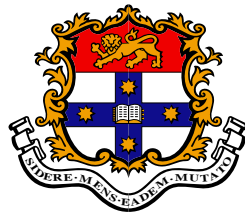# RECOGNISING, REPRESENTING AND MAPPING NATURAL FEATURES IN UNSTRUCTURED ENVIRONMENTS

## Fabio Tozeto Ramos

A thesis submitted in fulfillment
of the requirements for the degree of
Doctor of Philosophy

Australian Centre for Field Robotics
Department of Aerospace, Mechanical and Mechatronic Engineering
The University of Sydney

February, 2007

# Declaration

This thesis is submitted to the University of Sydney in fulfillment of the requirements for the degree of Doctor of Philosophy. This thesis is entirely my own work, and except where otherwise stated, describes my own research.

**Fabio Tozeto Ramos**

February, 2007

Reprinted with corrections and emendations March, 2008

# Abstract

Fabio Tozeto Ramos

The University of Sydney

Doctor of Philosophy

February 2007

## Recognising, Representing and Mapping Natural Features in Unstructured Environments

This thesis addresses the problem of building statistical models for multi-sensor perception in unstructured outdoor environments. The perception problem is divided into three distinct tasks: recognition, representation and association. Recognition is cast as a statistical classification problem where inputs are images or a combination of images and ranging information. Given the complexity and variability of natural environments, this thesis investigates the use of Bayesian statistics and supervised dimensionality reduction to incorporate prior information and fuse sensory data. A compact probabilistic representation of natural objects is essential for many problems in field robotics. This thesis presents techniques for combining non-linear dimensionality reduction with parametric learning through Expectation Maximisation to build general representations of natural features. Once created these models need to be rapidly processed to account for incoming information. To this end, techniques for efficient probabilistic inference are proposed. The robustness of localisation and mapping algorithms is directly related to reliable data association. Conventional algorithms employ only geometric information which can become inconsistent for large trajectories. A new data association algorithm incorporating visual and geometric information is proposed to improve the reliability of this task. The method uses a compact probabilistic representation of objects to fuse visual and geometric information for the association decision.

The main contributions of this thesis are: 1) a stochastic representation of objects through non-linear dimensionality reduction; 2) a landmark recognition system using a visual and ranging sensors; 3) a data association algorithm combining appearance and position properties; 4) a real-time algorithm for detection and segmentation of natural objects from few training images and 5) a real-time place recognition system combining dimensionality reduction and Bayesian learning.

The theoretical contributions of this thesis are demonstrated with a series of experiments in unstructured environments. In particular, the combination of recognition, representation and association algorithms is applied to the Simultaneous Localisation and Mapping problem (SLAM) to close large loops in outdoor trajectories, proving the benefits of the proposed methodology.

# Acknowledgements

The work presented in this thesis would not be possible without the input and support of several people directly or indirectly involved. I am deeply grateful to the people listed below whose contribution to this work is more important than what I can describe in these few lines.

It was a great privilege to be advised by Hugh Durrant-Whyte. I have learnt a lot from the many discussions we had either at the university or while having some pints in a local pub. I thank him for believing in my potential as a researcher, for bringing me to Australia and for supporting me over these 4 years. His insightful vision and thoughtful guidance were essential not only for this thesis but for my career as a researcher.

I am indebted to my examiners whose thoughtful questions and suggestions helped me in improving this thesis in the final version. Paul Newman provided me with an excellent and very detailed review. I am looking forward to spending some time in his group at Oxford. Eduardo Nebot has taught me not only how to barbecue in the traditional Argentinian way but how to keep myself motivated, always looking for new research directions. His thesis report is also very much appreciated. In the last months of my Ph.D I had the pleasure of sharing an office with Dieter Fox. I am very grateful to him for the discussions and suggestions that have originated new ideas and research topics that go beyond the material presented in this thesis.

As a Ph.D student at ACFR I had the chance to expose and discuss my ideas with very knowledgeable researchers who have provided significant inputs into this thesis. Specifically, I would like to mention my three main collaborators. It has been a pleasure working with Suresh Kumar. I am very grateful to him for being my coauthor in several publications and for the innumerous discussions on techniques for dimensionality reduction that constitute part of the contributions in this thesis. Ben Upcroft has been my friend since the begining of my Ph.D. Besides the innumerous schooners of beer we had, he has taught me writing styles, robotic programming and techniques for analysis of experiments. I am very thankful to him for being my coauthor over these years and for organising the statistical learning reading group; a source of ideas for many Ph.D students. Even before coming to Australia,

To Mom and Dad for the support and
encouragement throughout my life

# Contents

# List of Figures

# List of Tables

# Symbols and Notation

Unless otherwise stated, the following notation is used: Matrices are capitalised and vectors are in bold type. A set of vectors is both capitalised and bold type. There is no distinction in notation between probabilities and probability densities.

| Symbol | Meaning |
|---|---|
| $\cdot$ | dot product |
| $\otimes$ | outer product |
| $\succeq$ | used to indicate the constraint that a matrix is positive semi-definite, e.g. $K \succeq 0$ |
| $\lceil x \rceil$ | smallest integer $\geq x$ |
| $\| \cdot \|$ | a particular norm, usually $L_2$ |
| $|K|$ | determinant of matrix $K$ |
| $\nabla \mathbf{f_u}$ | partial derivative of $\mathbf{f}$ w.r.t $\mathbf{u}$ |
| $L^{\#}$ | pseudo-inverse transpose of $L$ |
| $\mathbf{z}^T$ | the transpose of vector $\mathbf{z}$ |
| $\mathbf{x}_{k|k-1}$ | the subscript means estimated at time $k$ given observations up to time $k-1$ |
| $d$ | dimensionality of the output space |
| $d_{r,s}^2$ | square of the distance between points $r$ and $s$ |
| $\text{diag}(W)$ | a vector containing the diagonal elements of matrix $W$ |
| $\mathcal{D}(\pi; \lambda)$ | Dirichlet distribution for $\pi$ parameterized by $\lambda$ |
| $D$ | dimensionality of the training data |
| $\delta_{ij}$ | Kronecker delta function |
| $\text{eig}(W)$ | eigenvector of matrix $W$ |
| $E$ | residual |
| $\eta$ | binary matrix $n \times n$ indicating neighbourhood |

| Symbol | Meaning |
|---|---|
| $\mathcal{F}_M$ | negative free energy |
| $\mathcal{G}$ | graph |
| $h(x)$ | differential entropy of random variable $x$ |
| $\mathbf{h}(\mathbf{x})$, $\mathbf{h_x}$ | observation model |
| $\mathbf{H}$ | centring matrix |
| $I$ | identity matrix or an image |
| $I(x,y)$ | input image |
| $I_C$ | colour space |
| $I_T$ | texture space |
| $\mathcal{I}(x,y)$ | mutual information between random variables $x$ and $y$ $J$ |
| $\log(z)$ | natural logarithm (base $e$) |
| $\mathcal{L}$ | log-likelihood or Lagrangian; clear from the context |
| $m$ | number of components in a mixture model |
| $M$ | probabilistic model |
| $n$ | number of training cases |
| $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \Sigma)$ | the variable $\mathbf{x}$ has a Gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\Sigma$ |
| $O(.)$ | big Oh |
| $p(.)$ | probabilistic function or probability |
| $\mathbf{P}_{k\mid k-1}$ | estimated covariance for the state vector at time $k$ |
| $\Phi(\mathbf{X})$ | embedding cost function for LLE |
| $q_{\boldsymbol{\theta}}(x)$ | free distribution on $x$ |
| $\mathbf{Q}_k$ | estimated covariance of the control inputs noise at time $k$ |
| $\mathbf{R}_k$ | covariance of the observation noise |
| $\mathbb{R}$ | the real numbers |
| $s$ | discrete hidden variable with multinomial distribution |
| $S$ | scatter matrix |
| $\mathbf{S}_k$ | innovation covariance matrix |
| $\mathcal{S}(\mathbf{x}; \rho, \boldsymbol{\Lambda}, \omega)$ | Student-t distribution for $\mathbf{x}$ parameterized by $\rho$, $\Lambda$ and $\omega$ |
| $\mathcal{T}(\mathbf{X})$ | sum of pairwise squared distances of points in $\mathbf{X}$ |
| $\text{Tr}(K)$ | trace of matrix $K$ |
| $\boldsymbol{\theta}$ | set of model parameters |
| $\tau(d^G_{r,s})$ | inner product matrix of $d^G_{r,s}$ |

| Symbol | Meaning |
|---|---|
| $\Upsilon$ | transformed points in semi definite embedding |
| $\mathbf{U}^k$ | control inputs up to time $k$ |
| $\mathbf{W}_k$ | Kalman gain matrix |
| $\mathcal{W}(\Gamma; \alpha, \mathbf{B})$ | Wishart distribution for matrix $\Gamma$ parameterized by $\alpha$ and $\mathbf{B}$ |
| $\mathbf{x}$ | hidden random variable (vector) |
| $\mathbf{x}_k$ | state vector at time $k$ |
| $\mathbf{x}_{v,k}$ | vehicle pose at time $k$ |
| $\mathbf{x}_{m,l}$ | map with landmark positions |
| $\chi^2$ | Chi-squared distribution |
| $\mathbf{Z}$ | $D \times n$ matrix of training inputs $\{\mathbf{z}_i\}_{i=1}^n$ or set of observations |
| $\mathbf{Z}^k$ | observations up to time $k$ |
| $\mathbf{z}_i$ | the $i$th training input or observation |

# Chapter 1

# Introduction

## 1.1  Motivation

This thesis is concerned with the problem of building stochastic models of unstructured environments from sensory information. These models are used to address the main perceptual tasks such as recognition, representation and mapping. The methods and algorithms described are applied to specific robotics problems in aerial, terrestrial and underwater domains. The combination of detection, representation and association results in more reliable robots that can operate robustly in complex natural environments. This chapter motivates the thesis and presents the main problems robots face when operating in unstructured domains.

Over the past ten years mobile robotics research has primarily focused on problems related to navigation such as localisation and map building. The first step in building an autonomous robot is to provide the ability to navigate safely in an unknown environment while keeping an internal estimate of its position with respect to the surrounding world. This must be achieved despite noisy measurements, irregular terrain, dynamic environments, different weather conditions and many other complexities.

When a robot estimates its position with respect to incrementally mapped environmental features, the problem is known as simultaneous localisation and mapping (SLAM). Conventional stochastic solutions to SLAM involve the computation of covariance matrices which are in general of complexity $O(n^2)$ in the number of features (or landmarks) in the map. Much effort has been made to reduce the complexity of SLAM algorithms. The best SLAM algorithms can now deal with many thousands of landmarks (Guivant and Nebot 2001b; Thrun et al. 2002; Paskin 2003; Bosse et al. 2004). However, in many cases, position information alone is not enough to navigate robustly. Association of map features becomes difficult as errors in position estimates increase. Furthermore, extraneous objects might exist in the

environment and their inclusion in the map can have disastrous consequences. For these reasons, reliable detection and association of landmarks plays a major role in autonomous localisation and mapping. When multiple sensors are combined to detect, associate and represent objects and landmarks, this becomes part of the wider *perception* problem.

Humans can detect and create internal representations of thousands of objects that enable them to classify observed objects as belonging to a particular class, despite variations in shape, colour or size. Such knowledge is acquired over many years of learning involving interpretation of sensory information and creation of models that are robust to complexity. Endowing robots with such capability would considerably enhance autonomy and reliability. In mapping and localisation, for example, a robot would be able to map only specific objects in the environment that are known to be reliable landmarks, and identify dynamic objects commonly found in real applications such as humans and cars.

The benefits of reliable perception in unstructured environments are analysed in two problems. In navigation, the robot has to detected and represent landmarks to create consistent maps. With a large number of detected landmarks the problem is how to associate them correctly despite position uncertainty. Furthermore, many tasks require descriptive maps that possess more information than just feature positions. For these tasks compact probabilistic representations can provide higher level models for decision making. These problems are detailed with illustrative examples below.

## 1.2   Navigation in Complex Environments

The robust identification of landmarks for localisation and mapping requires the classification of an object as static or dynamic and the selection of landmarks that are easier to identify and associate. Appearance properties such as shape, colour and texture can extract interesting features that are not apparent in purely geometric models. A major issue is the creation of accurate models that can account for the variability of appearance expected in unstructured environments. Appearance models can be divided into two classes: models for recognition encode a general description of the landmark class and must account for all variability within the class in both shape, colour and texture; models for association encode the information necessary to distinguish one particular landmark from others. Models can be generated as the robot navigates by creating unsupervised representations or learnt from training data. As it is difficult to provide extensive datasets that capture the variability of natural environments, this thesis concentrates on the unsupervised approach to building representations which is addressed in Chapter 4.

Most current outdoor robotics is based on sensors that provide geometric (range and bearing) information. Although these sensors are accurate, they only provide geometric

profiles of objects which are, in general, insufficient for recognition. Conversely, imaging sensors provide richer information such as shape, texture and colour. They are passive, do not consume much power and, in most cases, are less expensive than ranging sensors. The main problem stopping their applicability to outdoor robotics is the difficulty in interpreting the complex information provided.

As opposed to indoor robotics where large patterns such as walls and doors are easily identified with range sensors, outdoor applications are characterised by the lack of geometric structure. Moreover, the existence of far-field objects introduces another problem as landmarks may lie beyond the maximum range of the sensor. This can be seen in Figure 1.1 where a laser scan is plotted in a typical outdoor image. Because most objects are outside the maximum range of the sensor, only 10 (3%) out of 361 readings obtained from the laser scan are valied and useful measurements. These points are separated into three clusters with very similar spatial configuration. While two of these clusters are caused by reflections from trees, the third cluster results from a person. It can be seen that the identification of the person is difficult from only the range graph, but is possible from the associated image. This gives an idea of the importance of imagery in interpreting the world in outdoor robotics.

The main issue in the use of appearance information in unstructured environments is the difficulty in computing models able to encode the complexity of the data. The creation of appearance models for recognition involves learning generative or discriminative models from training data. The lack of structure in the environment imposes many difficulties, of which the need for extensive datasets is the most challenging. To address this issue, Bayesian inference is used in this thesis to show how few training examples can be used to create generative models for recognition and segmentation of natural features.

Although visual information can provide most of the necessary features for recognition, the task can be computed more efficiently when the search for the object is constrained to specific areas in the image. Range readings can be used to reduce the image area where objects are more likely to be found. An algorithm using this idea is presented in Chapter 3. This algorithm extracts shape information from laser readings and fuses it with appearance features to recognise landmarks in a discriminative fashion.

The ability to interpret complex environments is a key challenge for robots. Throughout this thesis algorithms are presented to address this problem combining multiple sensors and creating stochastic representations from training data.

## 1.2.1 Compact Probabilistic Representations

In many applications, a map with only the position of landmarks does not contain enough information for higher-level decision making. Sometimes it is desirable to have additional visual information that includes properties such as colour, shape and texture that can be

Figure 1.1: Typical image and laser scan from an outdoor environment. As most objects are further than the maximum range of the sensor, only 10 distances were obtained. They are insufficient for correct characterisation of the objects as shown in the image.

used to identify objects of interest for particular applications. Examples include rescue, inspection missions and underwater exploration.

Figure 1.2: Metric map augmented with landmark pictures. The additional visual information allows better perception and decision making.

Conventional SLAM algorithms create maps of point features representing the centre of landmarks. Those features are in general detected with range sensors and the centroid is extracted from the profile of the range measurements obtained. Whilst this type of map helps localisation, it does not contain the necessary information to, for example, classify trees according to their specie or to distinguish objects.

To illustrate this problem, Figure 1.2 depicts a metric map with pictures of identified landmarks. Map and pictures were generated by solving the SLAM problem in an urban environment. The inclusion of visual information allows better recognition and association of landmarks. Additionally, it provides information for higher level decision making, for tasks that go beyond navigation. From this example, it can be concluded that the combination of range sensors and imaging sensors for building maps improves the value of SLAM. Accurate range measurements and richer appearance information from cameras are thus complementary.

This thesis proposes a methodology to build compact probabilistic models to encode visual features. The framework is developed with non-linear dimensionality reduction techniques associated with statistical learning. The information from cameras can thus be incorporated into a SLAM framework to yield maps that have at the same time accurate position estimates and higher level feature information, such as landmark images obtained at different viewpoints. The probabilistic visual representation for natural features is described in

Figure 1.3: Gating procedure for data association. The robot with position $\hat{\mathbf{x}}_v$ and uncertainty represented by the ellipse has to associate an observation $\mathbf{z}$ to landmarks previously observed denoted by $\hat{\mathbf{x}}_{l,1}, \hat{\mathbf{x}}_{l,2}, \hat{\mathbf{x}}_{l,3}, \hat{\mathbf{x}}_{l,4}, \hat{\mathbf{x}}_{l,5}$. The gate defined by the dashed ellipse eliminates landmarks $\hat{\mathbf{x}}_{l,3}, \hat{\mathbf{x}}_{l,4}$ and $\hat{\mathbf{x}}_{l,5}$.

Chapter 4 of this thesis.

### 1.2.2   Robust Data Association

A major problem when performing localisation and mapping in large outdoor environments is reliable data association. As the uncertainty over the position of landmarks and vehicle grows, correct association can be very difficult. If for some circumstance an incorrect data association hypothesis is accepted and introduced in the estimation process, the filter may become inconsistent compromising the whole map. The traditional approach for data association is to use a technique known as *gating* (Blackman and Popoli 1999), explained in detail in Chapter 2.

Gating computes a hypothesis test to eliminate associations that are unlikely to be true considering the uncertainty in robot and landmark positions. This is illustrated in Figure 1.3. The uncertainty on robot pose is represented by the ellipse around its expected position denoted by $\hat{\mathbf{x}}_v$. At a particular instant, the robot makes a new observation $\mathbf{z}$. The problem is then to associate $\mathbf{z}$ to some of the landmarks previously detected $(\hat{\mathbf{x}}_{l,1}, \hat{\mathbf{x}}_{l,2}, \hat{\mathbf{x}}_{l,3}, \hat{\mathbf{x}}_{l,4}, \hat{\mathbf{x}}_{l,5})$ or classifying it as a new landmark. The dashed ellipse represent landmarks that are within the gate, and are possible associations for observation $\mathbf{z}$. This eliminates landmarks $\hat{\mathbf{x}}_{l,3}, \hat{\mathbf{x}}_{l,4}$

and $\hat{\mathbf{x}}_{l,5}$ from consideration. Among the remaining landmarks, the decision is to a associate $\mathbf{z}$ to either $\hat{\mathbf{x}}_{l,1}$ or $\hat{\mathbf{x}}_{l,2}$, or to a new landmark. Strategies to address this problem consider the nearest neighbour (NN) as the most likely. When the distance to the nearest neighbour is further than a defined distance, the measurement is considered as coming from a new landmark. Although this methodology is computationally efficient and provides accurate results for small trajectories it clearly fails when position uncertainties are large. Gating and NN become unreliable and not suitable for large scale SLAM problems.

The detection of loop closure in SLAM is a significant issue when a robot follows an extensive path. Observed landmarks are initialised in the map and have to be recognised and correctly associated when the robot re-observes them in a trajectory with loops. When navigating without closing the loop, the uncertainty of the robot location grows making further associations more difficult. Furthermore, the number of landmarks used in large-scale SLAM can be significant, increasing the complexity of the problem. This is illustrated in Figure 1.4 where SLAM is performed with hundreds of features. The uncertainty of the vehicle and feature positions grow making data association difficult.

This thesis address the problem of loop closure and data association in large-scale environments by integrating high-level appearance models into the data association process. A demonstration of these ideas in an extensive unstructured environment is provided in Chapter 6.

## 1.3  Contributions

The main contributions of this thesis are:

1. **Stochastic representation of objects through non-linear dimensionality reduction.** A non-linear statistical model encoding a neighbourhood-preserving dimensionality reduction is proposed as a representation for features in natural environments. This representation is able to distinguish similar objects such as trees and bushes. Efficient inference in this model is formulated and tested with inputs of thousands of dimensions. Inference operations result in mixture of Gaussians that can be integrated in a non-linear filtering scheme.

2. **Landmark recognition system using a laser and camera.** A new algorithm that combines laser and camera information for object detection in outdoor environments is presented. The algorithm is based on a combination of unsupervised and supervised dimensionality reduction methods that fuses information from these two sensors to discover the most discriminative dimensions. Having points mapped to a lower-dimensional space, logistic regression is applied for the final classification. The

Figure 1.4: SLAM over an extensive trajectory. The uncertainty of the vehicle and feature positions increases making data association difficult before closing the loop.

algorithm can process about 5 frames per second and is able to detect objects up to 30 metres from the robot. As opposed to most computer vision algorithms for object recognition that have to search for objects across the whole image, this new approach uses laser to identified regions of interest, and only information in these regions are processed. The resulting real-time algorithm works with high-resolution images to incorporate texture information.

3. **Data association algorithm combining appearance and position properties.** With a probabilistic representation both position and appearance can be used to associate measurements with landmarks using the normal gating technique. The result is a more robust data association algorithm that combines complementary clues. When data association using only position is statistically sufficient, appearance information is not required. When position information is not sufficient, the augmentation of landmark models with appearance information significantly helps in selecting the best hypothesis.

4. **Real-time detection and segmentation of natural objects from few training images.** A fully Bayesian learning methodology using variational calculus is derived for multivariate mixtures. Generative models for object and non-object image patches are learnt from less than 10 images. The algorithm can be used both to segment and classify natural objects and is tested in aerial, underwater and terrestrial domains. Results demonstrate that the algorithm is more stable and accurate than conventional maximum likelihood techniques.

5. **Large-scale deployment of the algorithms for a challenging outdoor simultaneous localisation and mapping problem.** The combination of detection, representation and data association significantly improve results over conventional SLAM algorithms while additionally providing statistical models of landmarks and images acquired at different distances and view points. The framework is tested in an unstructured environment where conventional approaches fail.

6. **Real-time place recognition system tested in indoor and outdoor environments.** The combination of dimensionality reduction and statistical modelling can be applied to the problem of place recognition from images. Generative models from places are learnt from a reduced dataset of images and labels. Images are divided into smaller patches and a classification scheme is proposed where instead of classifying patches individually, uses the whole set of patches in the image to draw its decision. Results are demonstrated in indoor and outdoor experiments.

## 1.4   Thesis Outline

This thesis is organised as follows:

### Autonomous Recognition and Mapping

In Chapter 2, the basic concepts of object recognition, dimensionality reduction and simultaneous localisation and mapping are presented. A review of previous work in these areas is provided and current linear and nonlinear dimensionality reduction techniques detailed. Algorithms for localisation and mapping are explained and recent research on combining vision and range sensors for simultaneous localisation and mapping emphasised.

### Recognition of Natural Features

Algorithms for detection and segmentation of natural features are presented in Chapter 3. The chapter starts by describing a Bayesian framework for learning generative models for

object recognition from images. These models can be trained from few images and additionally used for segmentation. When more training samples or other sensors are available the combination of dimensionality reduction and discriminative learning can provide a faster solution for object detection. Using these ideas, a method for classification from laser and visual information is proposed and tested for recognition of trees in an urban park.

## Stochastic Representation

The problem of representing natural features with probabilistic models is presented in Chapter 4. The probabilistic model described represents a regression function from raw sensor data to a low-dimensional space where essential properties are preserved. The representation has a form of a mixture of linear models with uncertainty encoded by a set of Gaussians. Experiments demonstrate the potential of the approach for intelligent information compression and abstraction.

## Place Recognition

The problem of place recognition from images is explored in Chapter 5. Given a set of images from particular places, the robot has to recognise its location. The problem can be seen as a multi-class classification task. However, rather than learning a single classifier, the approach creates generative models for each place. This has the advantage of being incremental, i.e. if classification is required for additional places, models already learnt can still be used and new data is incorporated to other models.

## Integrating Perception with Mapping

Perceptual models described in previous chapters are combined and applied to the problem of simultaneous localisation and mapping in Chapter 6. Reliable landmark recognition eliminates the problem of navigating in an environment with dynamic objects while landmark appearance representation can significantly improve data association. The combination of position and appearance information for data association is discussed in the chapter, with further outdoor experiments reporting the benefits of the new approach.

## Conclusions and Future Work

Chapter 7 concludes the thesis by analysing the experimental results obtained with the methodology proposed in Chapters 3, 4, 5 and 6. Directions for future work and open issues regarding representation from multi-sensory information and their application to robotics are then discussed.