For the complete discussion, see the corresponding slides & paper.

[1] Difference of Gradients (DoG) vs. Laplacian of Gradients (LoG)

» $DoG(I) = (I * G_{\sigma_1}) - (I * G_{\sigma_2})$
$\qquad\quad = I * (G_{\sigma_1} - G_{\sigma_2})$

If $I$ is a one-dimensional signal:

» LoG:

$$\Delta G = \nabla^2 G = \nabla \cdot \nabla G$$

$\nabla = $ divergence $\qquad \nabla G = $ gradient

$$\nabla = \left( \frac{\partial}{\partial x}, \frac{\partial}{\partial y} \right) \qquad \nabla G = \left( \frac{\partial G}{\partial x}, \frac{\partial G}{\partial y} \right) \quad \Big\} \text{ in 2D}$$
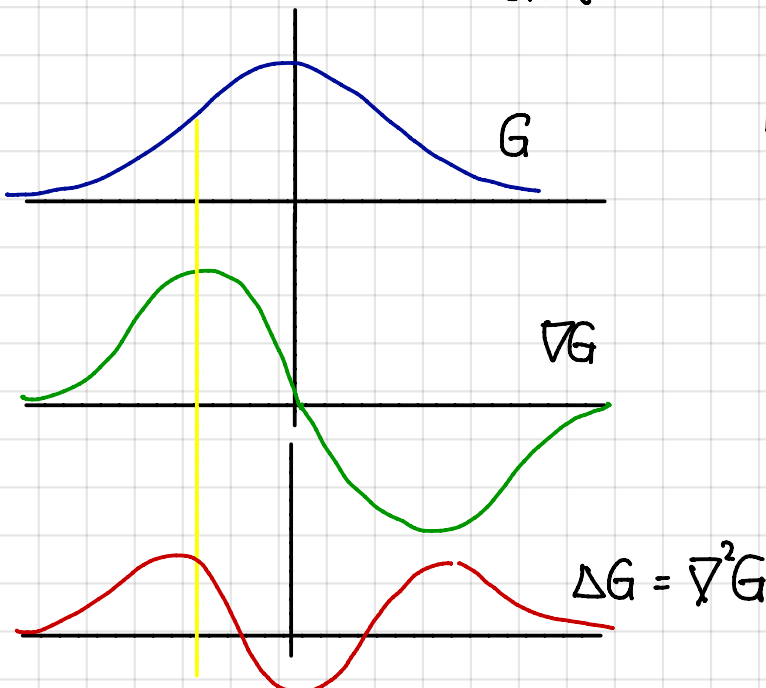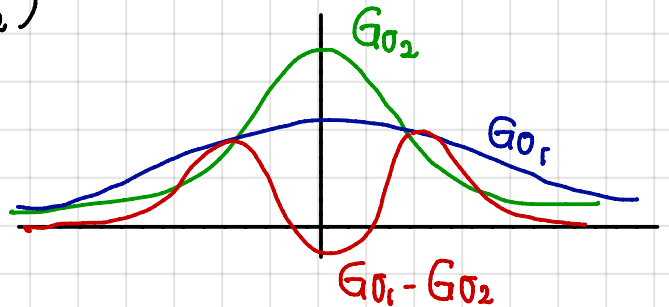
$$\Delta G = \nabla \cdot \nabla G = \frac{\partial}{\partial x}\left( \frac{\partial G}{\partial x} \right) + \frac{\partial}{\partial y}\left( \frac{\partial G}{\partial y} \right) = \frac{\partial^2 G}{\partial x^2} + \frac{\partial^2 G}{\partial y^2}$$

$G$

$\nabla G$

$$\Delta G = \nabla^2 G$$

$$LoG(I) = I * \Delta G$$
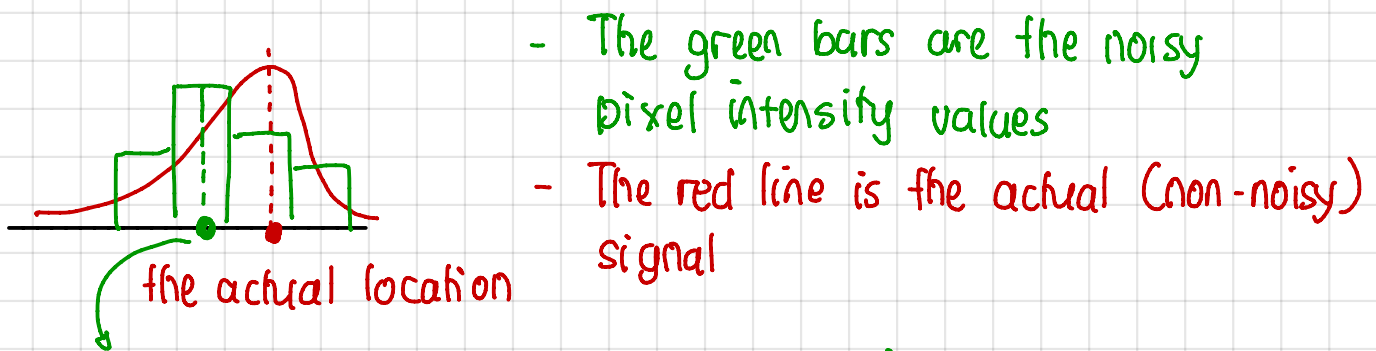
DoG is the approximation of LoG.

$$DoG(I) \approx LoG(I)$$

Motivation:

(1) Due to noise, the actual location of an extremum might be shifted.
(2) Even if we know the actual location, it can be not high enough (= low contrast).

Q : What does it mean by the actual location of an extremum?
How to calculate the actual location?

A :



- The green bars are the noisy pixel intensity values
- The red line is the actual (non-noisy) signal

the actual location

the predicted extremum location (using DoG)

» The actual signal is shifted by noise and discretization.

$\mathbb{D}$ = a 3x3 pixel patch where a keypoint located in the middle, and $\mathbb{D}$ is taken from one of the DoG images.
This 3x3 patch might be affected by noise.

$D(\bar{x})$ = the actual signal ; $\bar{x} = (x, y)$ → the actual location

Unfortunately, we don't know the actual signal location $(\bar{x})$.

Q : How to get $D(\bar{x})$, the signal we want to recover?
A : Through Taylor expansion:

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(x-a)}{2!}(x-a)^2 + \dots$$
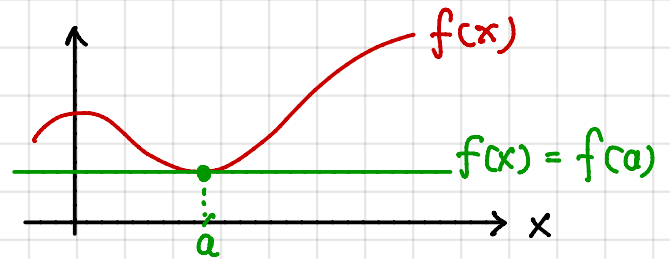
the unknown function

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(x-a)}{2!}(x-a)^2 + \ldots$$
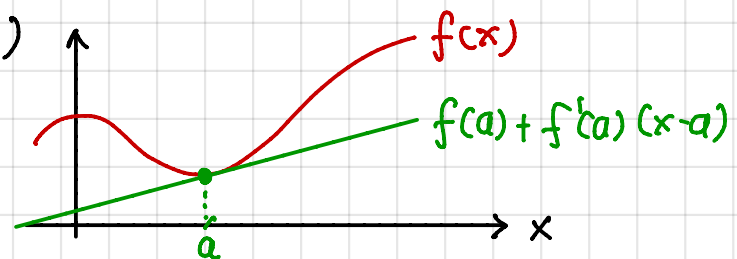
Meaning:

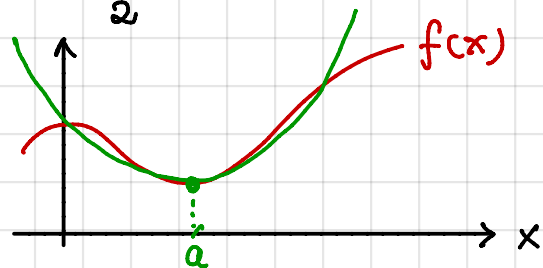1. When $x = a$ :   $f(x) = f(a)$

   The approximation is rough & basic



$f(x)$

$f(x) = f(a)$

2. When $f(x) = f(a) + f'(a)(x-a)$

   scalar values

   The approximation gets better.



$f(x)$

$f(a) + f'(a)(x-a)$

3. When $f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2}(x-a)^2$



$f(x)$

Notes:

- The approximation gets better when we include the higher order functions.

- 'a' is the point where we want to focus our approximation.

Using Taylor expansion in 2D:

$$D(\bar{x}) = D(0) + \frac{\partial D^T(\bar{x})}{\partial \bar{x}}\bigg|_{\bar{x}=0} \bar{x} + \bar{x}^T \frac{\partial^2 D(\bar{x})}{2\partial \bar{x}^2}\bigg|_{\bar{x}=0} \bar{x}$$

1×1     1×1        1×2  2×1     1×2   2×2     2×1

Note: → $D(\omega)$

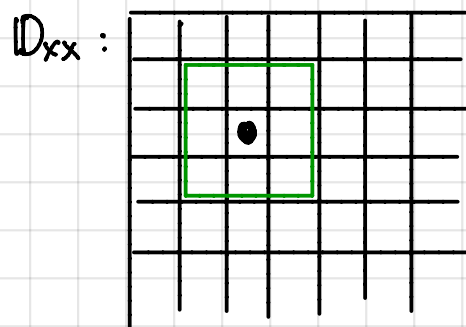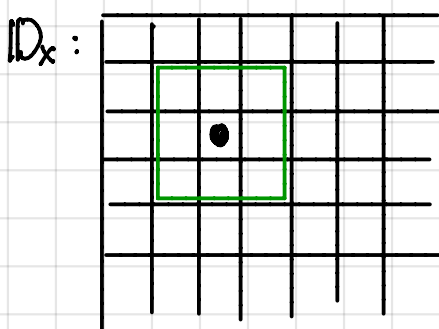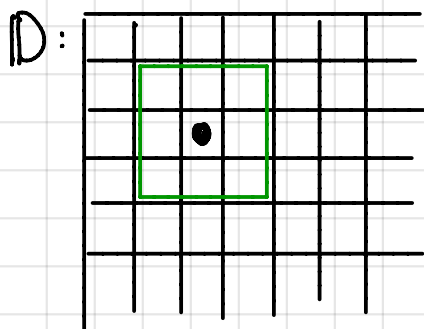$$D = \begin{matrix} -1 \\ 0 \\ +1 \end{matrix} \quad \begin{array}{ccc} & & \end{array}$$

-1   0   1

$$\frac{\partial D(\bar{x})}{\partial \bar{x}}\bigg|_{\bar{x}=0} = \begin{bmatrix} \partial D/\partial x \\ \partial D/\partial y \end{bmatrix}\bigg|_{\bar{x}=0}$$

$$\frac{\partial^2 D(\bar{x})}{\partial \bar{x}^2}\bigg|_{\bar{x}=0} = \begin{bmatrix} \dfrac{\partial^2 D}{\partial x^2} & \dfrac{\partial^2 D}{\partial x \partial y} \\ \dfrac{\partial^2 D}{\partial y \partial x} & \dfrac{\partial^2 D}{\partial y^2} \end{bmatrix}\bigg|_{\bar{x}=0}$$

Q: What are these $D$, $D_x$, $D_y$, $D_{xx}$, $D_{xy}$, $D_{yy}$?

A: $D$ is a DoG image (the output of the previous step), where there are a number of keypoints. For each of these keypoints, we take 3×3 pixels from $D$.

  $D_x$ (or $\frac{\partial D}{\partial x}$) & $D_y$ (or $\frac{\partial D}{\partial y}$) are the first derivative image of $D$.

$D$:           $D_x$:           $D_{xx}$:

we compute $D_y$, $D_{yy}$, and $D_{xy}$ (which is the same as $D_{yx}$) images in the same way.

To find the true extremum means:

$$\frac{\partial D(\bar{x})}{\partial \bar{x}} = 0 \qquad \rightarrow \text{ the output is a vector of } 2\times1.$$

$$\frac{\partial}{\partial \bar{x}} D(0) + \frac{\partial}{\partial \bar{x}}\left[ \frac{\partial \mathbb{D}^T}{\partial \bar{x}}\bigg|_{\bar{x}=0} \bar{x}\right] + \frac{\partial}{\partial \bar{x}}\left[ \bar{x}^T \frac{\partial^2 \mathbb{D}}{2\,\partial\bar{x}^2}\bigg|_{\bar{x}=0} \bar{x}\right] = 0$$
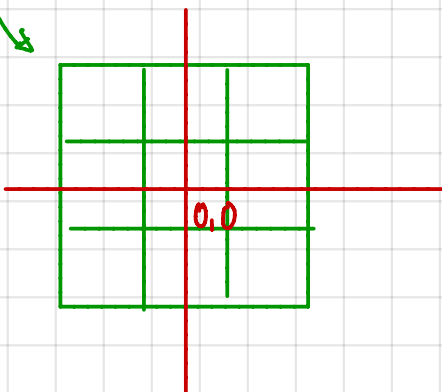
$$\underset{\longrightarrow \text{ constant} \longleftarrow}{}$$

$$0 + \frac{\partial \mathbb{D}(\bar{x})}{\partial \bar{x}}\bigg|_{\bar{x}=0} + \frac{\partial^2 \mathbb{D}(\bar{x})}{\partial \bar{x}^2}\bigg|_{\bar{x}=0} \bar{x} = 0$$

$$\boxed{\bar{x}^* = -\left( \frac{\partial^2 \mathbb{D}(\bar{x})}{\partial \bar{x}^2}\bigg|_{\bar{x}=0} \right)^{-1} \frac{\partial \mathbb{D}(\bar{x})}{\partial \bar{x}}\bigg|_{\bar{x}=0}}$$

If the offset of $\bar{x}^* > 0.5$

then : check the contrast based on the new location $\bar{x}^* \rightarrow \mathbb{D}(\bar{x}^*)$

else check the contrast based on $\mathbb{D}(0)$.
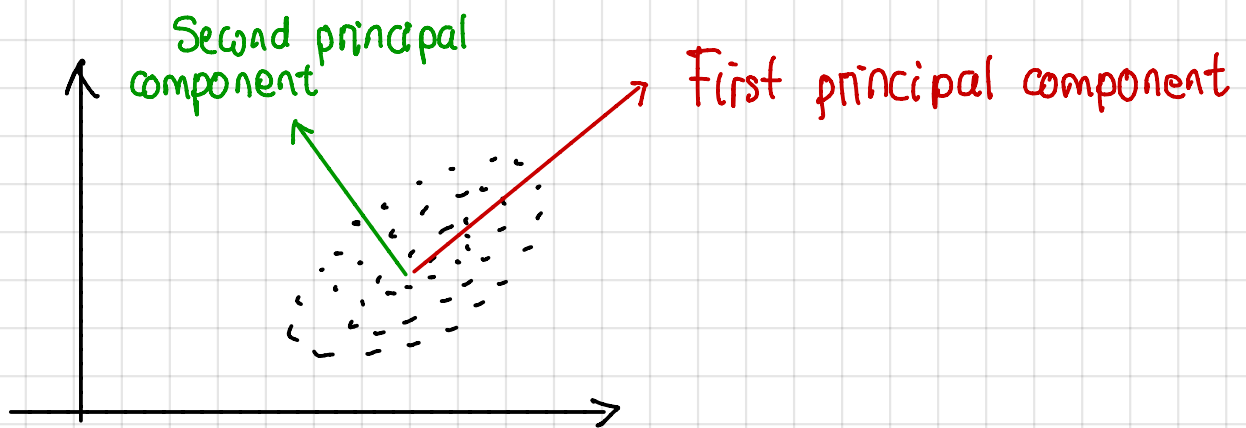


the distance from $(0,0)$ to the cell boundaries is $0.5$.

Based on the new location $\bar{x}^*$, we check if the extremum is high enough:

$$D(\bar{x}^*) = D[0] + \left.\frac{\partial}{\partial \bar{x}}(D^T(\bar{x}))\right|_{\bar{x}=0} \bar{x}^*$$

$$\underset{1\times1}{} \qquad \underset{1\times2}{} \qquad \underset{2\times1}{}$$

If $|D(\bar{x}^*)| < \underline{0.03}$ then reject!

it assumes the image intensity of the input is between $0 \sim 1$ (instead of $0 \sim 255$).

---

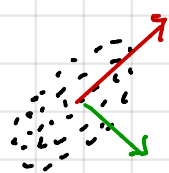# [6] First and Second Principal Components of Curvatures



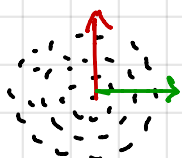$\alpha$ = The eigen value of the first principal component

$\beta$ = The eigen value of the second principal component

$\alpha \geqslant \beta$
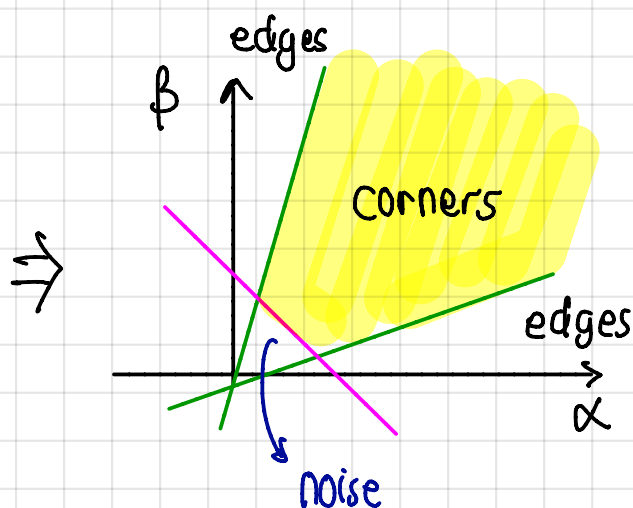
Eigenvalues ($\alpha$ & $\beta$) can indicate whether a pixel cluster is edge or corner:



$\alpha \gg \beta \rightarrow$ edge

$\alpha \approx \beta \rightarrow$ corner

Unlike the above illustration, our data is only 3×3 pixels (9 pixels), which is too sparse to calculate the principal axes.

↓

Solution: to use the ratio of $\alpha$ & $\beta$:

1. Hessian matrix:   $H = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{bmatrix}$

2.  $Tr(H) = D_{xx} + D_{yy} = \alpha + \beta$

   $Det(H) = D_{xx} D_{yy} - D_{xy}^2 = \alpha\beta$

3.  $\dfrac{Tr^2(H)}{Det(H)} = \dfrac{(\alpha+\beta)^2}{\alpha\beta}$     ∴ define: $\alpha = r\beta$

   $= (r\beta+\beta)^2 / r\beta^2 = \dfrac{(r+1)^2}{r}$

   if $\alpha = \beta$   $\rightarrow$ $r = 1$ :   $Tr^2(H)/Det(H) = 4$        $\rightarrow$ corner

   if $\alpha = 2\beta$   $\rightarrow$ $r = 2$ :   $Tr^2(H)/Det(H) = 9/2 = 4.5$

   if $\alpha = 10\beta$ $\rightarrow$ $r = 10$ :   $Tr^2(H)/Det(H) = 121/10 = 12.1$   $\rightarrow$ edge

Therefore:   if $Tr^2(H)/Det(H) < 12.1$
             then retain the keypoints, else reject them.

① Orientation assignment
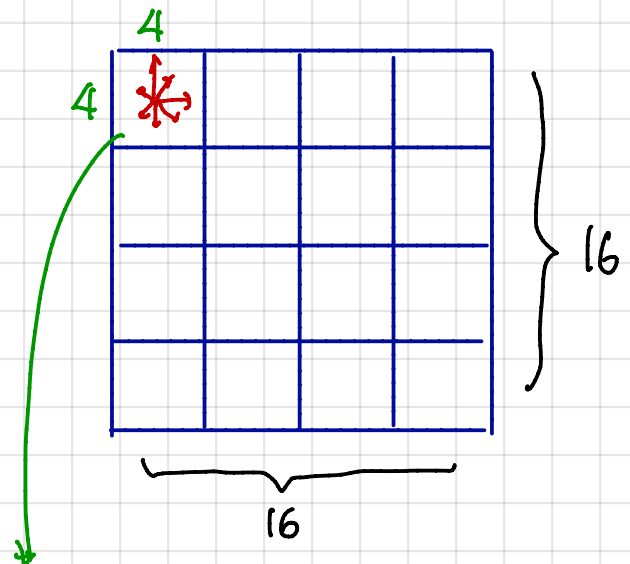
↓

Goal: to make the descriptor
invariant to image rotation.

↓

1. Extract 16×16 pixels surrounding
   a keypoint.
2. Create a histogram of orientations
   with 36 bins covering $360°$.

3. Choose the highest peak in the
   histogram, and any peaks above
   $80\%$, to calculate the orientation
   normalization.

② Descriptor:

one keypoint generates one
descriptor, which has a length
of 128.

↓

These 128 numbers are obtained
from 16×16 pixels:

4

4

16

16

For each block, we compute the histogram
of gradients with 8 bins (= orientations)

↓

Hence, for 1 block of 4×4 pixels we have
8 numbers. Thus, in total we have:

$$4 \times 4 \times 8 = 128.$$