

## Lecture 12-13: Feature Detection and Tracking

Lecturer: Luca Carlone

Scribes: -

**Disclaimer:** *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor(s).*

This lecture discusses how to:

- extract point features in an image,
- track features across two images, and
- perform descriptor-based feature matching (we mainly provide pointers).

The presentation mostly follows Chapter 4 in [4].

## -1.1 Corner detection

Consider a pixel  $\bar{\mathbf{x}} = [u, v]^T$  in an image and call  $W(\bar{\mathbf{x}})$  a rectangular window of given size centered at  $\bar{\mathbf{x}}$ . Moreover, call  $\mathcal{I}(\bar{\mathbf{x}})$  the intensity at pixel  $\bar{\mathbf{x}}$ . Then,  $\bar{\mathbf{x}}$  is a good “corner”, if shifting the window  $W(\bar{\mathbf{x}})$  in any direction  $\boldsymbol{\delta} = [\delta_u \ \delta_v]^T$  produces a new window  $W(\bar{\mathbf{x}} + \boldsymbol{\delta})$  which is different from  $W(\bar{\mathbf{x}})$ . Mathematically, a good corner is such that:

$$\min_{\|\boldsymbol{\delta}\|=\epsilon} \sum_{\mathbf{x} \in W(\bar{\mathbf{x}})} \|\mathcal{I}(\mathbf{x} + \boldsymbol{\delta}) - \mathcal{I}(\mathbf{x})\|^2 > 0 \quad (\text{and possibly as large as possible}) \quad (-1.1)$$

since if the minimum is zero, then there is a direction  $\boldsymbol{\delta} \neq 0$  (we constrain  $\|\boldsymbol{\delta}\| = \epsilon$ , for some small scalar  $\epsilon > 0$ ) that does not cause any change of appearance.

For a small displacement  $\boldsymbol{\delta}$ , we can take a first-order Taylor expansion of the cost and write:

$$\min_{\|\boldsymbol{\delta}\|=\epsilon} \sum_{\mathbf{x} \in W(\bar{\mathbf{x}})} \|\mathcal{I}(\mathbf{x}) + \nabla \mathcal{I}(\mathbf{x})^T \boldsymbol{\delta} - \mathcal{I}(\mathbf{x})\|^2 = \min_{\|\boldsymbol{\delta}\|=\epsilon} \sum_{\mathbf{x} \in W(\bar{\mathbf{x}})} \|\nabla \mathcal{I}(\mathbf{x})^T \boldsymbol{\delta}\|^2 = \quad (-1.2)$$

$$\min_{\|\boldsymbol{\delta}\|=\epsilon} \sum_{\mathbf{x} \in W(\bar{\mathbf{x}})} \boldsymbol{\delta}^T (\nabla \mathcal{I}(\mathbf{x}) \nabla \mathcal{I}(\mathbf{x})^T) \boldsymbol{\delta} = \min_{\|\boldsymbol{\delta}\|=\epsilon} \boldsymbol{\delta}^T \left( \sum_{\mathbf{x} \in W(\bar{\mathbf{x}})} \nabla \mathcal{I}(\mathbf{x}) \nabla \mathcal{I}(\mathbf{x})^T \right) \boldsymbol{\delta} \quad (-1.3)$$

where  $\nabla \mathcal{I}(\mathbf{x}) \in \mathbb{R}^{2 \times 1}$  is the gradient of the intensity function  $\mathcal{I}(\cdot)$  at  $\mathbf{x}$  (see also paragraph about “Image Gradients” below for practical considerations). Let us now define:

$$\mathbf{G} = \sum_{\mathbf{x} \in W(\bar{\mathbf{x}})} \nabla \mathcal{I}(\mathbf{x}) \nabla \mathcal{I}(\mathbf{x})^T \in \mathbb{R}^{2 \times 2} \quad (-1.4)$$

Then, since  $\boldsymbol{\delta}$  is a vector of magnitude  $\epsilon$ , we write it as  $\boldsymbol{\delta} = \epsilon \mathbf{d}$  where  $\mathbf{d}$  is a unit-norm vector. Substituting  $\boldsymbol{\delta} = \epsilon \mathbf{d}$  and the definition of  $\mathbf{G}$  in (-1.2) we get:

$$\min_{\|\boldsymbol{\delta}\|=\epsilon} \boldsymbol{\delta}^T \mathbf{G} \boldsymbol{\delta} = \min_{\|\mathbf{d}\|=1} \epsilon^2 \mathbf{d}^T \mathbf{G} \mathbf{d} = \epsilon^2 \lambda_{\min}(\mathbf{G}) \quad (-1.5)$$

where  $\lambda_{\min}(\mathbf{G})$  is the smallest eigenvalue of  $\mathbf{G}$  (if you do not understand how we get the smallest eigenvalue from the minimization, read this: [https://en.wikipedia.org/wiki/Rayleigh\\_quotient](https://en.wikipedia.org/wiki/Rayleigh_quotient)).

It is clear that  $\min_{\delta} \delta^T \mathbf{G} \delta$  is nonzero (for a nonzero  $\delta$ ) if and only if  $\mathbf{G}$  is positive definite. Moreover, the larger the eigenvalues of  $\mathbf{G}$ , the larger the minimum. Therefore, a pixel  $\bar{\mathbf{x}}$  is a good corner, if the corresponding matrix  $\mathbf{G}$  exhibits large eigenvalues.

Different corner detectors use different scores to quantify how “large” the eigenvalues of  $\mathbf{G}$  are. The popular *Harris corner detector* uses the following quantity (for a given small scalar  $k$ )

$$C(\mathbf{G}) = \det(\mathbf{G}) - k \operatorname{tr}(\mathbf{G})^2 \quad (\text{Harris corner score}) \quad (-1.6)$$

to score the quality (or “cornerness”) of a pixel [2].

Another popular choice of the cornerness score is simply:

$$S(\mathbf{G}) = \lambda_{\min}(\mathbf{G}) \quad (\text{Shi-Tomasi score}) \quad (-1.7)$$

and has been proposed by Shi and Tomasi in [5].

**An intuitive explanation of the Harris corner score** Calling  $\lambda_1$  and  $\lambda_2$  the two eigenvalues of  $\mathbf{G}$  and recalling that  $\det(\mathbf{G}) = \lambda_1 \lambda_2$  and  $\operatorname{tr}(\mathbf{G}) = \lambda_1 + \lambda_2$ , the Harris corner score can be rewritten as:

$$C(\mathbf{G}) = \lambda_1 \lambda_2 - k \lambda_1^2 - k \lambda_2^2 - 2k \lambda_1 \lambda_2 = (1 - 2k) \lambda_1 \lambda_2 - k(\lambda_1^2 + \lambda_2^2) \quad (-1.8)$$

Since  $k$  is typically small, it is clear that when  $\lambda_1, \lambda_2$  are small, the corner score will be small, consistently with what we expect.

Now consider the case in which  $\lambda_1 \gg \lambda_2$ . In this case we can write:

$$C(\mathbf{G}) = (1 - 2k) \lambda_1 \lambda_2 - k(\lambda_1^2 + \lambda_2^2) \approx -k \lambda_1^2 \quad (-1.9)$$

and the score becomes negative. This is the case in which we have an *edge*, rather than a corner.

Finally, consider the case in which  $\lambda_1 \approx \lambda_2 \gg 0$ :

$$C(\mathbf{G}) = (1 - 2k) \lambda_1 \lambda_2 - k(\lambda_1^2 + \lambda_2^2) \approx (1 - 2k) \lambda_1^2 - 2k \lambda_1^2 = (1 - 4k) \lambda_1^2 \quad (-1.10)$$

which is large, since we assumed  $\lambda_1 \gg 0$  and small  $k$ . Therefore, the score rewards large eigenvalues of  $\mathbf{G}$  when they have similar magnitude.

**Image gradients** When taking the Taylor expansion in (-1.2) we treated the intensity as a continuous function over the pixel domain, and adopted the standard definition of gradient (that we assume to be a column vector):

$$\nabla \mathcal{I}(\mathbf{x}) = \nabla \mathcal{I}(u, v) = \begin{bmatrix} \frac{\partial \mathcal{I}(u, v)}{\partial u} \\ \frac{\partial \mathcal{I}(u, v)}{\partial v} \end{bmatrix} \quad (-1.11)$$

However, in practice the image is formed by discrete pixels, and the notion of derivative needs to be replaced by a discrete counterpart. A common choice is to use *finite differences*:

$$\nabla \mathcal{I}(\mathbf{x}) = \nabla \mathcal{I}(u, v) \approx \begin{bmatrix} \frac{\mathcal{I}(u+h, v) - \mathcal{I}(u, v)}{h} \\ \frac{\mathcal{I}(u, v+h) - \mathcal{I}(u, v)}{h} \end{bmatrix} \quad (-1.12)$$

If you want to learn more about alternative choices, please take a look at [4, p. 99].

## -1.2 Feature Tracking

Assume we extracted  $N$  corner features from an image  $\mathcal{I}_1$ . Consider one of such corner features at pixel position  $\mathbf{x}_1$  and call  $\mathbf{p}$  the 3D point which projects to  $\mathbf{x}_1$ . The question we answer in this section is: *given a second (temporally continuous) image  $\mathcal{I}_2$ , compute the pixel projection of  $\mathbf{p}$  in  $\mathcal{I}_2$ , namely  $\mathbf{x}_2$* . Note that in this problem we are only given the images  $\mathcal{I}_1$  and  $\mathcal{I}_2$  and the pixel  $\mathbf{x}_1$  so we cannot use standard perspective projection to compute  $\mathbf{x}_2$  (e.g., the coordinates of  $\mathbf{p}$ , and the intrinsic and extrinsic parameters of the cameras picturing  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are unknown).

Another way to look at the problem is to think that the camera is moving between the capture of image  $\mathcal{I}_1$  and  $\mathcal{I}_2$ , and we want to compute the corresponding displacement of the pixel  $\mathbf{x}_1$  induced by the camera motion. Therefore, computing  $\mathbf{x}_2$  is the same as tracking the motion  $\boldsymbol{\delta} = [\delta_u, \delta_v]^T$  of pixel  $\mathbf{x}_1$  between image  $\mathcal{I}_1$  and  $\mathcal{I}_2$ :

$$\mathbf{x}_2 = \mathbf{x}_1 + \boldsymbol{\delta} \quad (-1.13)$$

Since we do not have any geometric information, our only hope is to compute the motion  $\boldsymbol{\delta}$  by comparing pixel intensities in the two images. More specifically, if we build a rectangular window  $W(\mathbf{x}_1)$  collecting a set of pixels around pixel  $\mathbf{x}_1$  in  $\mathcal{I}_1$ , we can try to compute  $\mathbf{x}_2$  by looking for a patch in  $\mathcal{I}_2$  whose appearance is similar to  $W(\mathbf{x}_1)$ .<sup>1</sup> Reasoning in terms of pixel displacement  $\boldsymbol{\delta}$ , this is equivalent to computing:

$$\min_{\boldsymbol{\delta}} \sum_{\mathbf{y} \in W(\mathbf{x}_1)} \|\mathcal{I}_1(\mathbf{y}) - \mathcal{I}_2(\mathbf{y} + \boldsymbol{\delta})\|^2 \quad (-1.14)$$

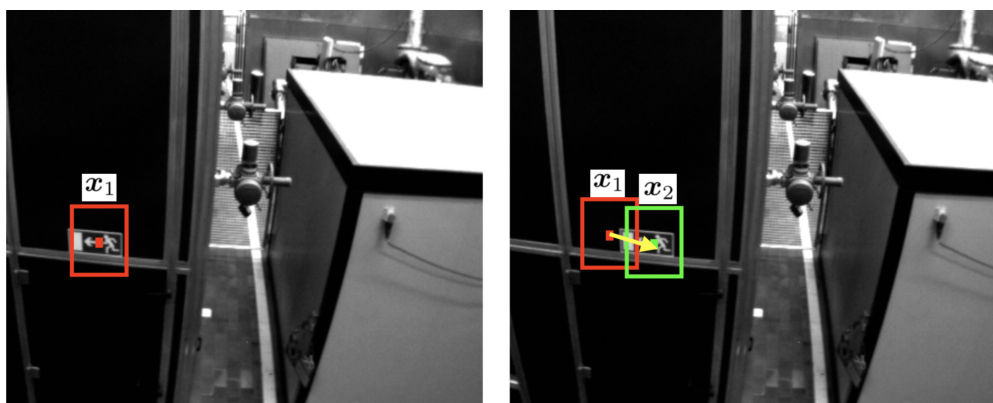


Figure -1.1: The goal of feature tracking is to compute the displacement of a given pixel  $\mathbf{x}_1$  (displacement shown as yellow arrow) between two images, due to the camera motion.

However, the minimization (-1.14) already took a strong assumption: every pixel in  $W(\mathbf{x}_1)$  is moving by the same amount  $\boldsymbol{\delta}$ . Unfortunately, this assumption, typically referred to as a *translational motion model* [4, Section 4.2.1], is only true if  $W(\mathbf{x}_1)$  is picturing a flat surface parallel to the image and the camera is moving parallel to it. This justifies the use of a more general deformation model, which assumes a more complex mapping of the pixels:

$$\min_{\mathbf{A}, \boldsymbol{\delta}} \sum_{\mathbf{y} \in W(\mathbf{x}_1)} \|\mathcal{I}_1(\mathbf{y}) - \mathcal{I}_2(\mathbf{A}\mathbf{y} + \boldsymbol{\delta})\|^2 \quad (-1.15)$$

<sup>1</sup>Note: this assumes that the illumination changes due to the motion of the camera are small, an assumption known as “brightness constancy” constraint.

which assumes an affine (pixel) motion model for the pixels in the window  $W(\mathbf{x}_1)$ . This model is a good approximation for small planar patches parallel to the image plane and undergoing small rotations [4, Section 4.2.1]. In order to compute the optimal displacement  $\boldsymbol{\delta}$ , we can take a Taylor expansion around the point  $\mathbf{A} = \mathbf{I}_2$  and  $\boldsymbol{\delta} = \mathbf{0}_2$  of the intensity  $\mathcal{I}_2(\mathbf{A}\mathbf{y} + \boldsymbol{\delta})$  (thought as a function over a continuous pixel space):

$$\mathcal{I}_2(\mathbf{A}\mathbf{y} + \boldsymbol{\delta}) \approx \mathcal{I}_2(\mathbf{y}) + \nabla\mathcal{I}_2(\mathbf{y})^\top[(\mathbf{A} - \mathbf{I}_2)\mathbf{y} + \boldsymbol{\delta}] \quad (-1.16)$$

which after substituting back into (-1.15) leads to a linear least squares model:

$$\min_{\mathbf{A}, \boldsymbol{\delta}} \sum_{\mathbf{y} \in W(\mathbf{x}_1)} \|\mathcal{I}_1(\mathbf{y}) - \mathcal{I}_2(\mathbf{y}) - \nabla\mathcal{I}_2(\mathbf{y})^\top[(\mathbf{A} - \mathbf{I}_2)\mathbf{y} + \boldsymbol{\delta}]\|^2 \quad (-1.17)$$

from which one can compute  $\mathbf{A}$  and  $\boldsymbol{\delta}$  in closed form.

In order to get some more insight, let us go back to the purely translational deformation model, whose Taylor expansion gives:

$$\mathcal{I}_2(\mathbf{y} + \boldsymbol{\delta}) \approx \mathcal{I}_2(\mathbf{y}) + \nabla\mathcal{I}_2(\mathbf{y})^\top\boldsymbol{\delta} \quad (-1.18)$$

After substituting (-1.18) back into (-1.15) we obtain:

$$\min_{\boldsymbol{\delta}} \sum_{\mathbf{y} \in W(\mathbf{x}_1)} \|\mathcal{I}_1(\mathbf{y}) - \mathcal{I}_2(\mathbf{y}) - \nabla\mathcal{I}_2(\mathbf{y})^\top\boldsymbol{\delta}\|^2 \quad (-1.19)$$

whose solution can be computed using linear least squares as:

$$\boldsymbol{\delta}^* = \left( \sum_{\mathbf{y} \in W(\mathbf{x}_1)} \nabla\mathcal{I}_2(\mathbf{y})\nabla\mathcal{I}_2(\mathbf{y})^\top \right)^{-1} \sum_{\mathbf{y} \in W(\mathbf{x}_1)} \nabla\mathcal{I}_2(\mathbf{y})(\mathcal{I}_1(\mathbf{y}) - \mathcal{I}_2(\mathbf{y})) \quad (-1.20)$$

Note that the matrix to be inverted is the same as the  $\mathbf{G}$  we computed to extract good corners in the previous section. Therefore, by construction, good corners are such that the matrix in (-1.20) is invertible.

### -1.3 Descriptor-based Feature Matching

Feature tracking typically works well if the images are not very different (e.g., collected by the same camera at consecutive time instants). However, feature tracking does not work well when the images are taken by radically different viewpoints. In these case, a *Descriptor-based Feature Matching* is typically preferred, where for each keypoint, one computes a *descriptor* vector. The descriptor can be understood as a unique “signature” of the feature.

Therefore, *descriptor-based feature matching* works as follows:

- given images  $\mathcal{I}_1$  and  $\mathcal{I}_2$ , it computes corners in both images as well as the corresponding descriptors for each corner;
- it matches corners in the two images if the corresponding descriptors are close enough (note: each descriptor is simply a large vector)

Many of these descriptors also come with a *feature detector*, which is typically a more sophisticated version of the Harris corners we saw earlier in this document. However, it is also common to attach detectors to Harris or Shi-Tomasi corners.

Common detector-descriptors are:

- SIFT [3]: [https://docs.opencv.org/3.4.3/da/df5/tutorial\\_py\\_sift\\_intro.html](https://docs.opencv.org/3.4.3/da/df5/tutorial_py_sift_intro.html) (this nice video provides an in-depth explanation: [https://www.youtube.com/watch?time\\_continue=3964&v=NPcMS49V5hg](https://www.youtube.com/watch?time_continue=3964&v=NPcMS49V5hg))
- SURF [1]: [https://docs.opencv.org/3.0-beta/doc/py\\_tutorials/py\\_feature2d/py\\_surf\\_intro/py\\_surf\\_intro.html](https://docs.opencv.org/3.0-beta/doc/py_tutorials/py_feature2d/py_surf_intro/py_surf_intro.html)
- ORB: [https://docs.opencv.org/3.0-beta/doc/py\\_tutorials/py\\_feature2d/py\\_orb/py\\_orb.html](https://docs.opencv.org/3.0-beta/doc/py_tutorials/py_feature2d/py_orb/py_orb.html)

All these (and more) are available in OpenCV, the most popular library for computer vision.

## References

- [1] H. Bay, T. Tuytelaars, and L. V. Gool. Surf: speeded up robust features. In *European Conf. on Computer Vision (ECCV)*, 2006.
- [2] C. Harris and M. Stephens. A combined corner and edge detector. *Proceedings of the 4th Alvey Vision Conference*, pages 147–151, August 1988.
- [3] D.G. Lowe. Distinctive image features from scale-invariant keypoints. *Intl. J. of Computer Vision*, 60(2):91–110, 2004.
- [4] Y. Ma, S. Soatto, J. Kosecka, and S.S. Sastry. *An Invitation to 3-D Vision*. Springer, 2004.
- [5] J. Shi and C. Tomasi. Good features to track. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 593–600, 1994.