

ECE276A: Sensing & Estimation in Robotics

Lecture 11: Visual Features and Optical Flow

Instructor:

Nikolay Atanasov: natanasov@ucsd.edu

Teaching Assistants:

Qiaojun Feng: qif007@eng.ucsd.edu

Tianyu Wang: tiw161@eng.ucsd.edu

Ibrahim Akbar: iakbar@eng.ucsd.edu

You-Yi Jau: yjau@eng.ucsd.edu

Harshini Rajachander: hrajacha@eng.ucsd.edu

UC San Diego

JACOBS SCHOOL OF ENGINEERING
Electrical and Computer Engineering

From Photometry to Geometry

- ▶ Suppose that instead of a lidar (which measures the positions of points in the world), we would like to use a camera to localize our robot and build a map of the environment
- ▶ **Image**: an array of positive numbers that measure the amount of light incident on the sensor
- ▶ How do we go from measurements of light (**photometry**) to measurements of positions of points in the world?

Correspondence



- ▶ **Corresponding points** in two views are image projections of the same geometric point in space
- ▶ **Correspondence problem:** establish which point in the second image corresponds to a given point $z_1 \in \mathbb{R}^2$ in the first image in the sense of being the same point in physical space
- ▶ **Idea:** look for a pixel $z_2 \in \mathbb{R}^2$ such that $I_2(z_2) \approx I_1(z_1)$

Correspondence

- ▶ **Matching windows:** a much more robust process of establishing correspondence is to compare not the brightness of individual pixels but that of small windows $W(z_1)$, $W(z_2)$ around the points
- ▶ **Aperture problem:** the brightness profile within the selected windows is not rich enough to allow us to recover the transformation of the pixel z_1 uniquely (e.g., blank wall)
- ▶ **Features:** points whose local regions are rich enough to allow solving the correspondence problem. Features establish a link between photometric measurements and geometric primitives.
- ▶ The window shape $W(z_1)$ and image values $I_1(y)$, $y \in W(z_1)$, associated with a pixel z_1 in the first image undergo a *nonlinear transformation* as a consequence of the change of viewpoint

Brightness constancy constraint

- ▶ Suppose we are imaging a point $m \in \mathbb{R}^3$ that emits light with the same energy in all directions (Lambertian) and radiance distribution $\mathcal{R}(m)$
- ▶ Suppose the camera is calibrated (i.e., $K = I_{3 \times 3}$) and the two camera frames are related by the rigid-body transformation $(R, p) \in SE(3)$.
- ▶ Let I_1 and I_2 be two images and $z_1, z_2 \in \mathbb{R}^2$ be the two pixels corresponding to m :

$$I_2(z_2) = I_1(z_1) \sim \mathcal{R}(m)$$

- ▶ From the projection equations, the point z_1 in image I_1 corresponds to the point z_2 in image I_2 if:

$$z_2 = g(z_1) := \frac{1}{\lambda_2} (\lambda_1 R z_1 + p)$$

where λ_1, λ_2 are the **unknown** scales/depths of the observed point m .

- ▶ **Brightness constancy constraint:** $I_1(z_1) = I_2(g(z_1))$

Local Deformation Models

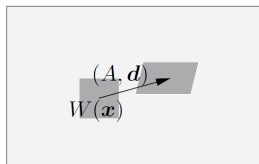
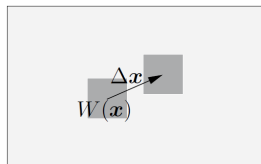
- ▶ The transformation g undergone by the entire image is determined by the scales λ_1, λ_2 of the visible surface and hence estimating g is as difficult as estimating the shape of the visible objects!
- ▶ Instead, we model the transformation only locally in a region $W(z)$:
 - ▶ **Translational model:** each point in the window undergoes the exact same translational motion $d \in \mathbb{R}^2$:

$$g(y) \approx y + d, \quad \forall y \in W(z)$$

This model is valid only in small windows and over short time durations but it is at the core of many feature matching and tracking algorithms.

- ▶ **Affine model:** each point in the window undergoes an affine transformation with parameters $A \in \mathbb{R}^{2 \times 2}$ and $d \in \mathbb{R}^2$:

$$g(y) \approx Ay + d, \quad \forall y \in W(z)$$



Matching Point Features

- ▶ Requiring that $I_1(z_1) = I_2(g(z_1))$ is too much to ask for due to the approximation of g and the presence of noise and occlusions
- ▶ **Correspondence problem:** an optimization problem that aims to determine the (translation or affine) parameters of the local transformation model of g :

$$\min_d \sum_{y \in W(z)} \|I_1(y) - I_2(y + d)\|_2^2 \quad \text{OR} \quad \min_{A,d} \sum_{y \in W(z)} \|I_1(y) - I_2(Ay + d)\|_2^2$$

- ▶ Our approximations of g are valid only locally in space and **time** so consider the continuous version of the brightness constancy constraint:

$$I_1(z) = I(z(t), t) \quad \underbrace{\approx}_{\text{brightness constancy}} \quad I_2(g(z)) \quad \underbrace{\approx}_{\text{translation model}} \quad I(z(t) + \nu dt, t + dt)$$

where dt is small and $\nu \in \mathbb{R}^2$ is the velocity of z

Continuous-Time Brightness Constancy

- ▶ **Brightness Constancy** (for the affine model):

$$I(z, t) \approx I(Az + \nu dt, t + dt)$$

- ▶ Linearizing the right-hand side around (z, t) :

$$I(Az + \nu dt, t + dt) \approx I(z, t) + \nabla_z I(z, t)^T (Az + \nu dt - z) + \frac{\partial I}{\partial t}(z, t) dt$$

leads to:

- ▶ Translational: $\min_{\nu} \sum_{y \in W(z)} \left\| \nabla_z I(y, t)^T \nu + \frac{\partial I}{\partial t}(y, t) \right\|_2^2$

- ▶ Affine: $\min_{A, \nu} \sum_{y \in W(z)} \left\| \nabla_z I(y, t)^T \left(\frac{(A - I)}{dt} y + \nu \right) + \frac{\partial I}{\partial t}(y, t) \right\|_2^2$

- ▶ **Aperture problem:** The brightness constancy equation $\left(\frac{\partial I}{\partial z} \nu + \frac{\partial I}{\partial t} = 0 \right)$ provides only one constraint for the two unknowns $\nu \in \mathbb{R}^2$.
- ▶ There are enough constraints on ν only when the brightness constancy constraint is applied to each y in a region $W(z)$ that contains “sufficient texture” and the motion ν is assumed constant in the region. 8

Feature Tracking and Optical Flow

- ▶ The brightness constancy equation ($\frac{\partial I}{\partial z} \nu + \frac{\partial I}{\partial t} = 0$) can be used to compute optical flow or track photometric features in a sequence of moving images
- ▶ **Optical flow**: the velocity ν of particle flowing through a given image location z
- ▶ **Feature tracking**: the computation of the velocity ν of a particle $z(t)$ moving through the image domain so that $z(t + dt) = z(t) + \nu dt$ (translational model)
- ▶ The only difference between optical flow and feature tracking is at the conceptual level, whether the vector ν is computed at fixed locations in the image or at moving points $z(t)$

Feature Tracking and Optical Flow

- ▶ To compute the velocity ν we need to solve:

$$\min_{\nu} \sum_{y \in W(z)} \left\| \nabla_z I(y, t)^T \nu + \frac{\partial I}{\partial t}(y, t) \right\|_2^2$$

- ▶ Letting $z = (u, v)$ and setting the gradient to zero results in:

$$\begin{aligned} 0 &= 2 \sum_{y \in W(z)} \left(\nabla_z I(y, t)^T \nu + \frac{\partial I}{\partial t}(y, t) \right) \nabla_z I(y, t) \\ &= 2 \sum_{y \in W(z)} \left(\begin{bmatrix} I_u^2(y) & I_u(y)I_v(y) \\ I_u(y)I_v(y) & I_v(y)^2 \end{bmatrix} \nu + \begin{bmatrix} I_u(y)I_t(y) \\ I_v(y)I_t(y) \end{bmatrix} \right) \\ &= 2 \left(\underbrace{\begin{bmatrix} \sum_y I_u^2(y) & \sum_y I_u(y)I_v(y) \\ \sum_y I_u(y)I_v(y) & \sum_y I_v(y)^2 \end{bmatrix}}_{G(z)} \nu + \underbrace{\begin{bmatrix} \sum_y I_u(y)I_t(y) \\ \sum_y I_v(y)I_t(y) \end{bmatrix}}_{b(z)} \right) \end{aligned}$$

- ▶ The optimal estimate of the image velocity at z is $\nu^* = -G(z)^{-1}b(z)$

Point Feature Selection

- ▶ For $G(z)$ to be invertible, the region $W(z)$ must have nontrivial gradients along independent directions, therefore resembling a “corner” structure.
- ▶ **Corner feature:** a pixel z such that the smallest singular value of $G(z)$ (equal to the eigenvalues for a symmetric matrix) is larger than some threshold τ
- ▶ **Harris corner detector:** A variation of the corner detector that thresholds the quantity:

$$\det(G) + k \operatorname{tr}^2(G) = \sigma_1 \sigma_2 + k(\sigma_1 + \sigma_2)^2 = (1 + 2k)\sigma_1 \sigma_2 + k(\sigma_1 + \sigma_2)^2,$$

where $k \in \mathbb{R}$ is a small scalar and σ_1, σ_2 are the singular values of G . Since k is small, both singular values of G need to be sufficiently large to pass the threshold.

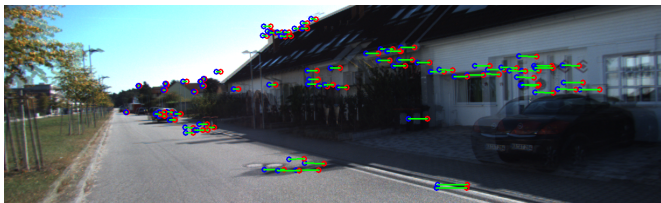
- ▶ More sophisticated techniques that utilize contours (or edges) and search for high curvature points in the detected contours are used in practice

Feature Tracking and Optical Flow

Algorithm 1 Basic Feature Tracking and Optical Flow

- 1: **Input:** Image I at time t
 - 2:
 - 3: Compute the image gradient (I_u, I_v)
 - 4: Compute $G(z) := \begin{bmatrix} \sum_{y \in W(z)} I_u^2(y) & \sum_{y \in W(z)} I_u(y)I_v(y) \\ \sum_{y \in W(z)} I_u(y)I_v(y) & \sum_{y \in W(z)} I_v^2(y) \end{bmatrix}$ at every pixel $z = (u, v)$
 - 5:
 - 6: (Feature tracking) select point features z_1, z_2, \dots such that $G(z_i)$ is invertible
 - 7: (Optical flow) select z_i on a fixed grid
 - 8:
 - 9: Compute $b(z) := \begin{bmatrix} \sum_{y \in W(z)} I_u(y)I_t(y) \\ \sum_{y \in W(z)} I_v(y)I_t(y) \end{bmatrix}$
 - 10:
 - 11: If $G(z)$ is invertible (guaranteed for point features), compute $\nu(z) = -G(z)^{-1}b(z)$
 - 12: Else $\nu(z) = 0$.
 - 13:
 - 14: (Feature tracking) at time $t + 1$, repeat the operation at $z + \nu(z)$
 - 15: (Optical flow) at time $t + 1$, repeat the operation at z
-

Feature Tracking and Optical Flow



Feature Tracking and Optical Flow

- ▶ The feature tracking/optical flow algorithm is very efficient when we use the translational deformation model
- ▶ When features are tracked over extended periods of time, however, the estimation error accumulates
- ▶ Instead of matching image regions between adjacent frames, one could match image regions between an initial frame and the current frame
- ▶ The simple translational deformation model is no longer accurate and we should use the affine deformation model
- ▶ Further reading:
 - ▶ J. Shi and C. Tomasi, "Good features to track," IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), pp. 593-600, 1994.

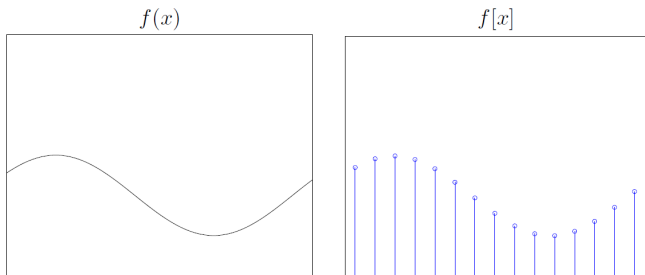
Image Gradient

- ▶ How do we compute the gradients $I_u(u, v, t)$, $I_v(u, v, t)$, and $I_t(u, v, t)$ needed for feature tracking/optical flow?
- ▶ We could approximate the derivatives using finite differences, e.g.,:

$$I_t(u, v, t) = I(u, v, t) - I_t(u, v, t - 1) \quad \text{OR} \quad I_t(u, v, t) = \frac{1}{2} (I(u, v, t + 1) - I_t(u, v, t - 1))$$

- ▶ To derive a more accurate approach we need to understand the relationship between a continuous signal $f(x)$ and its sampled version with period T :

$$f[x] = f(xT), \quad x \in \mathbb{Z}$$



Nyquist-Shannon Sampling Theorem

- ▶ If $f(x)$ is band limited, i.e., its Fourier transform satisfies $|F(\omega)| = 0$ for all $\omega > \omega_n$ (**Nyquist frequency**), it can be reconstructed exactly from a set of discrete samples at sampling frequency $\omega_s := \frac{2\pi}{T} > 2\omega_n$.
- ▶ The continuous signal $f(x)$ can be reconstructed by multiplying its sampled version $f[x]$ in the frequency domain with an ideal reconstruction filter $h(x)$ with Fourier transform:

$$H(\omega) = \begin{cases} 1, & \omega \in \left[-\frac{\pi}{T}, \frac{\pi}{T}\right] \\ 0, & \text{else} \end{cases} \quad h(x) = \mathbf{sinc} \left(\frac{\pi x}{T} \right), \quad x \in \mathbb{R}$$

- ▶ Multiplication in the frequency domain corresponds to convolution in the spatial domain, thus as long as $\omega_n < \frac{\pi}{T}$:

$$f(x) = f[x] * h(x), \quad x \in \mathbb{R}$$

Derivative of a Sampled Signal

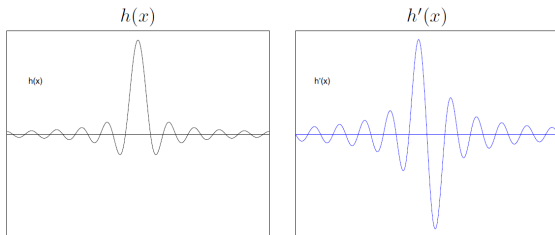
- ▶ Differentiating $f(x) = f[x] * h(x)$:

$$\frac{d}{dx} f(x) = \sum_{k=-\infty}^{\infty} f[k] \frac{d}{dx} h(x - k) = f[x] * \frac{dh}{dx}(x)$$

- ▶ Sampling the above result shows that the derivative of the sampled function $f'[x]$ can be computed as a convolution of the sampled signal $f[x]$ with the sampled derivative of the sinc function $h'[x]$:

$$f'[x] = f[x] * h'[x]$$

$$h'(x) = \frac{(\pi^2 x / T^2) \cos(\pi x / T) - \pi / T \sin(\pi x / T)}{(\pi x / T)^2}, \quad x \in \mathbb{R}$$

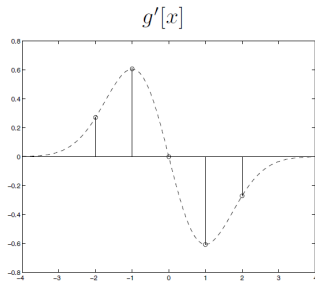
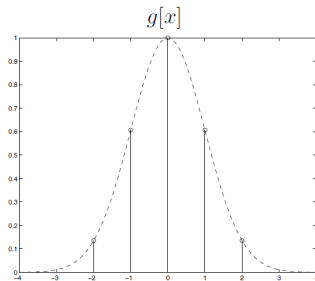


Five-tap Gaussian Filter

- ▶ The sinc function has infinite support and falls off very slowly away from the origin. Hence, the sinc convolution is not practically feasible and simple truncation yields undesirable artifacts.
- ▶ The derivative computation can be approximated by convolving with a Gaussian since it drops to zero much faster than the sinc:

$$g(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-x^2}{2\sigma^2}}$$

$$g'(x) = -\frac{x}{\sigma^2\sqrt{2\pi\sigma^2}} e^{\frac{-x^2}{2\sigma^2}}$$



$$g[x] = [0.1353 \quad 0.6065 \quad 1.0000 \quad 0.6065 \quad 0.1353]$$

$$g'[x] = [0.2707 \quad 0.6065 \quad 0 \quad -0.6065 \quad -0.2707]$$

Image Gradient

- ▶ In the case of images (2-D functions) the result is the same:

$$I(u, v) = I[u, v] * h(u, v) \quad h(u, v) = h(u)h(v) = \frac{\sin(\pi u/T) \sin(\pi v/T)}{\pi^2 uv/T^2},$$

- ▶ Note that $h(u, v) = h(u)h(v)$ is separable which leads to:

$$I_u[u, v] = I[u, v] * h'[u] * h[v] \quad I_v(u, v) = I[u, v] * h[u] * h'[v]$$

- ▶ The computation of the image derivatives is then accomplished as a pair of 1-D convolutions with filters obtained by sampling a continuous Gaussian function and its derivative:

$$I_u[u, v] = I[u, v] * g'[u] * g[v] = \sum_{k=-\omega/2}^{\omega/2} \sum_{l=-\omega/2}^{\omega/2} I[k, l] g'[u - k] g[v - l],$$

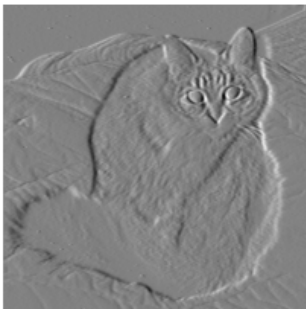
$$I_v[u, v] = I[u, v] * g[u] * g'[v] = \sum_{k=-\omega/2}^{\omega/2} \sum_{l=-\omega/2}^{\omega/2} I[k, l] g[u - k] g'[v - l]$$

- ▶ The number of samples is typically chosen as $\omega = 5\sigma$, imposing the fact that the window subtends 98.76% of the area under the Gaussian curve.

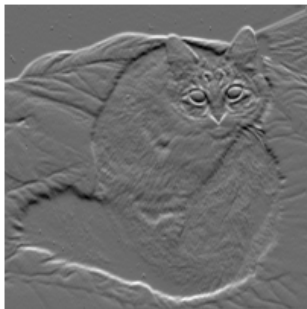
Image Gradient



I



I_u



I_v

Other Derivative Filters, Features, and Descriptors

- ▶ Other commonly used derivative filters:
 - ▶ **Interpolation filter:** $h[x] = \frac{1}{2}[1, 1]$ with derivative $h'[x] = \frac{1}{2}[1, -1]$
 - ▶ **Sobel filter:** $h[x] = \frac{1}{2+\sqrt{2}}[1, \sqrt{2}, 1]$ with derivative $h'[x] = \frac{1}{3}[1, 0, -1]$
 - ▶ **Gabor filter:** used for texture analysis
- ▶ Other features and descriptors (describe feature shape, color, texture):
 - ▶ **SIFT:** the Scale-Invariant Feature Transform (SIFT), introduced by David Lowe, is one of the most successful local image features/descriptors in the past decade. It makes the Harris corner scale invariant by using scale-space filtering via a Laplacian of Gaussian kernel (blob detector)
 - ▶ **SURF:** the Speeded-Up Robust Feature is a speeded-up version of SIFT which applies an approximate 2^{nd} derivative Gaussian filter at many scales along the axes and at 45° (more robust to rotation than Harris corners)
 - ▶ **FAST:** a Feature from Accelerated Segment Test detects corners by considering 16 pixels around the pixel y being tested and is several times faster than other corner detectors
 - ▶ **BRIEF:** a Binary Robust Independent Elementary Features speed up descriptor calculation and matching
 - ▶ **ORB:** Oriented FAST and Rotated BRIEF