

1 The digital image

- Problems of digital cameras sensors
- transmission interference
 - compression artefacts
 - spilling
 - scratches, sensor noise
 - bad contrast

1.1 Image as 2D signal

Signal: function depending on some variable with physical meaning
Image: continuous function
2 variables: xy -coordinates
3 variables: xy +time

Brightness is usually the value of the function, but other physical values are: Temperature, pressure, depth...

What is an image?

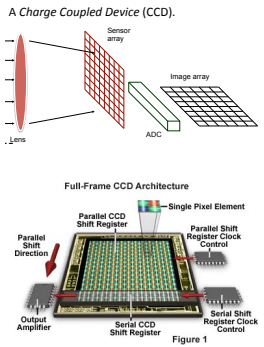
- A picture or pattern of a value varying in space and/or time
- Representation of a function $f: \mathbb{R}^n \rightarrow S$
- In digital form, e.g.: $I: \{1, \dots, X\} \times \{1, \dots, Y\} \rightarrow S$
- For greyscale CCD images, $n = 2, S = \mathbb{R}^+$

What is a pixel?

Not a little square! E.g. gaussian or cubic reconstruction filter.

1.2 Image sources

digital camera (CCD)



analog to digital Conversion

- The ADC measures the charge and digitizes the result.
- Conversion happens line by line.
- The charges in each photosite move down through the sensor array.

Blooming. Buckets have finite capacity. Photosite saturation causes blooming.

Bleeding or smearing during transit buckets still accumulate some charges. Due to tunneling and CCD Architecture. (Influenced by time "in transit" versus integration time. Effect is worse for short shutter times)

dark current CCDs produce thermally-generated charge. They give non-zero output even in darkness. Partly, this is the *dark current* and it fluctuates randomly. One can reduce it by cooling the CCD.

CMOS Has same sensor elements as CCD. Each photo sensor has its *own amplifier*. This leads to more noise (reduced by subtracting "black" image) and lower sensitivity (lower fill rate). The uses of standard CMOS technology allows to put other components on chip and "smart" pixels.

CCD vs. CMOS CCD: mature technology, specific technology, high production cost, high power consumption, higher fill rate, blooming, sequential read-out.

CMOS: recent technology, standard IC technology, cheap, low power, less sensitive, per pixel amplification, random pixel access, smart pixels, on chip integration with other components, rolling shutter (sequential read-out of lines)

1.3 Sampling

1D Sampling takes a function and returns a vector whose elements are values of that function at the sample points.

Undersampling "Missing" things between samples. Information lost

aliasing signals "traveling in disguise" as other frequencies. (Can happen in undersampling.)

1.4 Reconstruction

Inverse of sampling. Making samples back into continuous function. For output (need realizable method), for analysis or processing (need mathematical method), amounts to "guessing" what the function did in between.

Bilinear Interpolation

$$f(x, y) = (1 - a)(1 - b)f[i, j] + a(1 - b)f[i + 1, j] + abf[i + 1, j + 1] + (1 - a)bf[i, j + 1]$$

Nyquist frequency. Half the sampling frequency of a discrete signal processing system. Signal's max frequency (bandwidth) must be *smaller* than this.

sampling grids cartesian sampling, hexagonal sampling and non-uniform sampling

1.5 Quantization

real valued function will get digital values - integer values. Quantization is lossy and can't be reconstructed. Simple quantization uses equally spaced levels with k intervals.

usual quantization intervals

Grayscale image: 8 bit = $2^8 = 256$ gray-values. Color image RGB (3 channels): 8 bit/channel = $2^{24} = 16.7M$ colors. Nonlinear, for example log-scale.

1.6 Image Properties

Image resolution: Clipped when reduced. **Geometric resolution:** Whole picture but crappy when reduced. **Radiometric resolution:** Number of colors.

1.7 Image Noise

additive Gaussian Noise Common model $I(x, y) = f(x, y) + c$, where $c \sim N(0, \sigma^2)$. So that $p(c) = (2\pi\sigma^2)^{-1}e^{-c^2/2\sigma^2}$.

Poisson noise: (shot noise)

$$p(k) = \lambda^k e^{-\lambda} / k!$$
$$p(I) = \frac{1}{\sigma^2} \exp\left(-\frac{I^2 + f^2}{2\sigma^2}\right)$$
$$I_0\left(\frac{If}{\sigma^2}\right)$$

Multiplicative noise: $I = f + f \cdot c$

Signal to noise ration (SNR) $s = F/\sigma$ is an index of image quality, where

$$F = \frac{1}{XY} \sum_{x=1}^X \sum_{y=1}^Y f(x, y)$$

Often used instead: **Peak Signal to Noise Ratio (PSNR)** $\text{Speak} = F_{\max}/\sigma$

1.8 Colour Images

Consist of red, green and blue channel.

Prism (with 3 sensors) Separate light in three beams using dichroic prism. Requires 3 sensors and precise alignment. Gives good color separation. → high-end cameras

Filter mosaic Coat filter directly on sensor. "Demosaiicing" to obtain full colour & full resolution image. → low-end cameras

Filter wheel rotate multiple filters in front of lens. Allows more than 3 colour bands. → static scenes

new color CMOS sensor, foveon's X3 blue, green, red sensor, one above the other (descending) → better image quality

2 Image segmentation

TP fraction = true positive count/P,

$P = TP + FN$ and false positive fraction

FP fraction = false positive count/N, $N = FP + TN$. ROC curve always passes through (0, 0) and (1, 1).

MAP (Maximum A Posteriori) detector

Operating points choose an operating point by assigning relative costs and values to each outcome, VT_N, VT_P, CF_N, CF_P . V and C being values and costs. For simplicity, often $VT_N = VT_P = 0$.

Interim Summary Segmentation is hard. It is easier if you define the task carefully: (1) Segmentation task binary or continuous? (2) What are regions of interest? (3) How accurately must the algorithm locate the region boundaries?

Definition it partitions an image into regions of interest. It is the first stage in many automatic image analysis systems. A complete segmentation of an image I is a finite set of regions R_1, \dots, R_N , such that

$$I = \bigcup_{i=1}^N R_i \text{ and } R_i \cap R_j = \emptyset, \quad \forall i \neq j$$

segmentation quality the quality of a segmentation depends on what you want to do with it. Segmentation algorithms must be chosen and evaluated with an application in mind.

2.1 Thresholding

Is a simple segmentation process, produces a binary image B . It labels each pixel in or out of the region of interest by

comparison of the greylevel with a threshold T :

$$B(x, y) = \begin{cases} 1 & \text{if } I(x, y) \geq T \\ 0 & \text{if } I(x, y) < T. \end{cases}$$

Choosing T By trial and error. Compare results with ground truth. Automatic methods. (ROC curve)

Chromakeying Control Lighting! "Plain" distance measure (e.g.)

$$I_\alpha = |I - g| > T$$

$T \sim 20, g = \begin{pmatrix} 0 & 255 & 0 \end{pmatrix}^T$
Problems: Variation is *not* the same in all 3 channels. Hard alpha maske

$$I_{\text{comp}} = I_\alpha I_a + (1 - I_\alpha) I_b$$

Gaussian model per pixel (Like chromakeying.) mean $\mu \rightarrow I_\mu$, standard deviation $\sigma \rightarrow I_\Sigma$. $I_\alpha = |I - I_\mu| \geq I_\Sigma$

$T, T = \begin{pmatrix} 20 & 20 & 10 \end{pmatrix}$, $I_{bg} = \begin{pmatrix} 17 & 17 & 17 \end{pmatrix}$
background image. Or better (e.g.)

$$I_\alpha = \sqrt{(I - I_{bg})^T \Sigma^{-1} (I - I_{bg})} > T = 4$$

ROC Analysis Receiver operating Characteristic. An ROC curve characterizes the performance of a binary classifier. A binary classifier distinguishes between two different types of things.

Classification error. Binary classifiers make errors. Two types of input to a binary classifier: Positives, negatives. Four possible outcomes in any test: True positive, true negative, false negative, false positive.

ROC Curve. Characterizes the error trade-off in binary classification tasks. It plots the *true positive fraction*

TP fraction = true positive count/P,

$P = TP + FN$ and false positive fraction

FP fraction = false positive count/N, $N = FP + TN$. ROC curve always passes through (0, 0) and (1, 1).

MAP (Maximum A Posteriori) detector

Operating points choose an operating point by assigning relative costs and values to each outcome, VT_N, VT_P, CF_N, CF_P . V and C being values and costs. For simplicity, often $VT_N = VT_P = 0$.

Performance Assessment In real-life, we use two or even three separate sets of test data: (1) A *training set*, for tuning the algorithm, (2) A *validation set* for tuning the performance score, (3) An *unseen test set* to get a final performance score on the tuned algorithm.

Pixel connectivity Define neighbors, e.g. (for 2D) 4-neighborhood or 8-neighborhood

Pixel paths There are e.g. 4- and 8-connected paths. (p_i neighbor of p_{i+1}).

Connected regions A region is 4- or 8-connected if it contains a(n) 4- or 8-connected path between any two of its pixels.

2.2 Region Growing

- (1) Start from a seed point or region.
- (2) Add neighboring pixels that satisfy the criteria defining a region. (3) Repeat until we can include no more pixels.

```
function B = RegionGrow(I, seed)
[X,Y] = size(I);
visited = zeros(X,Y);
visited(seed) = 1;
boundary = emptyQ;
boundary.enQ(seed);
while (~boundary.empty())
    nextPoint = boundary.deQ();
    if (include(nextPoint, seed))
        visited(nextPoint) = 2;
        ForEach(x,y) in N(nextPoint)
            if (visited(x,y) == 0)
                boundary.enQ(x,y);
            visited(x,y) = 1;
        end
    end
end
```

2.2.1 Variations

seed selection (1) One seed point, (2) Seed region, (3) Multiple seeds.

seed selection (1) Greylevel thresholding, (2) Greylevel distribution model. E.g. include if $(I(x, y) - \mu)^2 < (n\sigma)^2$, $n = 3$. Can update μ and σ after every iteration, (3) color or texture information.

snakes A snake is an *active contour*. It's a polygon. Each point on contour moves away from seed while its image neighborhood satisfies an inclusion criterion. Often the contour has smoothness constraints, the algorithm iteratively minimizes an energy function:

$$E = E_{\text{tension}} + E_{\text{stiffness}} + E_{\text{image}}$$

2.3 Spatial relations

Markov Random Fields Markov chains have 1D structure. At every time, there is one state. This enabled use of dynamic programming. *Markov Random fields* break this 1D structure:

- Field of sites, each of which has a label, simultaneously.
- Label at one site dependent on others, no 1D structure dependencies.
- This means no optimal, efficient algorithms, except for 2-label problems. Minimize

$$\text{Energy}(y; \theta, \text{data}) = \sum_i \psi_1(y_i; \theta, \text{data}) + \sum_{\langle i, j \rangle \in \text{edges}} \psi_2(y_i, y_j; \theta, \text{data})$$

FG-BG segmentation The code does the following: • background RGB Gaussian model training (from many images) • shadow modeling (hard shadow and soft shadow) • graphcut foreground-background segmentation

Hit-and-miss transform $H = I \otimes S$ Searches for an exact match of the structuring element. Simple form of template matching.

Thinning $I \otimes S = I \setminus (I \otimes S)$

Thickening $I \otimes S = I \cup (I \otimes S)$

Sequential thinning/thickening With structuring elements S_1, \dots, S_n and

8-neighbor erode (Minkowsky subtraction) Erase any foreground pixel that has one eight-connected neighbor that is background

8-neighbor dilate (Minkowsky addition) Paint any background pixel that has one eight-connected neighbor that is foreground. *Applications:* Smooth region boundaries for shape analysis, remove noise and artefacts from an imperfect segmentation, match particular pixel configurations in an image for simple object recognition

structuring elements morphological operations take two arguments 1. a binary image 2. a structuring element Compare the structuring element to the neighborhood of each pixel. This determines the output of the morphological operation. The structuring element is also a binary array and has an origin.

$$I_1 \cup I_2 = \{x: x \in I_1 \text{ or } x \in I_2\},$$

$$I_2 \cap I_1 = \{x: x \in I_1 \text{ and } x \in I_2\},$$

$$I^C = \{x: x \notin I\},$$

$$I_1 \setminus I_2 = \{x: x \in I_2 \text{ and } x \notin I_2\}.$$

Erosion of binary image I by the structuring element S is defined by

$$I \ominus S = \{z \in E \mid S_z \subset I\}$$

S_z translation of S by vector z .

Dilation is $I \oplus S = \bigcup_{b \in S} I_b$.

Opening $I \circ S = (I \ominus S) \oplus S$.

Closing $I \bullet S = (I \oplus S) \ominus S$.

To remove holes in the foreground and islands in the background, do both opening and closing. Thesize and shape of the structuring element determine which features survive. In the absence of knowledge about the shape of features to remove, use a circular structuring element.

Granulometry Provides a size distribution of distinct regions or "granules" in the image. We open (opening as above) the image with increasing structuring element size and count the number of regions after each operation. Creates "morphological sieve".

```
function gSpec = granulo(I, T, maxRad)
% Segment the image I.
B = (I > T);
%Open the image at each structuring
element size up
%to a maximum and count the
remaining regions.
for x=1:maxRad
    O = imopen(B, strel('disk',x));
end
gSpec = diff(numRegions);
```

Linear Filtering Linear operations can be written:

$$I'_j = \sum_{i=1}^N \alpha_i I_i, \quad j = 1 \dots N$$

$I'_j =$ output of operation. k is kernel of the operation. $N(m, n)$ is a neighbourhood of (m, n) .

Correlation e.g. template matching. Linear operation: $I' = KI$

$I'(x, y) = \sum_{i,j \in N(x,y)} K(i, j) I(x + i, y + j)$

sequential thinning/thickening ♦

$I \bullet \{S_i: i = 1, \dots, n\} = ((I \bullet S_1) \bullet \dots \bullet S_n)$

2.4.1 Medial Axis Transform (MAT, skeletonization)

The skeleton and MAT are stick-figure representations of a region $X \in \mathbb{R}^2$. Start a grassfire at the boundary of the region, theskeleton is the set of points at which two fire fronts meet.

Skeleton Use structuring element

$$B = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

The n -th skeleton subset is

$$S_n(X) = (X \ominus_n B) \setminus [(X \ominus_n B) \ominus B]$$

where \ominus_n denotes n successive erosions. The skeleton is the union of all the skeleton subsets $S(X) = \bigcup_{n=1}^{\infty} S_n(X)$.

Reconstruction can reconstruct region X from its *skeleton subsets*.

$$X = \bigcup_{n=0}^{\infty} S_n(X) \oplus_n B$$

Applications and problems The skeleton/MAT provies a stick figure representing the region shap. Used in object recognition, in particula, character recognition. *Problems:* Definition of a maximal disc is poorly defined on a digital grid and is sensitive to noise on the boundary. Sequential thinning output sometimes preferred to skeleton/MAT.

Smoothing filters: $\frac{1}{10} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix}$.

Gaussian Kernel Idea: Weight contributions of neighboring pixels:

$$N_{\mu=0, \sigma}(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

Smoothing with a Gaussian instead of a box filter removes the artefact of the vertical and horizontal lines. Gaussian smoothing Kernel is *separable*! $N(x, y) = N(x)N(y)$. Amount of smoothing depends on σ and window size. Width $> 3\sigma$.

Scale space Convolution of a Gaussian with σ with itself is a gaussian with $\sigma\sqrt{2}$. Repeated convolution by a Gaussian filter produces the scale space of an image.

Gaussian filter top-5 (1) Rotationally symmetric. (2) Has a single lobe. → Neighbor's influence decreases monotonically. (3) Still one lobe in frequency domain. → No corruption from high frequencies (4) Simple relationship to σ (5) Easy to implement efficiently

Differential filters • Prewitt operator: $\begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$, • Sobel operator: $\begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}$.

High-pass filters • Laplacian operator: $\begin{pmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{pmatrix}$, • High-pass filter: $\begin{pmatrix} 1 & -1 & 1 \\ -1 & 4 & -1 \\ 1 & -1 & 1 \end{pmatrix}$.

Linear Filtering Linear operations can be written:

$$I'_j = \sum_{i=1}^N \alpha_i I_i, \quad j = 1 \dots N$$

$I'_j =$ output of operation. k is kernel of the operation. $N(m, n)$ is a neighbourhood of (m, n) .

Correlation e.g. template matching. Linear operation: $I' = KI$

$I'(x, y) = \sum_{i,j \in N(x,y)} K(i, j) I(x + i, y + j)$

Differentiation and convolution

$$\frac{\partial f}{\partial x} = \lim_{\varepsilon \rightarrow 0} \left(\frac{f(x + \varepsilon, y) - f(x, y)}{\varepsilon} \right)$$
$$\frac{\partial f}{\partial x} = \lim_{\varepsilon \rightarrow 0} \left(\frac{f(x + \varepsilon, y) - f(x, y)}{\varepsilon} \right)$$

which is obviously a convolution (-1)

Filters and templates Filters at some point can be seen as taking a -product between the image and some vector, the image is a set of dot products, filters look like the effects they are intended to find, filters find effects they look like.

Image sharpening Also known as enhancement. Increases the high frequency components to enhance edges.

$$I' = I + \alpha |K * I|,$$

where K is a high-pass filter kernel and $\alpha \in [0, 1]$.

Integral images integral images (also known as summed-area tables) allow to efficiently compute the convolution with a constant rectangle

$$I(x, y) = \int_0^x \int_0^y dy' I(x', y')$$
$$A = I(1),$$
$$A + C = I(3),$$
$$A + B = I(2),$$
$$A + B + C + D = I(4).$$
$$D = I(4) - I(2) - I(3) + I(1)$$

Also possible along diagonal.

Viola-Jones cascade face detection Very efficient face detection using integral images.

4 Image features

4.1 Template matching

Problem Locate an object, described by a template $t(x, y)$, in the image $s(x, y)$. *Example:* Passport photo as image and eyes to detect.

Method Search for the best match by minimizing mean -squared error $E(p, q)$

$$E(p, q) = \sum_{x, y=-\infty}^{\infty} [s(x, y) - t(x - p, y - q)]^2$$
$$= \sum_{x, y=-\infty}^{\infty} |s(x, y)|^2 + |t(x, y)|^2$$
$$- 2 \sum_{x, y=-\infty}^{\infty} s(x, y) \cdot t(x - p, y - q)$$

Equivalently, maximize *area correlation*

$$r(p, q) = \sum_{x, y=-\infty}^{\infty} s(x, y) \cdot t(x - p, y - q)$$
$$= s(p, q) * t(-p, -q)$$

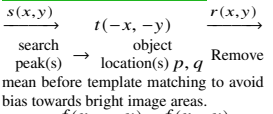
$\leq \sqrt{\left[\sum |s(x, y)|^2 \right] \cdot \left[\sum |t(x, y)|^2 \right]}$

where in the last step the Cauchy-Schwarz inequality was used. Equality \Leftrightarrow

$$s(x, y) = \alpha \cdot t(x - p, y - q) \quad \text{with } \alpha$$

Area correlation is equivalent to convolution of image $s(x, y)$ with impulse response $t(-x, -y)$.

Diagram of template matcher



4.2 Edge detection

Idea (continuous-space): Detect local gradient

$$\|\nabla(f(x, y))\| = \sqrt{(\partial_x f)^2 + (\partial_y f)^2}$$

Digital image: Use finite differences instead:

difference (-1) , **central difference** $(-1 \ 0 \ 1)$; **Prewitt** $\begin{pmatrix} -1 & 0 & 1 \\ -1 & [0] & 1 \\ -1 & 0 & 1 \end{pmatrix}$, **Sobel** $\begin{pmatrix} -1 & -1 & 1 \\ -1 & [0] & 1 \\ -1 & 1 & 1 \end{pmatrix}$, **Roberts** $\begin{pmatrix} [0] & 1 \\ -1 & 0 \end{pmatrix}$

Laplacian operator Detects discontinuities by considering second derivative

$$\nabla^2 f(x, y) = \partial_x^2 f(x, y) + \partial_y^2 f(x, y)$$

Isotropic (rotationally invariant) operator, zero-crossings mark edge location, discrete-space approximation by convolution with 3×3 impulse response

$$\begin{pmatrix} 0 & 1 & 0 \\ 1 & [-4] & 1 \\ 0 & 1 & 0 \end{pmatrix}, \text{ or } \begin{pmatrix} 0 & 1 & 0 \\ 1 & [-8] & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

Laplacian of Gaussian The Laplacian operator is very sensity to fine detail and noise, so blur it first with Gaussian. \rightarrow do it in one operator Laplacian of Gaussian (LoG)

$$\text{LoG}(x, y) = -\frac{1}{\pi\sigma^4} \left[1 - \frac{x^2 + y^2}{2\sigma^2} \right] \cdot e^{-\frac{x^2 + y^2}{2\sigma^2}}$$

4.2.1 Canny edge detector

(1) Smooth image with a gaussian filter
(2) Compute gradient magnitude and angle (Sobel, Prewitt,...)

$$M(x, y) = \sqrt{(\partial_x f)^2 + (\partial_y f)^2}$$
$$\alpha(x, y) = \arctan(\partial_y f / \partial_x f)$$

(3) Apply nonmaxima suppression to gradient magnitude image (4) Double thresholding to detect strong and weak edge pixels (5) Reject weak edge pixels not connected with strong edge pixels

Canny nonmaxima suppression

Quantize edge normal to one of four directions: horizontal, -45° , vertical, 45° . If $M(x, y)$ is smaller than either of its neighbors in edge normal direction \rightarrow suppress; else keep

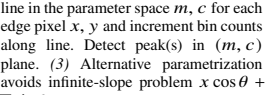
Double-thresh. of grad. magn.

strong edge: $M(x, y) \geq \theta_{\text{high}}$
weak edge: $\theta_{\text{high}} > M(x, y) \geq \theta_{\text{low}}$
Typical setting: $\theta_{\text{high}}, \theta_{\text{low}} = 2, 3$. Region labeling of edge pixels. Reject regions without strong edge pixels.

4.3 Feature detection

4.3.1 Hough transform

Problem: fit a straight line (or curve) to a set of edge pixels. Hough transform (1962): generalized template matching technique. (1) Consider detection of straight lines $y = mx + c$. (2) draw a line in the parameter space m, c for each edge pixel x, y and increment bin counts along line. Detect peak(s) in (m, c) plane. (3) Alternative parametrization avoids infinite-slope problem $x \cos \theta + y \sin \theta = \rho$



circle detection find circles of fixed radius r . For circles of undetermined radius, use 3d Hough transform for parameters (x_0, y_0, r)

Corner importance weight Give more importance to central pixels by using Gaussian weighting function

$$M = \sum_{x, y \in \text{window}} G(x - x_0, y - y_0, \sigma)$$
$$\begin{pmatrix} \partial_x f & \partial_x f \partial_y f \\ \partial_x f \partial_y f & \partial_y f \end{pmatrix}$$

Many applications benefit from features localized in (x, y) . Edges well localized only in one direction \rightarrow detect corners. Desirable properties of corner detector: (1) Accurate localization, (2) invariance against shift, rotation, scale, brightness change, (3) robust against noise, high repeatability

4.3.2 Detecting corner points

Robustness of Harris corner detector (1) Invariant to brightness offset: $f(x, y) \rightarrow f(x, y) + c$ (2) Invariant to shift and rotation (3) Not invariant to scaling

4.3.3 Most accurately localizable patterns

Local displacement sensitivity

$$S(\Delta x, \Delta y) = \sum_{x, y \in \text{window}} [f(x, y) - f(x - \Delta x, y - \Delta y)]^2$$
$$\approx f(x, y) + \partial_x f(x, y) \Delta x + \partial_y f(x, y) \Delta y$$
$$S(\Delta x, \Delta y) \approx \sum_{(x, y) \in \text{window}} \begin{pmatrix} \partial_x f & \partial_y f \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix} \begin{pmatrix} \partial_x f \\ \partial_y f \end{pmatrix}$$

Linear approximation for small $\Delta x, \Delta y$

$$\text{SSD} \approx \Delta^T M \Delta$$

Find points for which the following is large

$$\min \Delta^T M \Delta$$

for $\|\Delta\| = 1$. i.e. maximize eigenvalues of M .

Keypoint detection Often based on eigenvalues λ_1, λ_2 of M ("structure matrix"/"normal matrix"/"second-moment matrix")

$$M = \sum_{(x, y) \in \text{window}} \begin{pmatrix} (\partial_x f)^2 & \partial_x f \partial_y f \\ \partial_x f \partial_y f & (\partial_y f)^2 \end{pmatrix}$$

Measure of "corneriness"

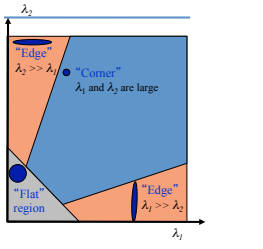
$$C(x, y) = \det(M) - k \cdot (\text{trace } M)^2$$
$$= \lambda_1 \lambda_2 - k \cdot (\lambda_1 + \lambda_2)^2$$

Represent function on a new basis. Basis elements have the form $e^{-i2\pi(u.x + v.y)}$. The Fourier transform is

$$\hat{f}(u, v) = \int_{\mathbb{R}^2} dx dy f(x, y) e^{-i2\pi(u.x + v.y)}$$

Basis functions of Fourier transform are eigenfunctions of linear systems.

Important functions



Corner importance weight Give more importance to central pixels by using Gaussian weighting function

$$M = \sum_{x, y \in \text{window}} G(x - x_0, y - y_0, \sigma)$$
$$\begin{pmatrix} \partial_x f & \partial_x f \partial_y f \\ \partial_x f \partial_y f & \partial_y f \end{pmatrix}$$

Compute subpixel localization by fitting parabola to *cornerness function*

Robustness of Harris corner detector (1) Invariant to brightness offset: $f(x, y) \rightarrow f(x, y) + c$ (2) Invariant to shift and rotation (3) Not invariant to scaling

4.3.4 Lowe's SIFT features

Recover features with position, orientation and scale.

Position (1) Look for strong responses of DoG filter, (2) only consider local maxima.

Scale (1) Look for strong responses of DoG filter, (2) only consider local maxima in both position and scale. (3) Fit quadratic around maxima for subpixel accuracy.

Orientation (1) Create histogram of local gradient directions computed at selected scale. (2) Assign canonical orientation at peak of smoothed histogram. (3) Each key specifies stable 2D coordinates (x, y, scale, orientation)

SIFT descriptor (1) Thresholded image gradients are sampled over 16×16 array of locations in scale space. (2) Create array of orientation histograms (3) 8 orientations $\times 4 \times 4$ histogram array = 128 dimensions

5 Fourier Transform

5.1 Aliasing

One can't shrink an image by taking every second pixel. If we do, characteristic errors appear. Typically, small phenomena look bigger; fast phenomena can look slower. Common phenomena (1) Wagon wheels rolling the wrong way in movies. (2) Checkerboards misrepresented in ray tracing (3) Striped shirts look funny on color television.

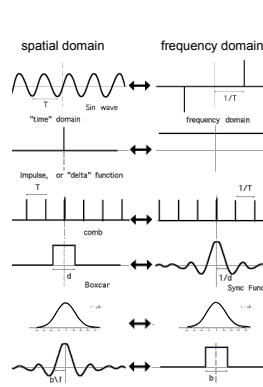
5.2 Definition

Represent function on a new basis. Basis elements have the form $e^{-i2\pi(u.x + v.y)}$. The Fourier transform is

$$\hat{f}(u, v) = \int_{\mathbb{R}^2} dx dy f(x, y) e^{-i2\pi(u.x + v.y)}$$

Basis functions of Fourier transform are eigenfunctions of linear systems.

Important functions



Convolution theorem (1) The Fourier transform of the convolution of two functions is the product of their Fourier transforms

$$\hat{f} \cdot \hat{g} = \widehat{f * g}$$

The Fourier transform of the product of two functions is the convolution of the Fourier transforms

$$\hat{f * g} = \mathcal{F}(f \cdot g)$$

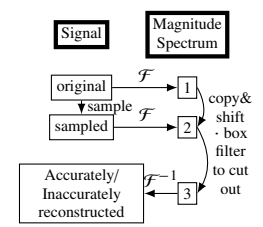
5.3 Sampling

Go from continuous world to discrete world, from function to vector. Samples are typically measured on regular grid. We want to be able to approximate integrals sensibly \rightarrow Delta function

$$S_{2D}(f(x, y)) = \sum_{i, j=-\infty}^{\infty} f(x, y) \delta(x - i, y - j)$$
$$= f(x, y) \sum_{i, j=-\infty}^{\infty} \delta(x - i, y - j),$$

with S = Sample operator.

FT of sampled signal



In the figure above the accuracy depends on the overlapping wave functions in "2". The box filter then can't cut out appropriately the magnitude spectrum to get a proper result in "3". This leads to an inaccurately reconstructed signal.

Proper sampling To avoid this effect, this is the procedure:

original signal $\xrightarrow{\text{lp filtering}}$ lp filt. sign $\xrightarrow{\text{sample}}$ sampl. sign. $\xrightarrow{\text{reconstr.}}$ reconstr. sign

Smoothing as low-pass filtering

The message of the FT is that high frequencies lead to trouble with sampling. Solutions suppress high frequencies before sampling. A filter whose FT is a box is *bad*, because the filter kernel has infinite support. Common solution: use a Gaussian.

Nyquist sampling theorem Nyquist theorem: The sampling frequency must

be at least twice the highest frequency. $\omega_s \geq 2\omega$. If this is not the case, the signal needs to be bandlimited before sampling, e.g. with a low-pass filter.

5.4 Image Restoration

Pixelization Possibilities: Square pixels, Gaussian reconstruction filter, Bilinear interpolation, perfect reconstruction filter.

Motion blurring Each light dot is transformed into a short line along the x_1 -axis:

$$h(x_1, x_2) = \frac{1}{2\ell} [\theta(x_1 + \ell) - \theta(x_1 - \ell)] \delta(x_2)$$

Noise Gaussian blurring kernel:

$$h(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{x_1^2 + x_2^2}{2\sigma^2}\right)$$

$$\mathcal{F}[\hat{h}](u, v) \cdot \mathcal{F}[h](u, v) = 1$$
$$(\hat{h} * h)(x) = \delta(x)$$

$\hat{h}(x)$ may be determined more easily in Fourier space:

$$E[c_i^2] = [R_{cc}] = [AR_{ff}A^*]_{ii}$$

To determine $\mathcal{F}[\hat{h}]$, we need to estimate (1) the distortion model $h(x)$ (point spread function) or $\mathcal{F}[h](u, v)$ (modulation transfer function) (2) the parameters of $h(x)$, e.g. for defocussing

Motion Blur FT Eq. 4 Problem: $\mathcal{F}[\hat{h}(u) = 1/h(u)]$. sinc has many zeroes and these frequencies can't be recovered! Solution: Regularized reconstruction filter

$$\hat{F}[\hat{h}](u, v) = \frac{\mathcal{F}[h]}{\|\mathcal{F}\|^2 + \varepsilon}$$

Singularities are avoided by the regularization ε .

Space-time super-resolution One can put two movies of the same thing and merge their frames for space and time super-resolution.

Spatial super-resolution

• lens + pixel = low-pass filter (ediseder to avoid aliasing) • Low-res images = $D * H * G * (\text{desired high-res-image})$. D:decimate, H:lens+pixel, G: Geometric warp • Simplified case for translation: $LR = (D * G) * (H * HR)$. G is shift-invariant and commutes with H. First compute H HR, then deconvolve HR with H. • Super-resolution needs to restore attenuated frequencies. Many images improve S/N ratio $\sim \sqrt{n}$, which helps. Eventually Gaussian's double exponential always dominates.

6 Unitary transforms

Digital image as a matrix:

$$f = \begin{bmatrix} f(0,0) & \dots & f(N-1,0) \\ \vdots & \ddots & \vdots \\ f(0,L-1) & \dots & f(N-1,L-1) \end{bmatrix}$$
$$= f_{j,y}$$

or as a vector

$$\mathbf{f} = \begin{pmatrix} f(0, 0) \\ \vdots \\ f(N - 1, L - 1) \end{pmatrix}$$

General approach (1) Sort samples $f(x, y)$ of an $M \times N$ image (or rectangular block in the image) into column vector of length $M \times N$. (2) Compute transform coefficients $\mathbf{c} = A\mathbf{f}$ where A is a matrix of size $(MN)^2$. (3) Transform A is unitary, iff $A^{-1} = A^*$ (4) If A is real-valued, i.e. $A = \bar{A}$, transform is orthonormal.

Energy conservation $\|\mathbf{c}\|^2 = \mathbf{c}^* \mathbf{c} = \mathbf{f}^* A^* A \mathbf{f} = \|\mathbf{f}\|^2$

Image collection f_i one image, $F = f_1, \dots, f_n$,

Auto-correlation function

$$R_{ff} = E[f_i f_i^*] = F F^* / n$$

energy distribution Energy is conserved, but often will be unevenly distributed among coefficients. Autocorrelation matrix:

$$R_{cc} = E[\mathbf{c}\mathbf{c}^*] = E[\mathbf{A}\mathbf{f}\mathbf{f}^*\mathbf{A}^*] = \mathbf{A}R_{ff}\mathbf{A}^*$$

Mean squared values ("average energies") of the coefficients c_i are on the diagonal of R_{cc} .

$$E[c_i^2] = [R_{cc}] = [\mathbf{A}R_{ff}\mathbf{A}^*]_{ii}$$

Eigenmatrix of autocorrelation matrix Definition: Eigenmatrix Φ of autocorrelation matrix R_{ff} . (1) ϕ is unitary (2) The columns of Φ form a set of eigenvectors of R_{ff} , i.e.,

$$R_{ff}\Phi = \Phi\Lambda,$$

where $\Lambda = \text{diag}(\lambda_0, \dots, \lambda_{MN-1})$. (3) R_{ff} is symmetric nonnegative definite, hence $\lambda_i \geq 0$ for all i (4) R_{ff} is normal matrix, i.e. $R_{ff}^* R_{ff} = R_{ff} R_{ff}^*$, hence unitary eigenmatrix exists.

$$\arg \min_i D_i = \|\mathbf{I}_i - \mathbf{I}\|$$

Computationally expensive, i.e. requires presented image to be correlated with every image in the database!

6.1 Karhunen-Loeve Transform

Strongly correlated samples with equal energies \rightarrow uncorrelated samples, most of the energy in first coefficient.

Properties (1) Unitary transform with matrix $A = \Phi^*$ where the columns of Φ are ordered according to decreasing eigenvalues. (2) Transform coefficients are pairwise uncorrelated

$R_{cc} = \mathbf{A}R_{ff}\mathbf{A}^* = \Phi^* R_{ff} \Phi = \Lambda$

(3) Energy concentration property: No other unitary transform packs as much energy into the first J coefficients, where J is arbitrary. Mean squared approximation error by choosing only first J coefficients is minimized.

Optimal energy concentration

(1) To show optimum energy concentration property, consider the truncated coefficient vector $\mathbf{b} = \mathbf{I}_J \mathbf{c}$, where \mathbf{I}_J contain ones on the first J diagonal positions, else zeros. (2) Energy in first J coefficients for arbitrary transform A

$$E = \text{tr}(\mathbf{R}_{bb}) = \text{tr}(\mathbf{I}_J \mathbf{R}_{cc} \mathbf{I}_J)$$

$$= \text{tr}(\mathbf{I}_J \mathbf{A} R_{ff} \mathbf{A}^* \mathbf{I}_J) = \sum_{k=0}^{J-1} a_k^T R_{ff} \bar{a}_k$$

where a_k^T is the k -th row of A . (3) Lagrangian cost function to enforce unit-length basis vectors

$$L = E + \sum_{k=0}^{J-1} \lambda_k (1 - a_k^T \bar{a}_k)$$

Differentiating L with respect to a_j yields necessary condition

$$R_{ff} \bar{a}_j = \lambda_j \bar{a}_j, \quad \forall j < J$$

6.2 Basis images and eigenimages (EI)

For a unitary transform, the inverse transform $\mathbf{f} = A^* \mathbf{c}$ can be interpreted in terms of the superpositions of "basis images" (columns of A^*) of size MN . If the transform is a KL transform, the basis images, which are the eigenvectors of the autocorrelation matrix R_{ff} , are called "eigenimages". If energy concentration works well, only a limited number of eigenimages is needed to approximate a set of images with small error. These eigenimages form an optimal linear subspace of dimensionality J .

EI for recognition To recognize complex patterns (e.g., faces), large portions of an image (say of size MN) might have to be considered. High dimensionality of "image space" means high computational burden for many recognition techniques. Transform $\mathbf{c} = \mathbf{W}\mathbf{f}$ can reduce dimensionality from MN to J by representing the image by J coefficients. Idea: tailor a KLT to the specific set of images of the recognition task to preserve the salient features.

Simple recognition Simple Euclidean distance (SSD) between images. Best match wins

$$\arg \min_i D_i = \|\mathbf{I}_i - \mathbf{I}\|$$

Computationally expensive, i.e. requires presented image to be correlated with every image in the database!

Eigenspace matching Let \mathbf{I}_i be the input image, \mathbf{I} the database. The "character" of the face $\mathbf{I} = J - \langle \mathbf{I} \rangle$, with J being any image (set). Do KLT (aka PCA) transformation

$$\hat{\mathbf{I}}_i \mapsto \mathbf{p}_i, \quad E^* \hat{\mathbf{I}}_i = \mathbf{p}_i.$$

$$\sim \hat{\mathbf{I}}_i \approx E \mathbf{p}_i,$$

$$\sim \mathbf{I}_i - \mathbf{I} = \hat{\mathbf{I}}_i - \hat{\mathbf{I}} \approx E(\mathbf{p}_i - \mathbf{p})$$

$$\sim \|\mathbf{I}_i - \mathbf{I}\| \approx \|\mathbf{p}_i - \mathbf{p}\|,$$

with closest rank-k approximation property of SVD. Approximate

$$\arg \min_i D_i = \|\mathbf{I}_i - \mathbf{I}\| \approx \|\mathbf{p}_i - \mathbf{p}\|$$

6.3 Eigenfaces (EF)

Concatenate face pixels into "observation vector", \mathbf{x} .

EI for recognition (1) Input image, (2) normalize, (3) subtract mean face, (4) KLT, (5) Find most similar \mathbf{p}_i , (6) similarity measure, (7) rejection system, (8) result of identification.

Limitations of EFs Differences due to varying illumination can be much larger than differences between faces!

6.4 Fisherfaces/LDA

Training data: For eigenfaces distance of difference of illumination are within individual variance. Key idea: Find directions where ratio of between/within individual variance are maximized. Linearly project to basis where dimension with good signal to noise ratio are maximized.

Fisher linear discriminant analysis

Eigenimage method maximizes “scatter” within the linear subspace over the entire image set - regardless of classification task

E_opt = argmax_E (det(ERE*))

Fisher linear discriminant analysis: Maximize between-class scatter, while minimizing within-class scatter,

F_opt = argmax_F (det(FRW*)) / (det(FRF*))

R_B = sum_i N_i (u_i - u)(u_i - u)*

R_W = sum_{i=1,...,C} (Gamma_i - u_i)(Gamma_i - u_i)*

. N_i are the samples in class i and u_i is the mean in class i. Solution: Generalized eigenvectors w_i corresponding to the k largest eigenvalues {lambda_i | i = 1, ..., k}, i.e.

R_B w_i = lambda_i R_W w_i, i = 1, ..., k

Problem: within-class scatter matrix RW at most of rank L - c, hence usually singular. Apply KLT first to reduce dimension of feature space to L - c (or less), proceed with Fisher LDA in low-dimensional space.

Eigenfaces vs. Fisherfaces

Eigenfaces conserve energy but the two classes e.g. in 2D are no longer distinguishable. FLD (Fisher LDA) separates the classes by choosing a better 1D subspace. Fisher faces are much better in varying illuminations.

Varying illumination (FF) All images of same Lambertian surface with different illumination (without shadows) lie in a 3D linear subspace. Single point source at infinity

f(x, y) = a(x, y) (ell^T n(x, y)) L,

a(x, y) surface albedo, L light source intensity. Superposition of arbitrary number of point sources at infinity is still in same 3D linear subspace, due to linear superposition of each contribution to image. Fisher images can eliminate within-class scatter.

Appearance manifold approach

For everyobject, (1) sample the set of viewing conditions (2) use these images as feature vectors (3) apply a PCA over all the images (4) keep the dominant PCs (5) sequence of views for one object represent a manifold in space of projections (6) what is the nearest manifold for a given view?

Object-pose manifold Appearance changes projected on PCs (1D pose changes). Sufficient characterization for recognition and pose estimation.

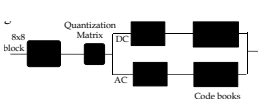
6.5 JPEG image compression

We don't resolve high frequencies too well...let's use this to compress images...JPEG!

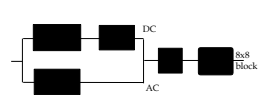
Concept Block-based discrete cosine transform (DCT)

JPEG Encoding and Decoding

Encoding:



Decoding:



DCT

A variant of discrete fourier transform: Real numbers, fast implementation. Block sizes: (1) small block: faster, correlation exists between neighboring pixels (2) better compression in smooth regions The first coefficient B(0, 0) is the DC component, the average intensity. The top-left coefficients represent low frequencies, the bottom right high frequencies.

Entropy Coding (Huffman code)

symbol	prob.	code	binary fraction
Z	0.5	1	0.1
Y	0.25	01	0.01
X	0.125	001	0.001
W	0.125	000	0.000

The code words, if regarded as a binary fraction, are pointers to the particular interval being coded. In Huffman code, the code words point to the base of each interval. The average code length is H = -sum p(s) log2 p(s) -> optimal.

7 Scale-space representations

From an original signal f(x) generate a parametric family of signals f^t(x), where fine-scale information is successively suppressed.

7.1 Image pyramid

Level 0: 1 x 1, Level 1: 2 x 2, Level 2: 4 x 4, Level J - 1: N/2 x N/2, Level J (base): N x N.

7.2 Applications

- (1) Search for correspondence: look at coarse scales, then refine with finer scales
- (2) Edge tracking: a “good” edge at a fine scale has parents at a coarser scale
- (3) Control of detail and computational cost in matching: e.g. finding stripes; terribly important in texture representation

7.3 Pyramids

Gaussian pyramid Smooth with gaussians, because “gaussian^2” = another gaussian. Progressively blurred and subsampled versions of the image. Adds scale invariance to fixed-size algorithms.

Laplacian Pyramid

<+> Shown the information added in gaussian pyramid at each spatial scale. Useful for noise reduction & coding.

Wavelet/QMIF

Bandpassed representation, complete, but with aliasing and some non-oriented subbands. Recursive application of a two-band filter bank to the lowpass band of the previous stage yields octave band splitting.

Steerable pyramid

Shown components at each scale and orientation separately. Non-aliased subbands. Good for texture and feature analysis.

7.4 Haar Transform

Two major sub-operations: (1) Scaling captures info at different frequencies (2) Translation captures info at different locations Can be represented by filtering and downsampling. Relatively poor energy compaction.

8 Optical Flow

In Visual Computing, people seem like to use

I_# = dI/d#, u = dx/dt, v = dy/dt

where the subindex means a derivative if and only if we are talking about I.

8.1 Applications

- 1 tracking 2 structure from motion 3 stabilization 4 compression 5 Mo-saicing

8.2 Brightness constancy

Definition of Optical Flow “Apparent motion of brightness patterns”. Ideally, the optical flow is the projection of the three-dimensional velocity vectors on the image.

Caution required 1 Uniform, rotating sphere OF = 0 2 No motion, but changing lighting OF != 0

8.3 Mathematical formulation

I(x, y, t) brightness at (x, y) at time t.

Brightness constancy assumption

I(dx/dt dt, y + dy/dt dt, t + dt) = I(x, y, t)

Optical flow constraint equation

dI = dI/dx dx/dt + dI/dy dy/dt + dI/dt = 0

8.4 The aperture problem

The motion of an edge seen through an aperture is (in some cases) inherently ambiguous. E.g. the edge is physically moving upwards, but the edge motion alone is consistent with many other possible motions, and in this case the edge e.g. appears to move diagonally.

8.5 Optical Flow meaning

Estimate of observed projected motion field. Not always well defined! Compare: 1 Motion Field (or Scene Flow), projection of 3-D motion field 2 Normal Flow: observed tangent motion 3 Optic Flow:

Apparent motion of the brightness pattern: Apparent motion of the brightness pattern (hopefully equal to motion field) Consider barber pole illusion

Planar motion Ideal motions of a plane, X, Y being the horizontal and vertical direction and Z normal to the image plane: 1 translation in X 2 translation in Y 3 rotation around Z 4 rotation around Y

8.6 Regularization: Horn & Schunck algorithm

The Horn-Schunck algorithm assumes smoothness in the flow over the whole image. thus, it tries to minimize distortions in flow and prefers solutions which show more smoothness. The flow is formulated as a global energy functional which is the sought to be minimized. This function is given for two-dimensional image streams as Eq. 5: The associated ELE are

dL/dx - d/dx dL/d(dx/dx) - d/dy dL/d(dy/dx) = 0,

dL/dy - d/dx dL/d(dx/dy) - d/dy dL/d(dy/dy) = 0.

this gives

dI/dx (dI/dx dx/dt + dI/dy dy/dt + dI/dt) - alpha^2 dx = 0,

dI/dx (dI/dx dx/dt + dI/dy dy/dt + dI/dt) - alpha^2 dy = 0.

with Delta = d^2/dx^2 + d^2/dy^2

Remarks

- 1 Coupled PDE solved using iterative methods and finite differences x = Delta x - lambda (dI/dx dx/dt + dI/dy dy/dt + dI/dt) dI/dx, y = Delta y - lambda (dI/dx dx/dt + dI/dy dy/dt + dI/dt) dI/dy.
- 2 More than two frames allow a better estimation of I.
- 3 Information spreads from corner-type patterns.
- 4 Errors at boundaries
- 5 Example of regularisation: selection principle for the solution of illposed problems.

8.7 Lucas-Kanade: Integrate over a Patch

The Lucas-Kanade method assumes that the displacement of the image contents between two nearby instants (frames) is small and approximately constant within a neighborhood of the point p under consideration. thus the optical flow equation can be assumed to hold for all pixels within a window centered at p. Namely, the local image flow (velocity) vector (x, y) must satisfy

dI(qk)/dx x + dI(qk)/dy y = -dI(qk)/dt

for k = 1, ..., n and qk the pixels inside the window. These equations can be written in matrix form

Av = b

where x = (x y)^T, v = (x y)^T and

A_ij = dI(qi)/dx_j, b_i = -dI(qi)/dt

Eq. 1 is overdetermined, so do compromise solution by the least squares principle

ple Eq. 6

8.8 Gradient-Based Estimation

Assume brightness constancy. Let f1(x) and f2(x) be 1D signals (images) at two time instants. Let f2 = f1(x - delta), where delta denotes translation.

f1(x) - f2(x) = delta f1'(x) + O(delta^2) ~ delta approx (f1(x) - f2(x)) / f1'(x)

Assume displaced image well approximated by first-order Taylor series

I(x + u, t + 1) approx I(x, t) + u . grad I(x, t) + I_t(x, t)

Insert Eq. 2 in Eq. 3 to get

grad I(x, t) . u + I_t(x, t) = 0.

This is called the gradient constraint equation.

8.9 Pyramid/Coarse-to-fine

- Limits of the (local) gradient method: 1 Fails when intensity structure within window is poor 2 Fails when displacement is large (typical operating range is motion of 1 pixel per iteration!). Linearization of brightness is suitable only for small displacements. 3 Brightness is no strictly constant in images. Actually less problematic than it appears, since we can pre-filter images to make them look similar.

8.10 Parametric motion models

Global miton models offer:

- 1 More constrained solutions than smoothness (Horn-Schunck)
- 2 Integration over a large area than a translation-only model can accommodate (Lucas-Kanade)

9 Questions

- I_comp = I_alpha I_a + (1 - I_alpha) I_b
- MAP, Maximum a posteriori detector.
- graph cuts
- Solve MRFs with graph cuts
- impulse respons t(-x, -y)
- Canny nonmaxima suppression
- Entropy Coding (Huffman code)
- Aperture problem: normal flow
- Lucas-Kanade: Iterative refinement/local gradient method
- Coarse-to-fine-estimation

A Big equations

F[h](u, v) = 1/(2*ell) * integral from -ell to ell of dx1 exp(-i*2*pi*u*x1) * integral from -inf to inf of dx2 delta(x2) exp(-i*2*pi*v*x2) = sinc(2*pi*ell*u)

E = double integral dx dy [(dI/dx dx/dt + dI/dy dy/dt + dI/dt)^2 + alpha^2 (||grad x|| + ||grad y||)^2]

v = (sum_i w_i I_x(q_i)^2 / sum_i w_i I_x(q_i) I_y(q_i), sum_i w_i I_x(q_i) I_y(q_i) / sum_i w_i I_y(q_i)^2)^-1

.(-sum_i w_i I_x(q_i) I_t(q_i) / -sum_i w_i I_y(q_i) I_t(q_i))

Eq. 1 is overdetermined, so do compromise solution by the least squares principle