

Chapter 1

Outline

Campbell (2016)

In this thesis we will examine the use of information and probability theory measures to track the pose, shape, and appearance of non-rigid articulated objects with multiple independently moving cameras. We demonstrate a hierarchical framework that uses the Hilbert-Schmidt independence criterion, to estimate first the shape and appearance, then the pose, then the activity of the dominant actor. At each level different features are employed, however the same kernel function and the same Hilbert-Schmidt norm are used to estimate the independence.

The most elementary level involves a single static camera tracking a single moving non-articulated, rigid object. This requires only that the appearance be tracked. Note that the appearance is defined as the combination of the gray-level intensity, the distance function, and the motion. This example is essentially TLD, which is a very robust and accurate tracker. I have made a contribution to show that my algorithm is an extension to TLD in this case. My tracker models a target as a collection of points. Each point is naturally endowed with x, y coordinates. The features for each point are the intensity at that pixel, the distance to the nearest edge, and the motion of that point from the previous frame. In each frame we find the set of points that maximises the HSIC with the reference (as well as maximising the self similarity and minimising the similarity with the negative training example). It is important to note that the HSIC fails for uniform distributions.

To track the pose we use some sort of optimisation scheme to search through potential poses. Options are: hierarchical search, PSO, LM.

1. Single camera tracking - basically done. 2. Multi camera tracking with known calibration - could be done without a prior labelling of the appearance (i.e just using

the MI scheme from previously and then creating a model on the fly.) 3. unknown calibration - investigate the additional parameters that need to be optimised. 4. activity recognition - constant alpha working (temporal synch + clustering / recognition) 5. DTW - allows unknown alpha, therefore may improve the performance.

1. Introduction

- (a) Mixed reality
- (b) Pose, shape, appearance
- (c) Articulated objects
- (d) Markerless motion capture
- (e) Appearance modelling
- (f) Activity recognition
- (g) Independently moving cameras

2. Probability & Information theory

- (a) Probability spaces
- (b) Random variables
- (c) Distributions
- (d) Moments
- (e) Expectation, variance, covariance
- (f) Entropy & uncertainty
- (g) Mutual information

3. Hilbert Schmidt independence criterion

- (a) Euclidean spaces
- (b) Functional analysis
- (c) Hilbert spaces
- (d) Reproducing property
- (e) RKHS
- (f) Kernels - universality / characteristic
- (g) Covariance Measures

- (h) HSIC measure

4. Appearance tracking

5. Pose estimation

- (a) Bottom-up
- (b) Top-down: Analysis by synthesis
- (c) Calibration free - motion parameter estimation

6. Activity recognition

- (a) Temporal synchronisation
- (b) Dynamic time warping

7. Easter eggs

- (a) Reference to kiwi (bird) #
- (b) Cite Einstein
- (c) Cite Erdos
- (d)

1.1 Information Theory

What is information theory? Where did it start? What questions does it seek to answer? What are the important concepts? – ¿ capacity of a channel to transmit information – ¿ what capacity does the system have to transmit information about data from two images?

The focus of this chapter is on the mathematical concept of information theory, which was first proposed by Claude Shannon in his seminal 1948 paper "A Mathematical Theory of Communication". Central to information theory is the concept of *channel capacity* which describes the ability of a *channel* to transmit information from a source to a destination. In particular, the noisy channel coding theorem specifies an upper bound on the rate at which information can be transmitted through a channel in the presence of noise. Notably, this upper bound depends only on the statistical characteristics of the channel.

1.1.1 Entropy & Uncertainty

How are entropy and uncertainty related? Make the relationship between uncertainty and variance etc quite clear. These links are going to be important later when we talk about HSIC, as it is dependent on the covariance, which is functionally related to the notions of uncertainty and therefore entropy.

1.2 Pose estimation

1.2.1 Bottom up

The principle difference between bottom-up and top-down is the use of the model. In BU estimation we do not know the model of the target. As we track a series of points, however, we can estimate the joint positions by computing the minimum spanning tree on the geodesic distance graph. We could do a graph cut type thing. Initially we want to find which edges are connected to which others. There is an edge between two points if the distances are preserved. We then want to find the joint positions. Similar to Steve's work.

1.2.2 Top down

This is the main approach we are using. We have a known model of the target (i.e all the limb lengths and connections) that we wish to fit to the images(s) the best. For any particular pose we can evaluate the cost, then we just take the best pose. Note that we can think of the model as something that acts to preserve geodesic distances between points, i.e. it lets them move, but only in a way that preserves the shape of the target.

1.2.3 Calibration free

This is currently future work. Track the pose independently in both images. Then we want to find the 3DT pose matrix that best matches the two 2D pose matrices. The principle question is how to propose any given 3DT matrix? probably hard given all the constraints etc.

Chapter 2

Mathematical Preliminaries

The principle contribution of this thesis is to demonstrate the utility of probability and information theory measures for addressing the correspondence problem in certain computer vision applications. In particular we will examine how the *mutual information* (MI) and the *Hilbert Schmidt independence criterion* (HSIC) can be used to evaluate the relationship between two objects. We begin with an introduction to key concepts in probability theory, then demonstrate how these concepts are central to the development of information theory. We conclude with a discussion of the mutual information. In the following chapter we introduce the reader to Hilbert spaces, then demonstrate how Hilbert spaces relate to earlier concepts in probability and information theory.

2.1 Elements of Probability Theory

2.1.1 Probability spaces

A probability space is a triple (Ω, \mathcal{F}, P) where Ω is a sample, or outcome space, \mathcal{F} is the event space, and $P : \mathcal{F} \rightarrow \mathbb{R}$ is a function that maps events to probabilities. The sample space refers to the set of all outcomes that could occur, while the element $A \in \mathcal{F}$ is a subset of the event space such that $A_i \subseteq \Omega$. The function $P : \mathcal{F} \rightarrow \mathbb{R}$ assigns a value to elements of the event space that describe the likelihood of that event occurring. The probability function must satisfy the following properties (2013grayentropy):

- *nonnegative*:

$$P(A) \geq 0, \quad \forall A \in \mathcal{F}; \quad (2.1)$$

- *normalised*:

$$P(\Omega) = 1; \quad (2.2)$$

- *countably additive:*

If $A_i \in \mathcal{F}, i = 1, 2, \dots$ are disjoint, then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i). \quad (2.3)$$

The first condition states that probability values must be greater than or equal to zero, while the second and third conditions collectively state that

In the context of computer vision, for example, the sample space is defined by the application. For a set of images the sample space will change depending on factors such as the width and height of the image, the number of channels used to describe each pixel and the range of values each channel can take. The sample space for the set of 8-bit gray-scale images with a fixed width and height is given by the set of all possible configurations of pixel intensities.

2.1.2 Random variables & Distributions

(2013durretprobability) Given our probability space we define a random variable \mathcal{X} to be a function that maps the sample space to some measure space, often the space of real numbers, i.e. $\mathcal{X} : \Omega \rightarrow \mathbb{R}$. The random variable \mathcal{X} is then said to induce a distribution on \mathbb{R} by letting the probability that $X = a$ be the measure of the set $\{\omega \in \Omega : \mathcal{X}(\omega) = a\}$, which is denoted by $P(\mathcal{X} = a)$. Intuitively, a random variable is an entity that can take on the value of any element of the sample space with probability given by P . Drawing a single sample from a distribution is then akin to selecting a particular value for the random variable, according to the probability space. In this work we will denote by $P(x)$ the probability that \mathcal{X} takes on the particular value x , i.e. that $P(\mathcal{X} = x)$. Given two random variables, \mathcal{X} , and \mathcal{Y} , we can define the joint probability $P(\mathcal{X}, \mathcal{Y})$ to be the probability distribution for the set of N ordered pairs $\{(x_i, y_i) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$.

We may choose our random variable to be a gray scale image, for example. In the absence of any domain information it may be the case that every pixel value is equally likely to occur, and therefore we describe our random variable as having a uniform distribution. In contrast, we may assume that our image is drawn from a set of images of the natural world, so that it may contain a plant or an animal. In this case there is some restriction on the particular pixel values that are likely to occur, and we therefore say that our random variable is governed by some unknown probability distribution function that we may wish to know more about. We may be interested in knowing,

for instance, whether our image contains a kiwi (reference). In our sample space from which we draw images of the natural world there is a probability function that describes the particular configurations of pixels that give rise to images that appear to be kiwi's. Given an image, we may wish to know whether the distribution that generated the current image is equivalent to the true distribution that generates images of kiwi's. This example is obviously far from trivial, as one cannot simply prescribe a function for describing any and all flightless birds from New Zealand. Describing probability functions and their relationship to each other has been the subject of an extraordinary amount of work. In computer science we typically call this problem pattern recognition, which reflects the fact that our goal is to detect patterns and make inferences on the state of the world.

2.1.3 Moments

In order to understand probability distribution functions we need a language to describe them. This is achieved by computing the *moments* of the distribution. For a discrete univariate probability density function, with an expected value μ_1 , the n^{th} moment of the distribution is given by

$$\mu_n = \sum_{i=1} (x_i - \mu_1)^n P(x_i). \quad (2.4)$$

The expectation value, or first moment, of the distribution is given by

$$\mu_1 = \sum_{i=1} x_i P(x_i), \quad (2.5)$$

which is the sum of the individual elements, weighted by their respective probabilities. For a sample of size N the arithmetic mean is given by

$$\bar{x} = \frac{1}{N} \sum_{i=1} x_i. \quad (2.6)$$

The second moment of the distribution, is given by

$$\mu_2 = \sum_{i=1} (x_i - \mu_1)^2 P(x_i), \quad (2.7)$$

which is the expected value of the squared deviation from the mean. For a discrete sample the second moment is defined as the variance and is given by

$$\sigma = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x}). \quad (2.8)$$

2.2 Information Theory

Since Claude Shannon's seminal 1948 paper "A Mathematical Theory of Communication", the field of information theory has developed rapidly. [Give some examples etc].

Much of the early work concerning information theory was intended to be used to develop methods for telephony, and in particular to understand theoretical limits for which a message could be transmitted across a noisy channel. Shannon sought a function that would describe the uncertainty associated with a particular random variable. Shannon recognised that the *entropy*, previously found in statistical mechanics as Boltmann's entropy equation (reference), was a suitable function for relating the probability distribution of a random variable to its uncertainty (Shannon (2001)). For a sample $X = \{x_i\}$ from a random variable \mathcal{X} with corresponding probability distribution $P = \{p_1, p_2, \dots, p_N\}$ the Shannon entropy is

$$H(X) = \sum_i^N p_i \log(p_i). \quad (2.9)$$

Further, Renyi (1961) demonstrated that the Shannon entropy is the limiting case as $\alpha \rightarrow 1$ in a continuous family of functions that satisfy

$$H_\alpha(X) = \frac{1}{1-\alpha} \log \sum_i^N p_i^\alpha. \quad (2.10)$$

The Renyi entropy is defined for $\alpha = 2$ as

$$H_2(X) = -\log \sum_i^N p_i^2. \quad (2.11)$$

The joint entropy of the random variables X , and Y , is given by

$$H(X, Y) = \sum_i^N \sum_k^N p_{i,k} \log(p_{i,k}), \quad (2.12)$$

where $p_{i,k} = P(x_i, y_k)$ and $Y = \{y_i\}$ is a sample from the random variable \mathcal{Y} . One of the key characteristics of the entropy is that it is minimised for events that are either very likely or very unlikely to occur. As a consequence, the entropy is maximised when

we are most uncertain about the outcome of an event. Given a Bernoulli trial (such as flipping a coin) the entropy is maximised when the probability of either outcome is equal, as can be seen in (bernoulli entropy figure).

One of the principle measures to utilise the entropy is the mutual information (reference (shannon?)), defined by

$$I(X, Y) = H(X) + H(Y) - H(X, Y). \quad (2.13)$$

The mutual information is the reduction in uncertainty in X , given observations of Y (and vice versa). Notably, the mutual information can be used to characterise the dependency between X and Y . We can express the mutual information in terms of probabilities as

$$I(X, Y) = \sum_i^N \sum_k^N p_{i,k} \log \left(\frac{p_{i,k}}{p_i p_k} \right), \quad (2.14)$$

which is zero if and only if \mathcal{X} and \mathcal{Y} are independent, i.e. if and only if $p_{i,k} = p_i p_k$.

We can see that the mutual information will be greater when the respective entropies are maximised and when the joint entropy is maximised. Consider, for example, the outcome of two Bernoulli trials, such as flipping two coins T_1 and T_2 in succession. In order for the mutual information to be maximised we wish to reduce our uncertainty about the outcome of T_2 after observing T_1 by as much as possible. This occurs when our initial uncertainty, i.e our uncertainty about the outcome of T_1 , is maximised and our uncertainty about T_2 given T_1 is minimised. We can see that this occurs, for example, when $P(T_1 = heads) = 0.5$ and $P(T_1 = heads, T_2 = heads) = 1$ (with all other joint probabilities equal to zero).

In order to compute the mutual information we need a mechanism to compute the probabilities p_i . It is common to use either histograms or Parzen window techniques to approximate the underlying densities. Given samples $X = \{x_i\}$ and $Y = \{y_i\}$, the histogram is constructed by tallying the frequency of occurrence of each particular value. We then compute the entropy as per 2.11.

We compute the Parzen density estimate of the mutual information according to the framework laid out by Wells, Viola, Atsumi, Nakajima, and Kikinis (1996). The Parzen window technique (Duda, Hart, and Stork (2012)) approximates a probability distribution as a superposition of functions centered on each of the elements of a sample. Any differentiable function that integrates to one can be used in Parzen density estimation. In this work, as in Wells *et al.* (1996), we will use the Gaussian function:

$$G_\psi(x) = \frac{1}{\sqrt{2\pi|\psi|}} \exp\left(\frac{1}{2}x^\top \psi^{-1}x\right), \quad (2.15)$$

where ψ is an empirically chosen covariance matrix. In order to approximate the entropy we create two new samples A and B from X , with $N_A, N_B < N$ their respective sizes. The entropy is approximated according to

$$H(X) \approx H^*(X) = \frac{-1}{N_B} \sum_{x_i \in B} \log \frac{1}{N_A} \sum_{x_j \in A} G_\psi(x_i - x_j), \quad (2.16)$$

whereby an individual probability is approximated according to

$$p(x) \approx p^*(x) = \frac{1}{N_A} \sum_{x_j \in A} G_\psi(x - x_j). \quad (2.17)$$

In his paper, Shannon (2001) states that "The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point".

2.3 Functional Analysis

The principle theoretical framework we are using in this thesis is that of reproducing kernel Hilbert spaces (RKHS). Specifically, we are interested in the Hilbert-Schmidt independence criterion (HSIC), which measures the cross-covariance between two RKHS's associated with two random variables. We begin this section with an introduction to Hilbert spaces and to the more general theory of a function space. We then introduce the reproducing property and demonstrate its utility in evaluating functions from a Hilbert space. Finally, we discuss how the theory of RKHS motivates the use of the HSIC to measure independence.

2.3.1 Hilbert spaces

Most readers will be familiar with the concept of the space of real numbers, \mathcal{R} , which is an example of a vector space. We can characterise this space as a set of dimensions, whereby points in the space take on a single real value from each of the dimensions and for which operations such as addition and scalar multiplication are permitted. Furthermore, we can endow points in the space with certain properties, such as a norm, which typically represents the size, or magnitude of the vector from the origin to the point. Interestingly, we can also represent *functions* as a vector space. Specifically,

a function space \mathcal{F} is defined as a collection of functions which supports the following definitions:

Definition 2.3.1. (from gdurret03) An inner product is a function $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathcal{R}$ that satisfies, for every $f, g \in \mathcal{F}$ and $\alpha \in \mathcal{R}$:

- Symmetric: $\langle f, g \rangle = \langle g, f \rangle$
- Linear: $\langle f, g \rangle = \langle g, f \rangle$
- Symmetric: $\langle f, g \rangle = \langle g, f \rangle$

Definition 2.3.2. (from gdurret03) A norm is a non-negative function $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathcal{R}$ that satisfies, for every $f, g \in \mathcal{F}$ and $\alpha \in \mathcal{R}$:

- Symmetric: $\langle f, g \rangle = \langle g, f \rangle$
- Symmetric: $\langle f, g \rangle = \langle g, f \rangle$
- Symmetric: $\langle f, g \rangle = \langle g, f \rangle$

Note that there exist a rich class of both inner products and norms that characterise many different types of spaces. The Euclidean norm, $\|x\| = \sqrt{\sum x_i^2}$, and the Euclidean dot product $\langle x, x \rangle = xx^T$ are simple examples. It is important to note that dealing with the abstract notion of a norm or inner product between functions is no more trivial in theory than dealing with the norm or inner product between vectors of real numbers (although it may well be in practice).

In this work we consider a function space known as a Hilbert space.

Chapter 3

Pose Estimation

In this chapter we discuss our approach to detecting and tracking the pose of an articulated object through a series of images. We begin with a discussion of our approach and demonstrate its performance in a number of different test cases. We conclude with a discussion of prior work in the fields of object recognition, segmentation, and pose estimation.

3.1 Background

Augmented and virtual reality, which we refer to in this thesis as *mixed reality*, are progressing at a rapid pace. This advancement is due largely to improvements in hardware and in our understanding of the human visual system. Collectively, these advancements mean that mixed reality devices are becoming ergonomic and are capable of delivering content that seamlessly blends with the surrounding world. In particular, this progress means that it is now reasonable to develop mixed reality applications for users interacting in dynamic, unpredictable environments. Consequently, there are two situations for which pose tracking is required. The first is to track the user as she navigates an environment in order to deliver the current pose of the device, necessary for delivering content. The second situation involves tracking both people and objects that occupy the environment surrounding the user. Many mixed reality platform providers (references?) address the first problem (of finding the position of the headset) by tracking the internal dynamics of the headset with embedded inertial sensors. For better performance and for more complex behaviours many providers also require that the user operates within a controlled environment that features a number of cameras around its perimeter. These cameras track the user as they move

around the environment, which allows for more interactivity between the user and the application, and improves the positional tracking system within the application, which improves the display performance. In most cases, however, this visual tracking system relies on distinctive visual markers placed on the user. In addition, few systems include capabilities to track external actors. Applications would become far more immersive if mixed reality platforms were capable of tracking arbitrary objects in their environment.

3.1.1 Terminology

In this thesis we refer to an *actor* as any object, whether it is a person, a car, or a [endemic NZ bird]. We differentiate between the cases where the actor is *rigid* or *non-rigid*, and whether they are *articulated* or not. A human hand, for instance, is a non-rigid articulated object, due to the fact that the fingers can move independently (typically under certain constraints). In contrast, a solid object such as a desk is not articulated (if one ignores any drawers or movable appendages). A human hand is non-rigid because the skin undergoes complex non-linear motions in response to an underlying rigid motion. In this work however we ignore these non-linear effects and treat the hand as a rigid body object. Similarly we treat a walking person as a rigid-body, articulated actor.

From a formal perspective we use the following definition:

Definition 3.1.1. Given an object represented as a graph over a collection of nodes, the object is rigid and non-articulated if the Euclidean distances between the nodes are fixed. A rigid object is articulated if the Euclidean distances between nodes are free to change, while the Geodesic distances on the mesh remain constant.

Actors are represented as a graph $G = (V, E)$. We have that $V = \{v_i\}$ is a set of nodes, $E = \{e_{ij}\}$ is a set of edges and $e_{ij} > 0$ denotes an edge between node i and node j . Our graph is weighted, such that that weight on each edge represents the Euclidean distance between it's two nodes, i.e. $e_{ij} = d(i, j)$. The geodesic distances between any two nodes is the sum of the edges that lie on the shortest path between them. Intuitively this works in the same way as a skeleton. At any point in time separate joints that are not connected through a bone may be any arbitrary Euclidean distance from each other. If we were to find the shortest path between the ends of each of our thumbs, however, we would find that the distance we travel along the bones remains constant at any point in time. A piece of future work, therefore, is to explore whether

we can build an unsupervised model by exploiting this property. This idea will be discussed further in Chapter x.

[probably a good place for a diagram]

In order to reduce the complexity of the problem we assume that any actor we wish to track can be specified by a graph G ahead of time. This dramatically reduces our search space, as many of the nodes in our graph are now constrained to move in a particular fashion.

Given a model of the actor, represented as a graph, how do we represent the pose? There are a number of alternative measures for this, such as exponential maps (), quaternions, and Euler angles. We choose instead to use rotation and translation matrices as these are convenient and we are only dealing with small graphs, so we are not operating under memory constraints.

[Go into more detail about how the model is put together and how the rotation matrices work etc. + homogeneous coordinates.]

We wish to find the current pose of the actor in every frame of the video. We are choosing in this work to rely on top-down approaches, whereby we place a model into the scene and record the degree to which it fits the image. We then adjust the pose of the model in such a way that after a certain number of iterations we can be confident the pose of the model will align with the true underlying actor.

We begin by specifying the true location of the actor in the first frame, and from this build our initial model of the actor. Evaluating the model at any given pose means sampling, for every point in the model, the image intensity at that pixel, the distance to the nearest edge and the optical flow information at that point. The similarity between any current pose of our model and the set of positive (or negative) examples is given by the Hilbert-Schmidt independence criterion, as outlined in section x. Using the HSIC as a measure of the similarity between two random variables is advantageous as it captures all of the functional dependencies between the two samples.

[example of sampling procedure]

References

- Campbell, J. (2016). *How to do Anything*. Dunedin, New Zealand: Non-existent Publishers.
- Duda, R. O., Hart, P. E., and Stork, D. G. (2012). *Pattern classification*. John Wiley & Sons.
- Renyi, A. (1961). On measures of entropy and information.
- Shannon, C. E. (2001). A mathematical theory of communication. *ACM SIGMOBILE Mobile Computing and Communications Review*, 5(1), 3–55.
- Wells, W. M., Viola, P., Atsumi, H., Nakajima, S., and Kikinis, R. (1996). Multi-modal volume registration by maximization of mutual information. *Medical image analysis*, 1(1), 35–51.