\ &= [\mathbf{x}_i^T, \mathbf{t}_j^T, \mathbf{s}_k^T]^T \end{aligned} \quad (3.58)$$

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 \mathbf{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 \mathbf{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 \mathbf{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 \mathbf{Q} ($H \times P$) に拡張して、次式で表す。

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

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

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

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 \mathbf{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、 $Y (H \times M)$ は H 個のイベントにおける M 種類のオブジェクト（表 9.2 では $M=12$ ）の出現頻度を表すデータ行列である。 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 $C (N_B \times M_O)$ （この時点では道路の右か左かはわからない）から歩行者（表 9.2 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。

$$Y = \Psi(C) \quad (3.67)$$

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。

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

ただし、 \hat{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、 \hat{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^{\wedge}_{max} と最小値 f^{\wedge}_{min} 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_I 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。

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

ただし、

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

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

である。

車載器、センサー、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。

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

ただし、

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

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， 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．

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．

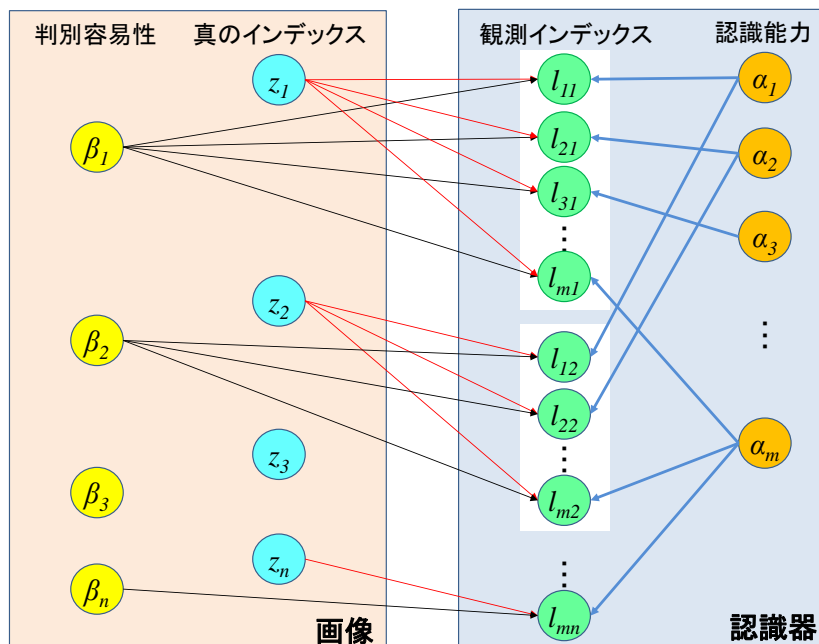


図 3.9 グラフィカルモデル

Fig. 3.9 Graphical Model..

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：

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

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 3.9 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：

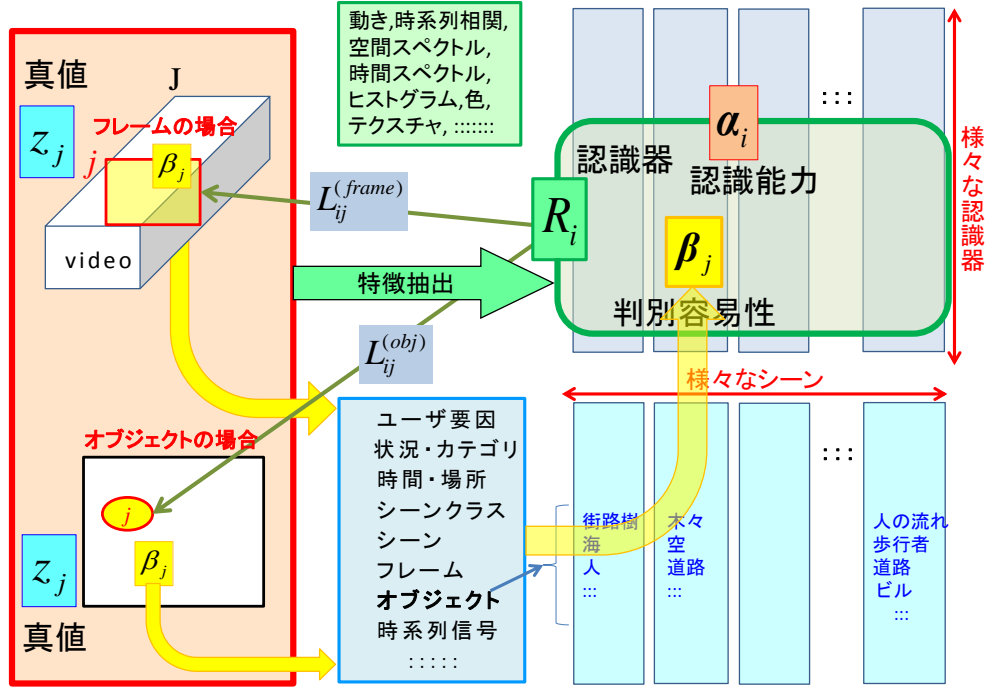


図 3.10 映像と認識器の関係

Fig. 3.10 Relation between images and recognizers.

【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：

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

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：

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

ただし、 $\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\beta}$ 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 $\boldsymbol{\alpha}, \boldsymbol{\beta}$ の値を見つける．

以上、認識対象をオブジェクトインデックスとして説明した．これをさらに多元化した図式が図 3.10 である．ここでは判別対象を {ユーザ要因, 状況・カテゴリ, 時間・場所, シーンクラス, シーン, フレーム, オブジェクト, 時系列信号} に拡張した．

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．