

4F12: Computer vision

Summarized from I. Budvytis & S. Albanie lectures, Michaelmas 2021

Oussama Chaib

October 2021

1 Introduction

1.1 Preamble

Threefold goal to the course:

1. Define computer vision
2. Motivation to study computer vision
3. Ways computer vision problems can be formulated so they can be addressed by computer algorithms.

1.2 What is computer vision?

First questions include the definition of computer vision and motivation to study it. Vision is our most powerful sense (more than 50% of the brain's cortex is devoted to vision and related tasks). Discovering from images what is in the scene and so on.

Definition – Computer vision is a collection of algorithms which capture process and interpret images of a scene to extract useful information about the scene.

Examples – Reconstructing objects from a scene (3D shape), discovery of location of image, semantic segmentation of the image (each pixel assigned a label, i.e: person, car, building etc.)

All these examples correspond to three different tasks in computer vision (also known as **3 R's**):

1. Recognition
2. Registration
3. Representation

Computer vision isn't exactly the same as image processing and pattern recognition. Image processing: converting the input to output (more desirable image). Pattern recognition: a certain class of computer vision, require examples.

Why study computer vision? Curiosity, lots of applications of vision in real life.

Some examples: autonomous vehicles, automation, human-computer interaction, augmented reality, security, medical imaging, 3D modelling and measurements.

Starting from biology – The eye is, in many ways, an optical camera, while the brain would be similar to a collection of computer algorithms (complex human neural network where each part of the brain is in charge of a certain thing). An electric impulse is created at the retina,

sent to the brain neural network, then a signal is sent to the spinal cord/body so the body can react.

Why don't we directly copy the brain? It's quite complex, we don't know everything about it, might be hard to replicate and costly as well. It took a certain level of constraints for the brain to evolve to what it is today. Direct copying of the brain may also not be very revealing.

Solution: Instead of directly copying the brain, we could get inspired by it! Deep learning methods for example.

Building a computer vision algorithm – Several parts:

- **Camera :** Convert light to electric signals, convert the encoded signal into a 3D array (with different channels for colors for example, i.e: RGB). Image formation is a **many-to-one** mapping. Inverse imaging is quite challenging since the images only tell us along which ray a feature lies and not how far along the ray (called image depth).

These ambiguities are typically used in art to create optical illusions (i.e. Ames room with slanted wall).

Images are unstructured 2D arrays. Before using them, they should be processed to reduce the amount of information so it is directly relevant to solving the problem (these are called image features, i.e: edges, blobs...).

Image feature – Piece of info about a certain region of an image that contains certain desirable properties. We decide which features to extract based on what problem we would like to solve.

How do we reduce ambiguities? Drawing a set of constraints. i.e: one solution is to use multiple images instead of one. We can introduce assumptions about the world in the imaging scene (i.e: a face contains two eyes and a nose and a mouth). Deep learning is interesting in that sense because it will help determine those assumptions/constraints from labeled and unlabeled data.

- **Feature extraction:** What do they look like? How can they be obtained? Example: edges: discontinuities in images (filtering by smoothing kernel, then looking for largest gradient magnitudes), corners: areas in image with large changes in local curvature.
- **Perspective projection:** Establishing a camera model. Accounts for the position of the camera, perspective projection, position of CCD array on image plane. How 3D coordinates of a point in a scene is matched to a pixel coordinate in an image? All of this is studied within the framework of **perspective geometry**. Lines from planes parallel to image plan all intersect in vanishing points.

Most algorithms tasked with extraction of 3D information from images make very few assumptions about the scene or 3D image. Some types:

- **Shape from texture:** Simple assumptions: either homogenous (brick wall) or isotropic textures. Can infer orientation of surfaces from analyzing how the texture statistics vary over the image.
- **Stereo vision:** Assuming the scene is viewed from two cameras. The constraints are in this case relaxed. Features should be matched using two different images (known as **correspondance problem**) as long as the two cameras are calibrated. Otherwise, we can only recover 3D features using **projective ambiguity** (to be seen later). The 3D structure can be reconstructed to scale (to be seen later too).

- **Structure from motion:** Similar to stereo vision, but instead of using two pictures we allow the camera to move and collect several images/collect a sequence from different viewpoints.
- **Shape from contour:** Instead of tracking 2D points (structure from motion), apparent contours of the the image are tracked. Not covered in this class.
- **Shape from shading:** Unlike the two previous methods, doesn't require camera motion. Uses knowledge from how light reflects in scenes/objects. Also not covered.

1.3 Geometrical and statistical frameworks

Algorithms introduced in previous section belong to a certain **geometrical framework**. Three key steps in geometric frameworks:

1. Reduce the information content of an object by extracting useful features (edges, corners, blobs).
2. Model the imaging process (usually using perspective projection) and express using projective transformations).
3. Invert the transformation using as many images and constraints as possible to extract 3D structure and motion.

Not all problems can be studied using geometrical framework (example: determining the family of cats in an image). In this case, **statistical frameworks** perform better. They often involve some learning process from a collective set of images with the ability to estimate the confidence of your model.

We will look closely at one type of deep learning architecture: **convoluted neural networks (CNN)**. They have multiple layers of feature responses which are obtained by filtering/convolution and non-linear activation functions. The weights of each filter are learned by training the algorithm (millions or billions of parameters!). They are very effective in many computer vision tasks and need large data sample.