



midterm review

16-385 Computer Vision
Carnegie Mellon University (Kris Kitani)

Filtering and
convolution

filter

I	I	I
I	I	I
I	I	I

A 10x10 grid representing an output layer. The grid consists of 100 empty cells arranged in 10 rows and 10 columns. A single cell at the top-left corner, corresponding to index $(1, 1)$, is highlighted with a thick blue border.

$$h[m, n] = \sum_{k,l} g[k, l] f[m + k, n + l]$$

output k, l filter image (signal)

Filtering vs Convolution

filtering
(cross-correlation)

$$h = g \otimes f$$

convolution

$$h = g \star f$$

output **filter** **image**

$$h[m, n] = \sum_{k,j} g[k, l] f[m + k, n + l]$$

What's the
difference?

$$h[m, n] = \sum_{k,j} g[k, l] f[m - k, n - l]$$

Filtering vs Convolution

filtering
(cross-correlation)

$$h = g \otimes f$$

convolution

$$h = g \star f$$

output **filter** **image**

$$h[m, n] = \sum_{k,j} g[k, l] f[m + k, n + l]$$

filter flipped
vertically and
horizontally

$$h[m, n] = \sum_{k,j} g[k, l] f[m - k, n - l]$$

Filtering vs Convolution

filtering
(cross-correlation)

$$h = g \otimes f$$

convolution

$$h = g \star f$$

output filter image

$$h[m, n] = \sum_{k,j} g[k, l] f[m + k, n + l]$$

filter flipped
vertically and
horizontally

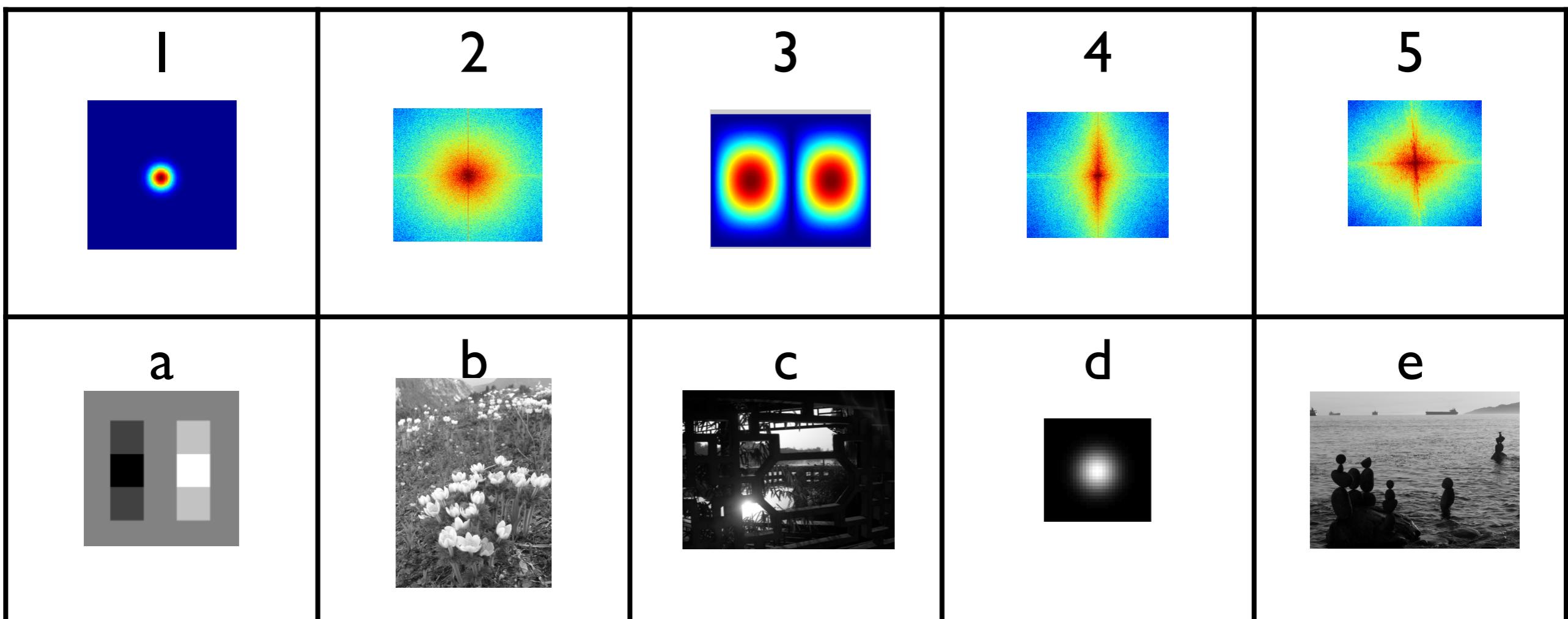
$$h[m, n] = \sum_{k,j} g[k, l] f[m - k, n - l]$$

Suppose g is a Gaussian filter.
How does convolution differ from filtering?

Recall...

$\frac{1}{16}$	1	2	1
2	4	2	
1	2	1	

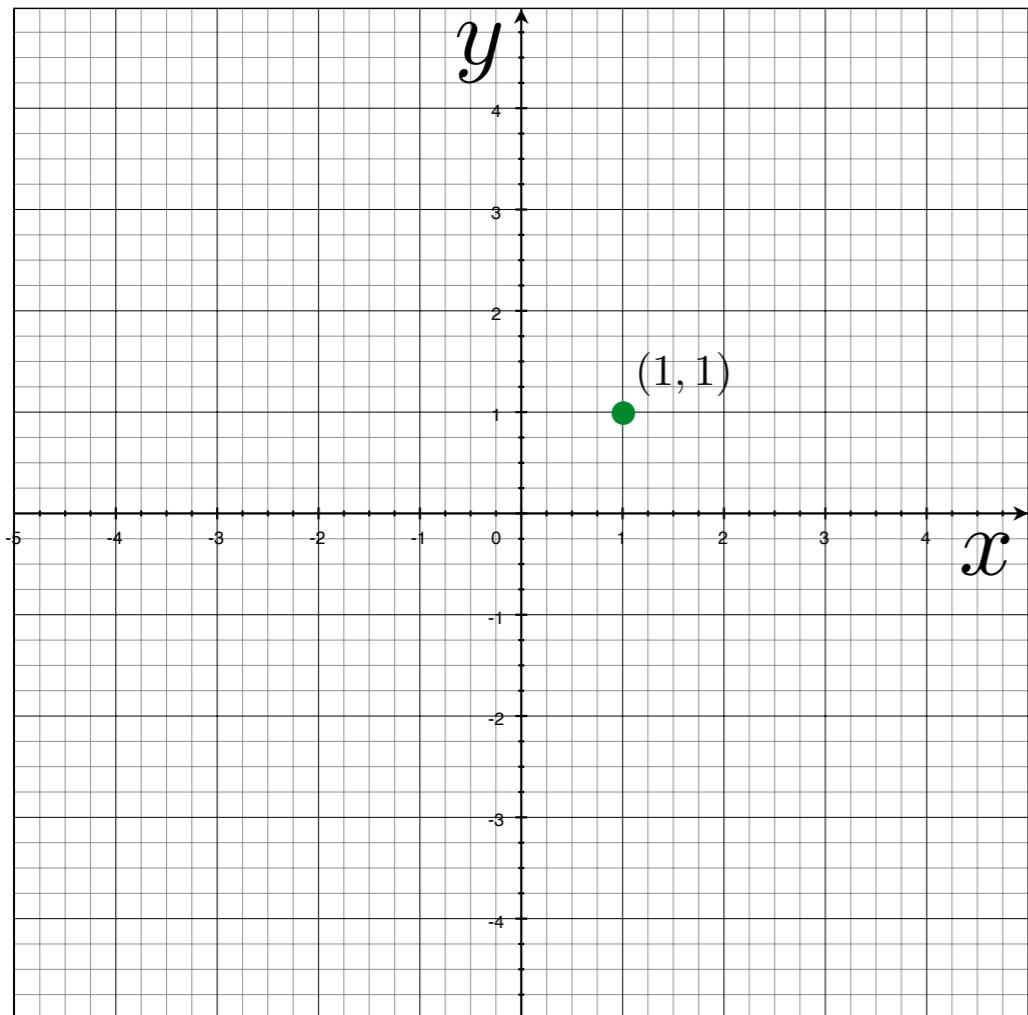
Match the image to the Fourier magnitude image:



Hough Transform

Image and parameter space

variables
 $y = mx + b$
parameters



a point becomes
a line

variables
 $y - mx = b$
parameters

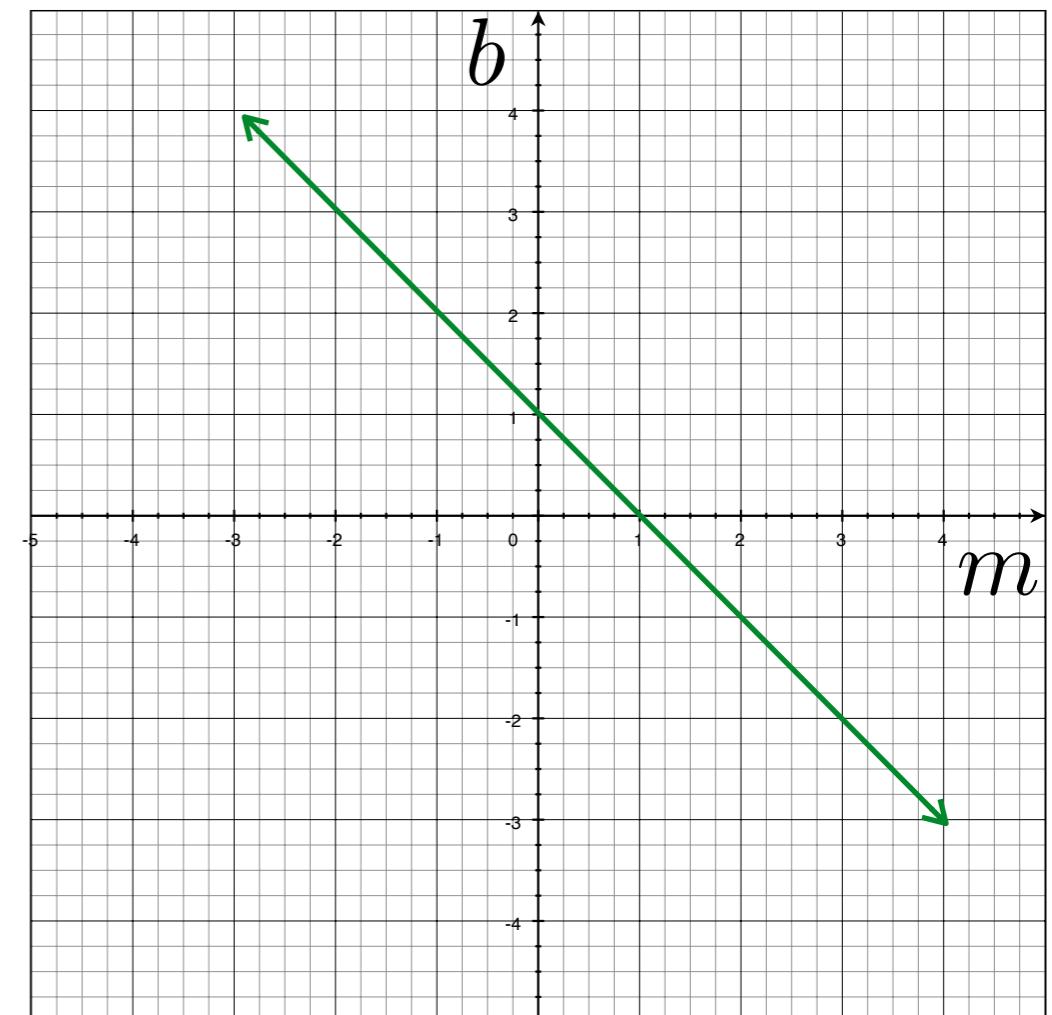


Image space

Parameter space

Problems with parameterization

What's wrong with the parameterization (m, c) ?

How big does the accumulator need to be?

$$A(m,c)$$

Problems with parameterization

What's wrong with the parameterization (m, c) ?

How big does the accumulator need to be?

$$A(m,c)$$

The space of m is huge!

$$-\infty < m < \infty$$

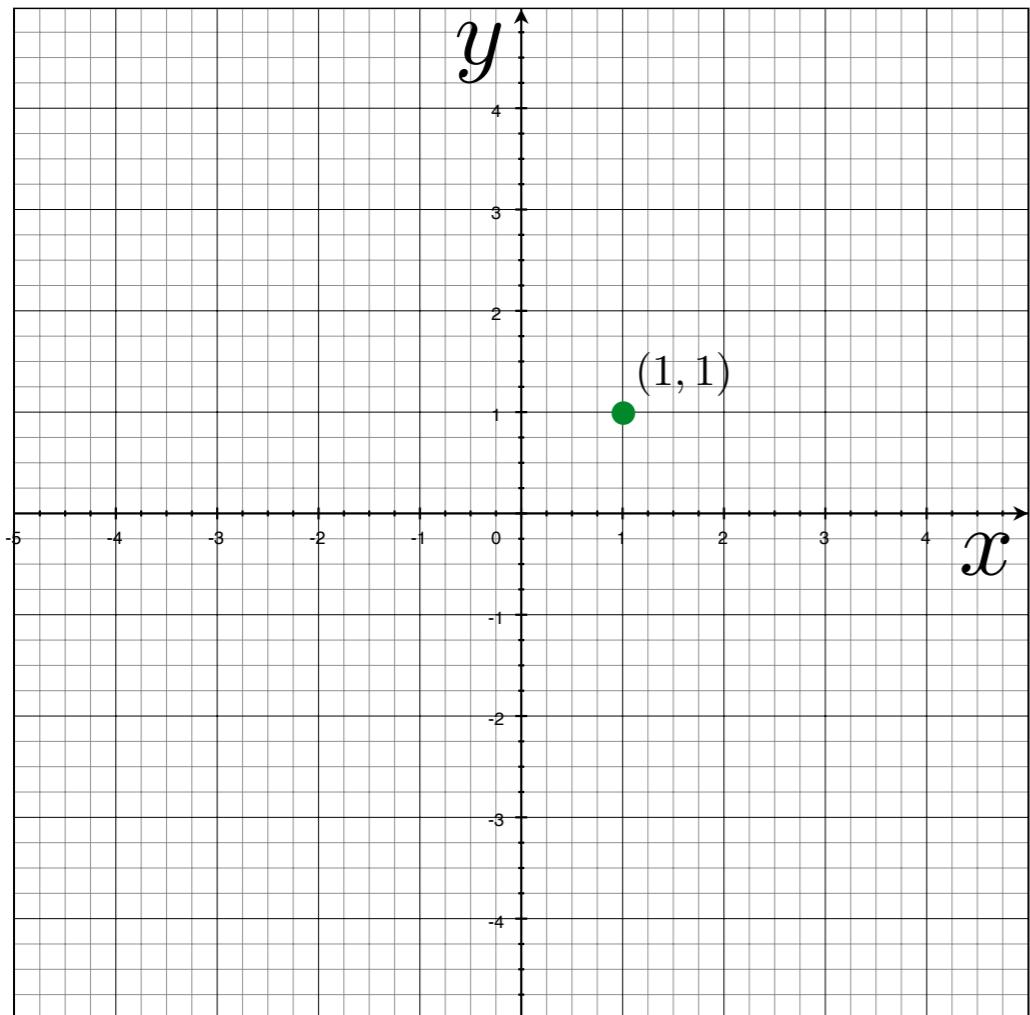
The space of c is huge!

$$-\infty < c < \infty$$

Image and parameter space

$$y = mx + b$$

variables
parameters



a point becomes
a line

$$x \cos \theta + y \sin \theta = \rho$$

parameters
variables

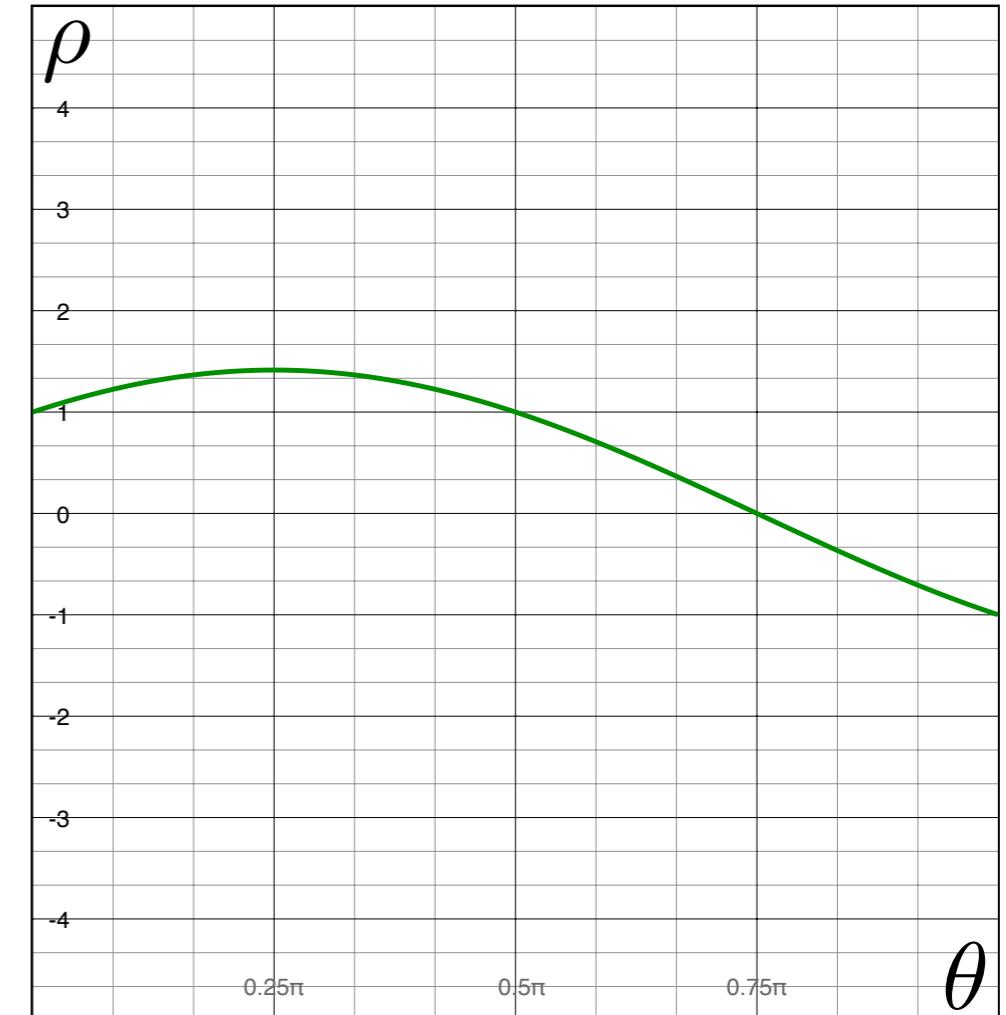
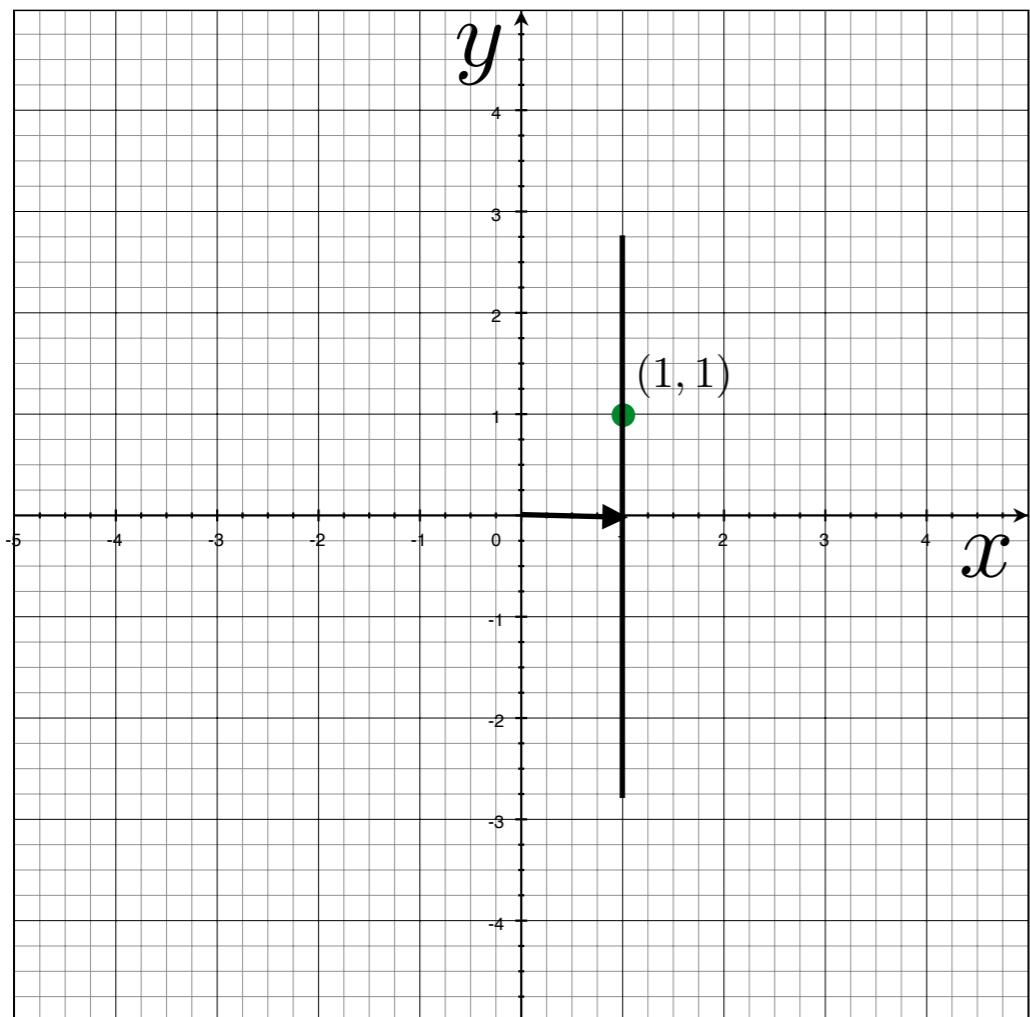


Image space

Parameter space

Image and parameter space

variables
 $y = mx + b$
parameters



a line
becomes
a point

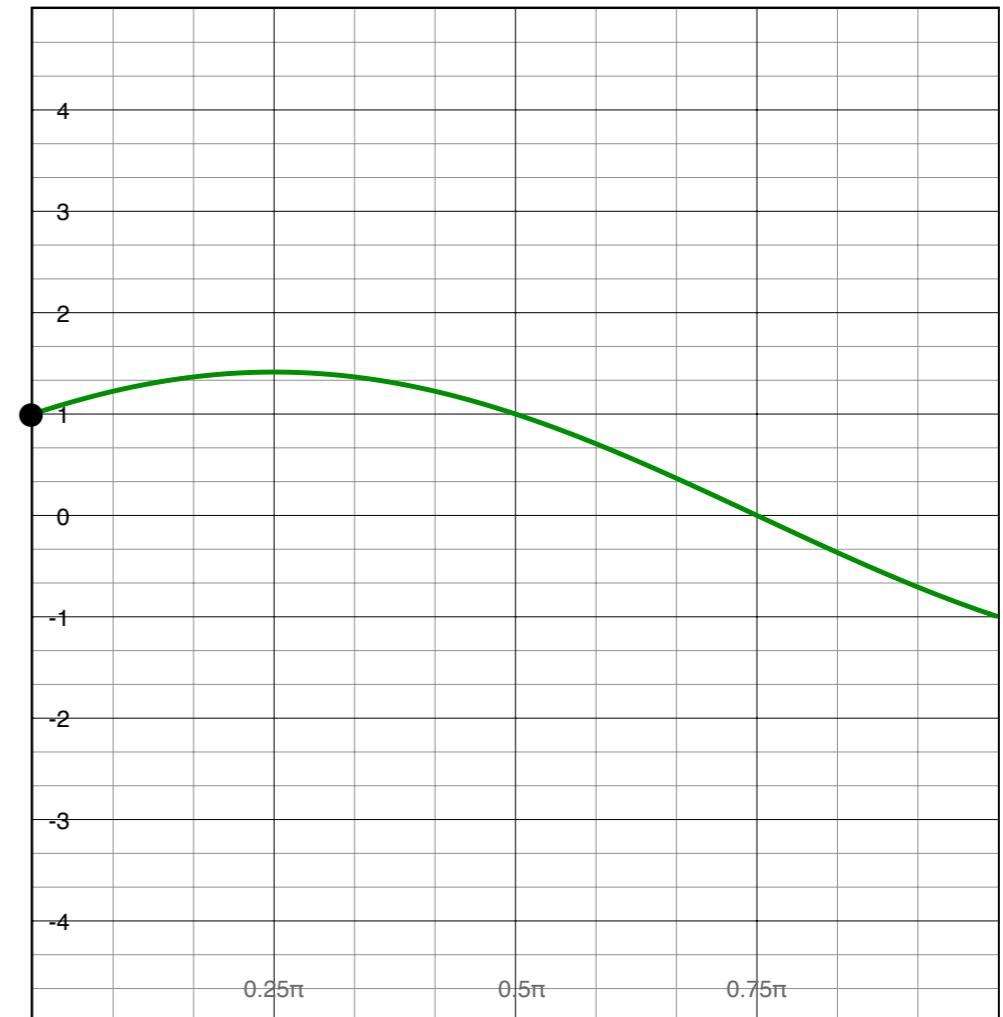


Image space

Parameter space

Image and parameter space

variables
 $y = mx + b$
parameters

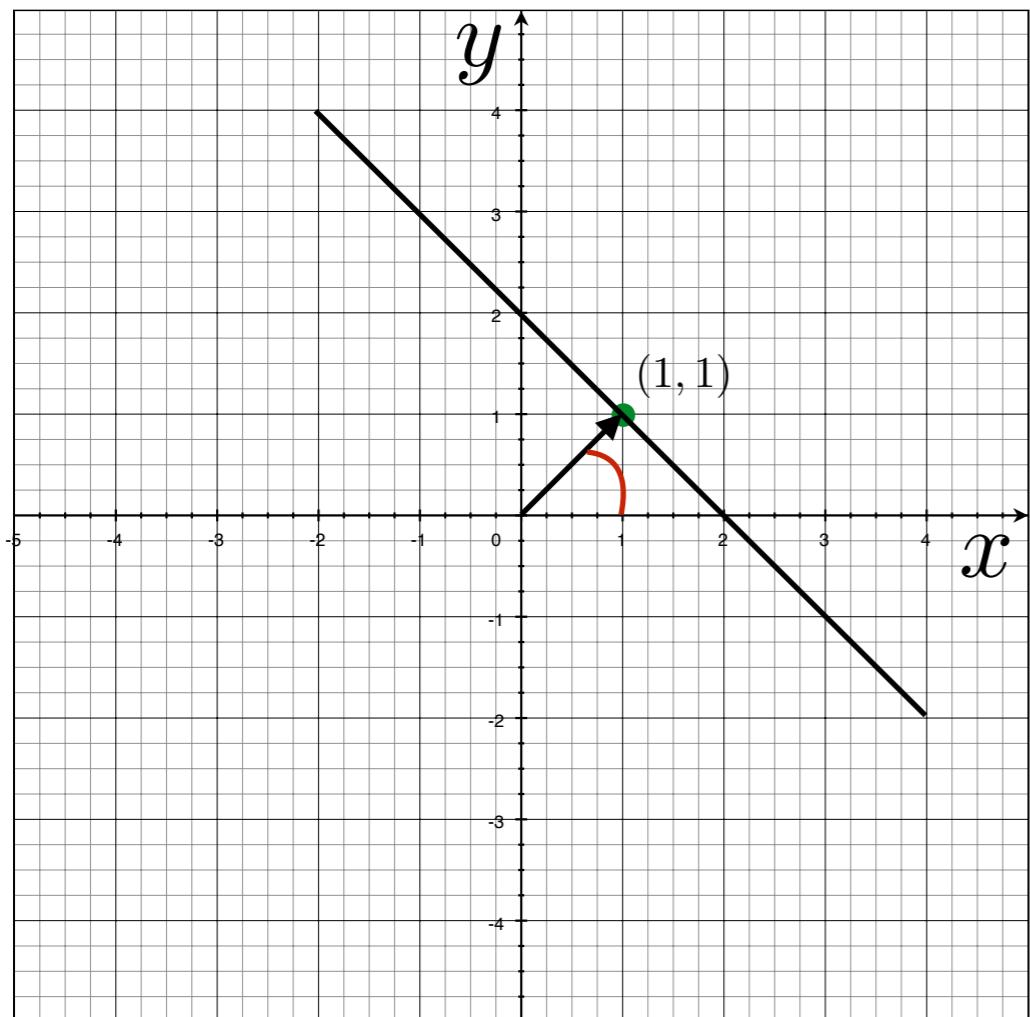
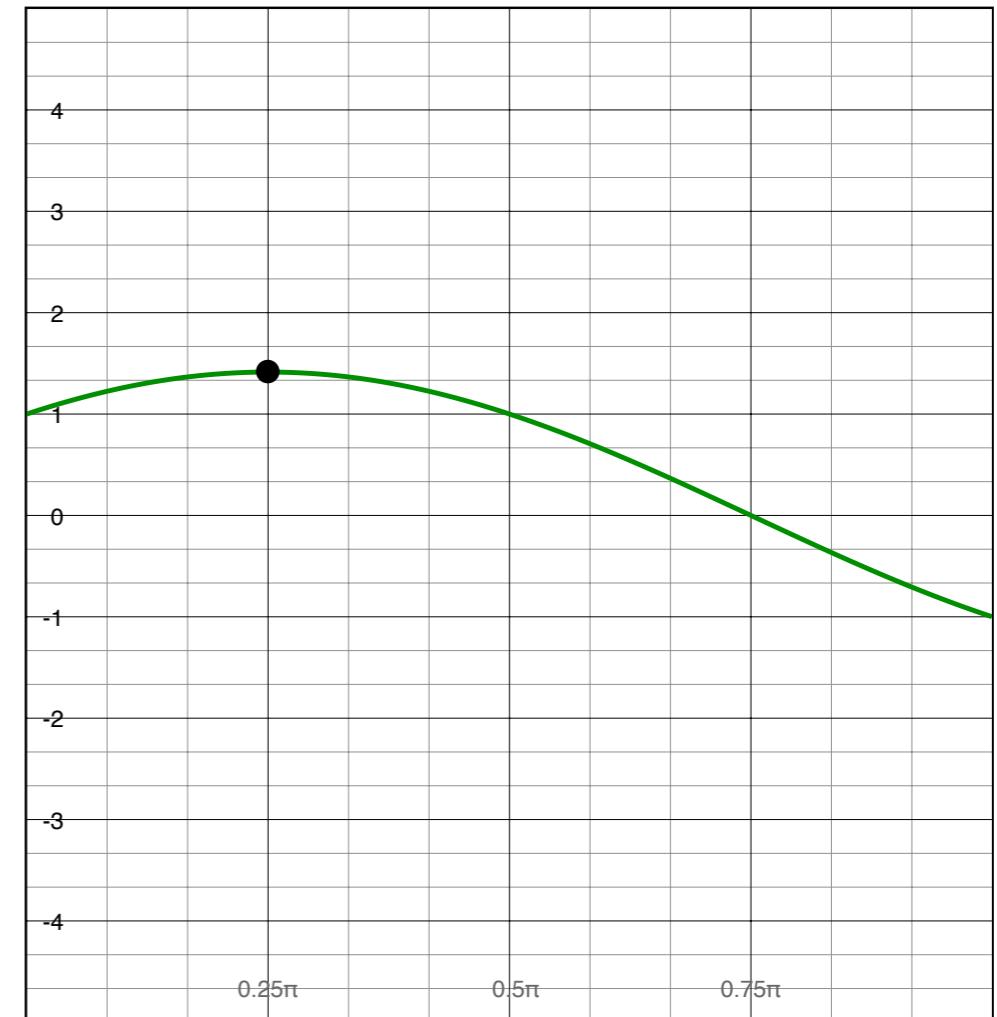


Image space

a line becomes a point



Parameter space

Image and parameter space

variables
 $y = mx + b$
parameters

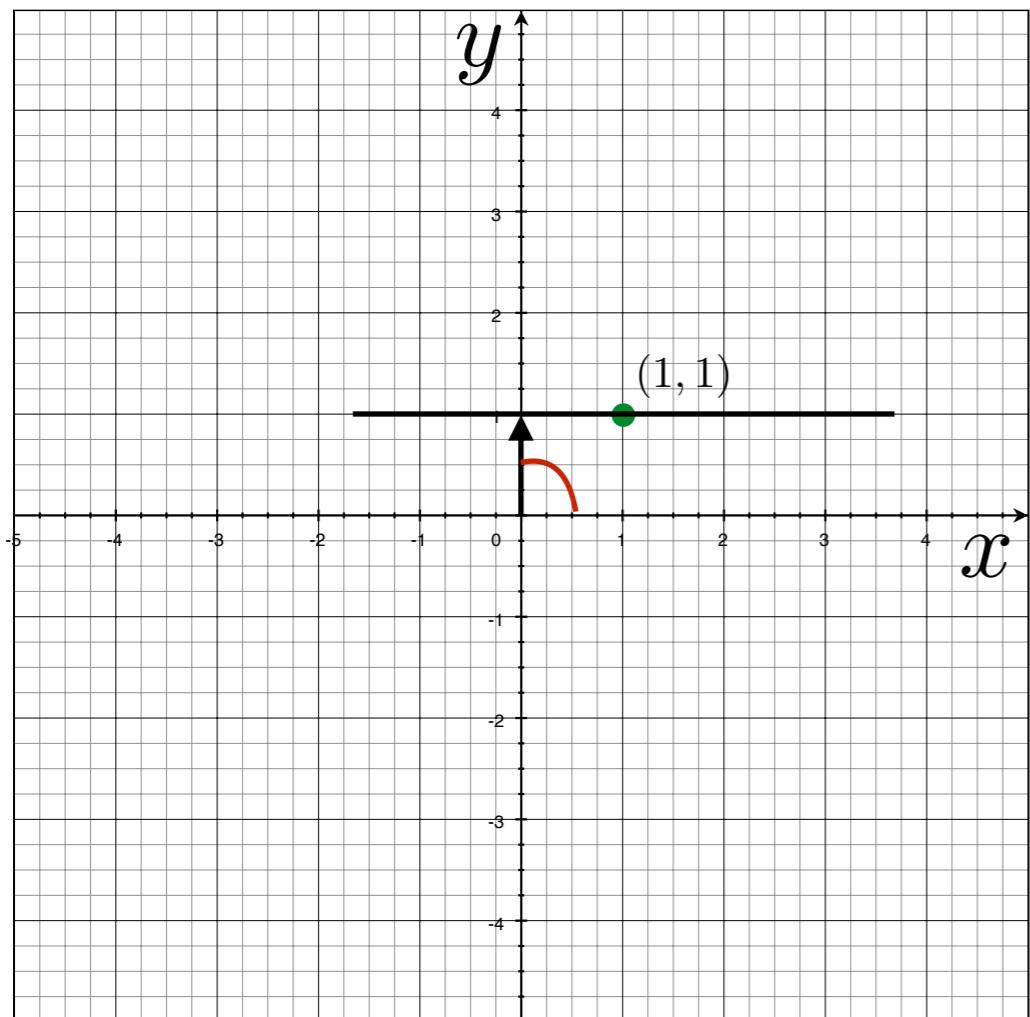
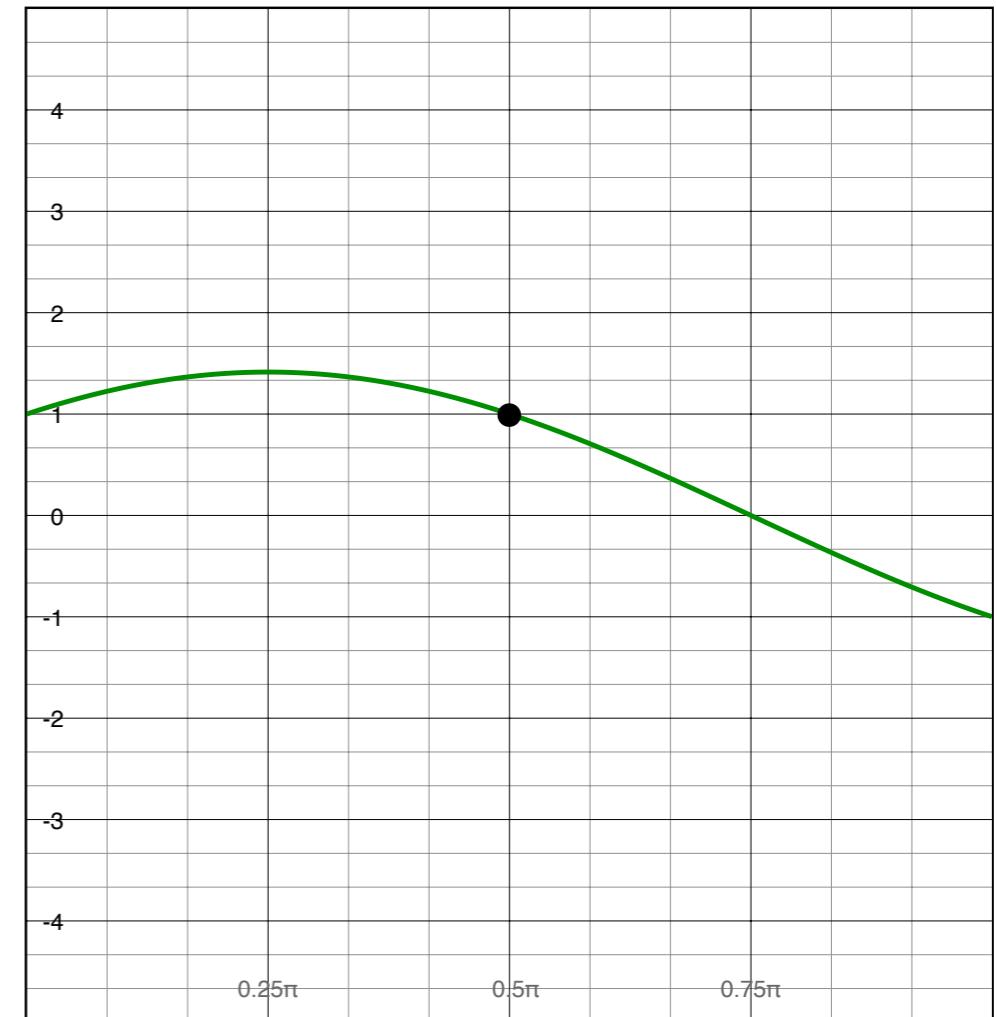


Image space

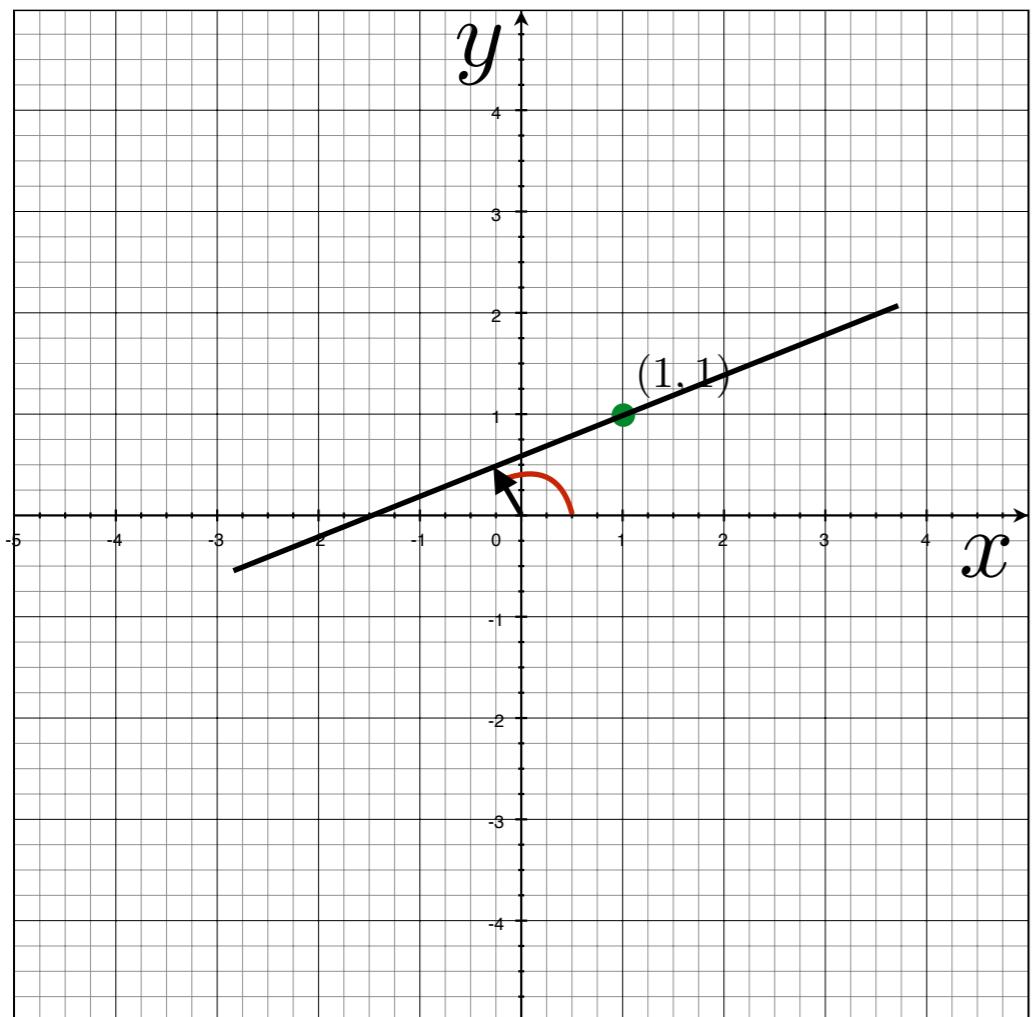
a line becomes a point



Parameter space

Image and parameter space

variables
 $y = mx + b$
parameters



a line
becomes
a point

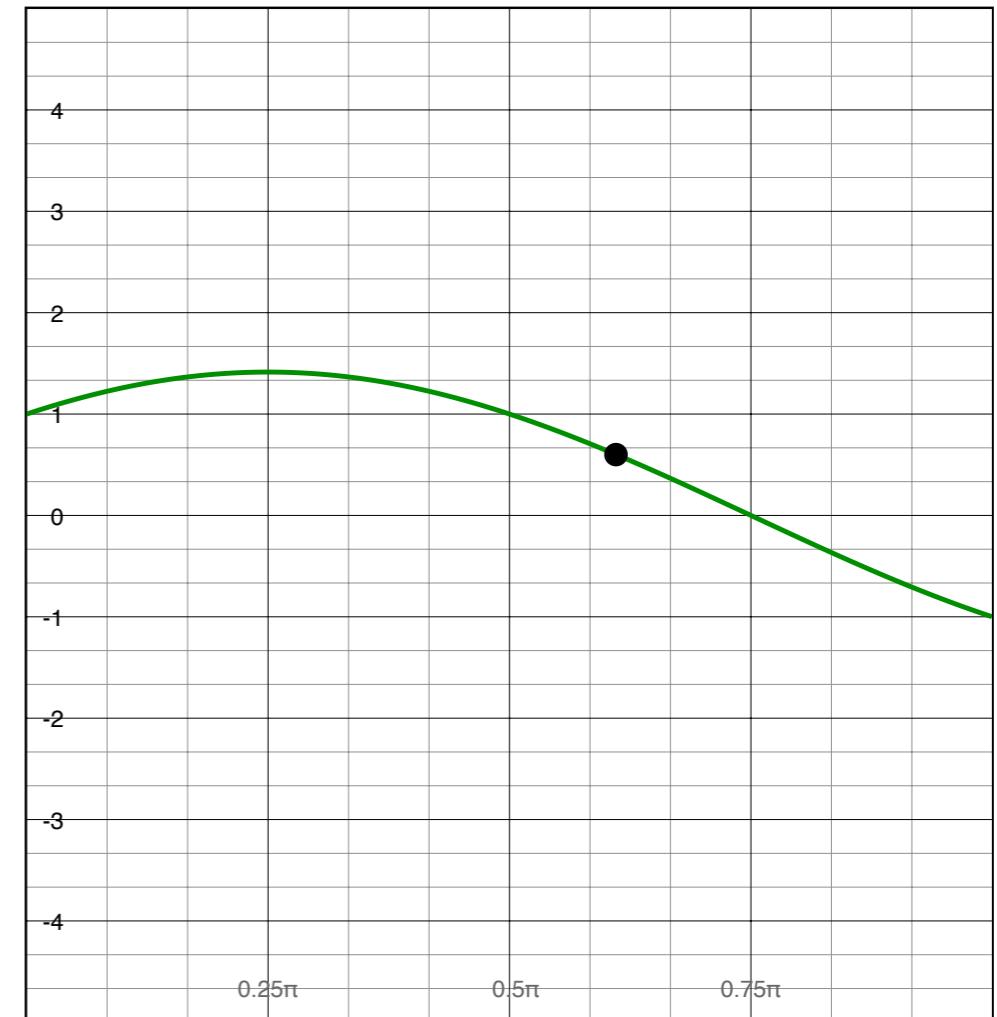
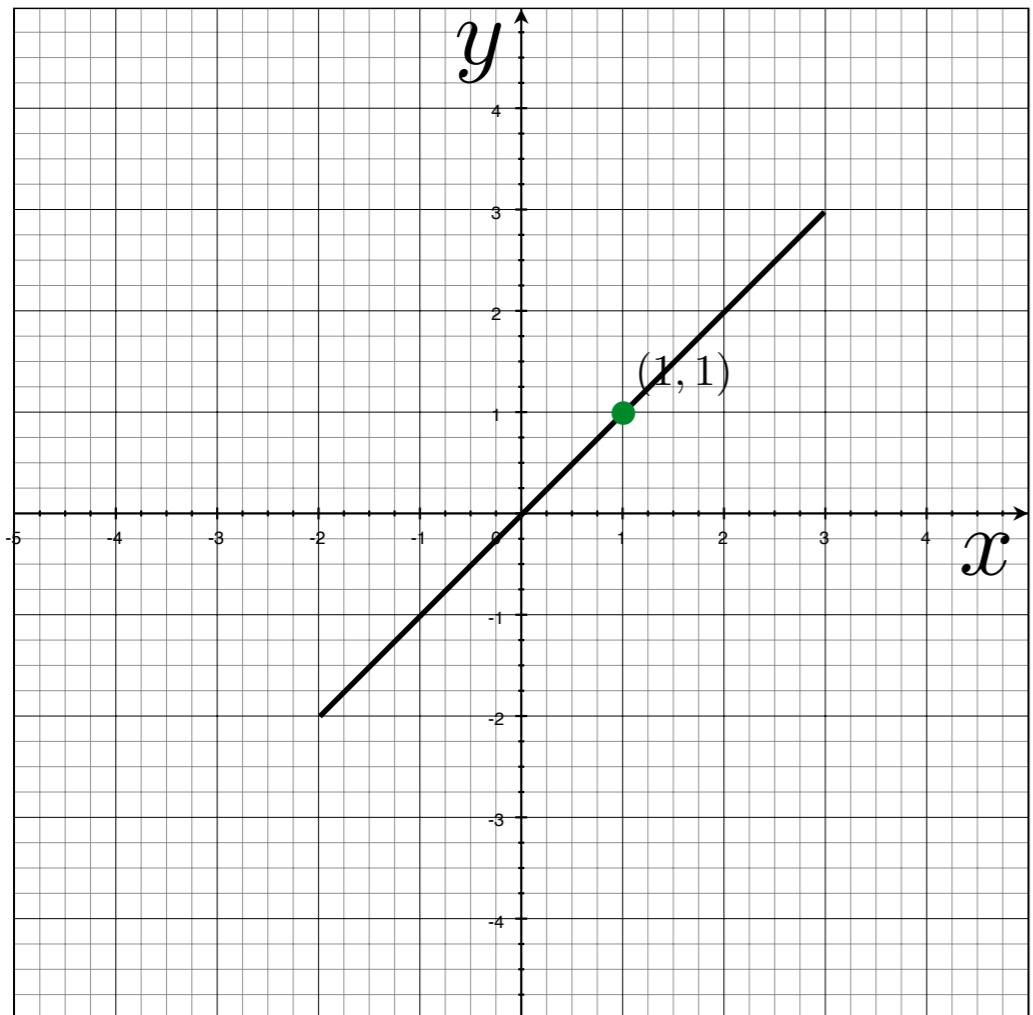


Image space

Parameter space

Image and parameter space

variables
 $y = mx + b$
parameters



a line
becomes
a point

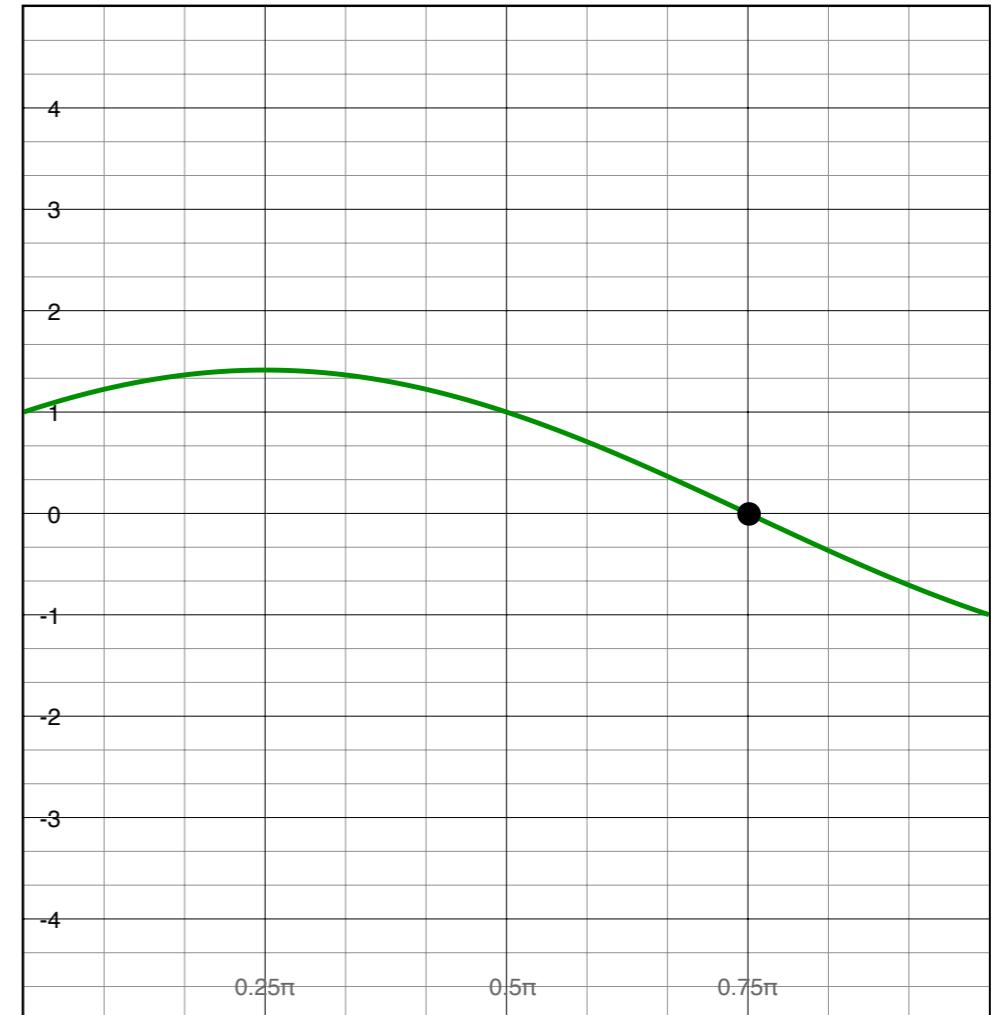
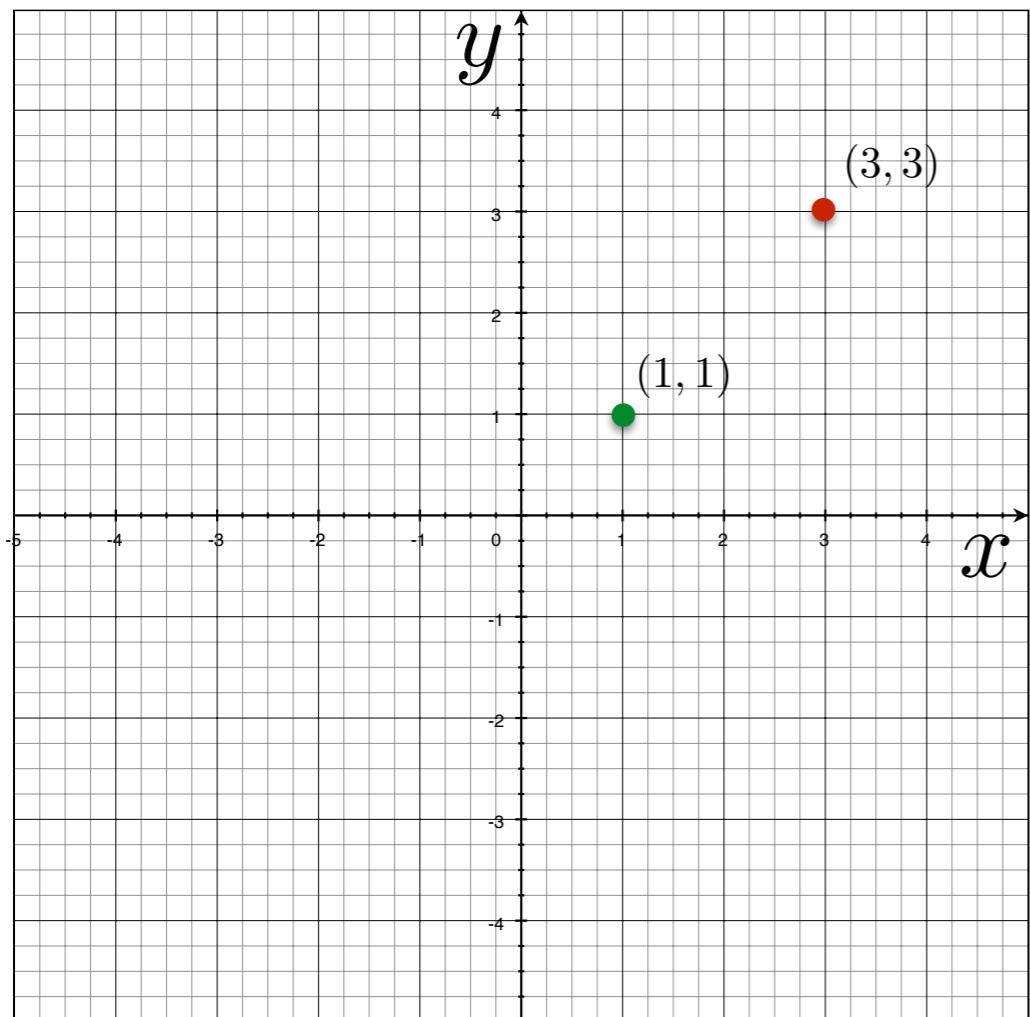


Image space

Parameter space

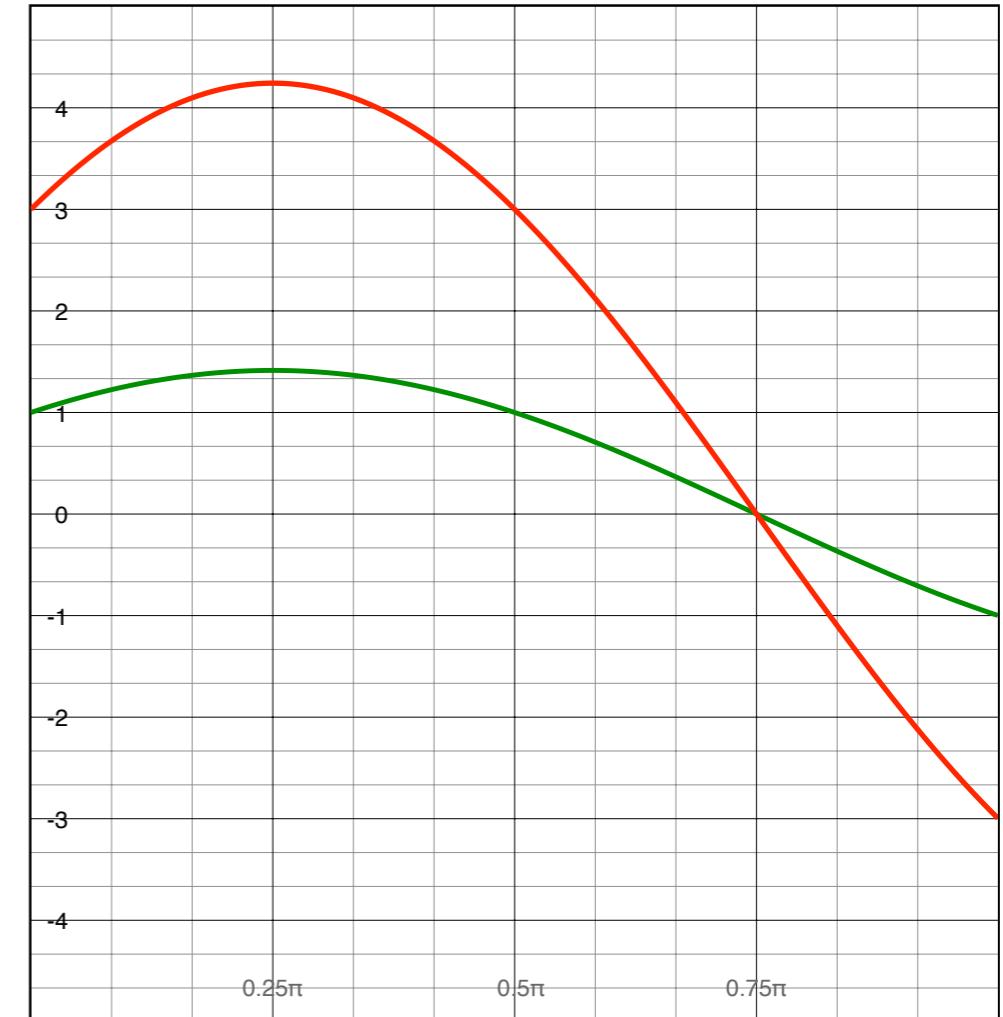
Image and parameter space

variables
 $y = mx + b$
parameters



two points
become
?

Image space

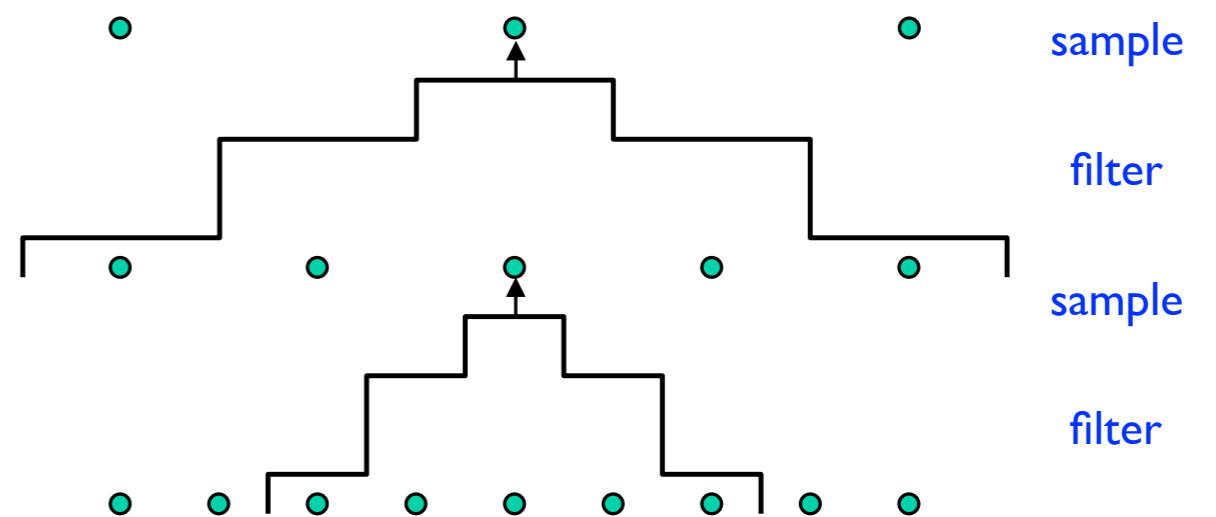


Parameter space

Image Pyramids

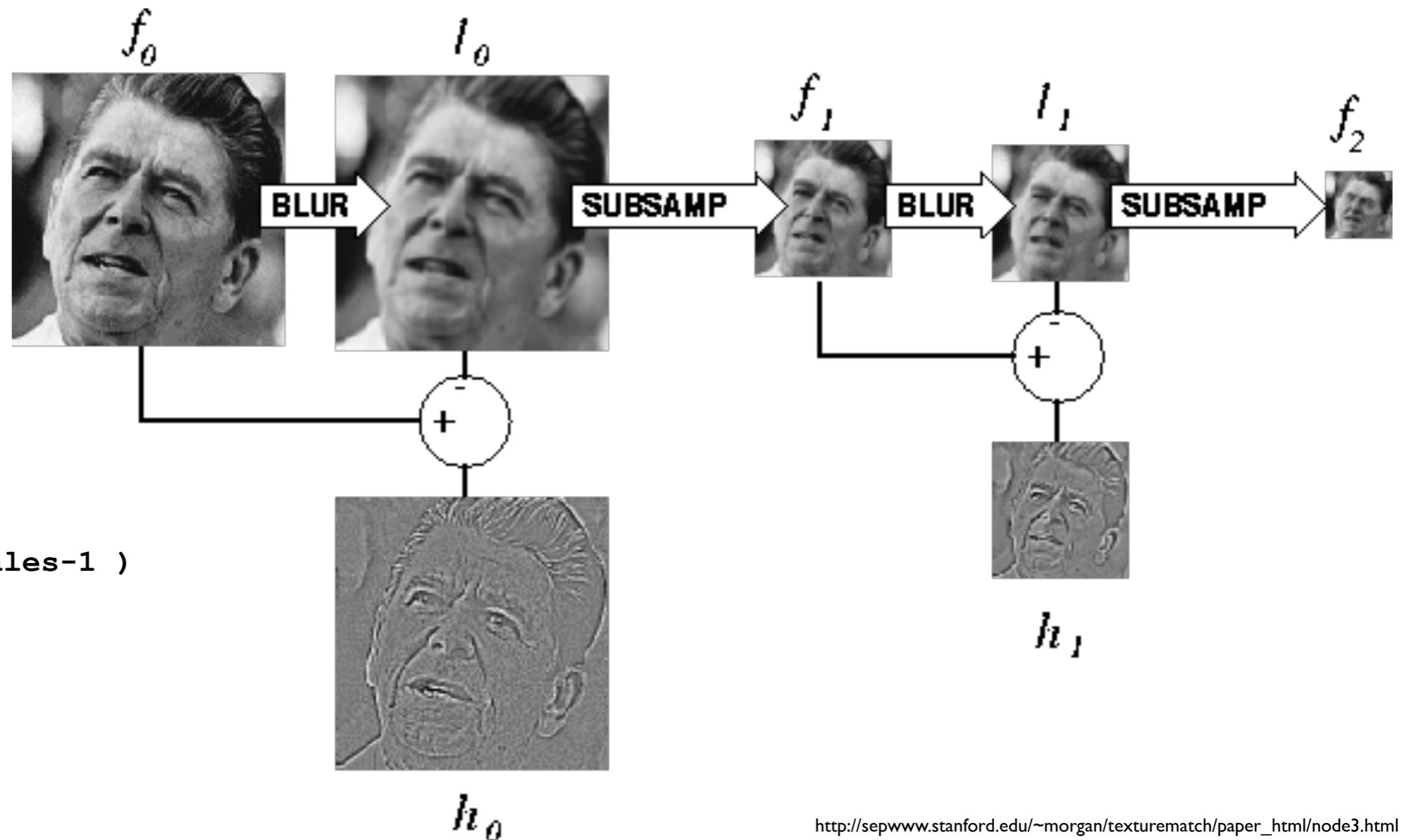
Constructing a Gaussian Pyramid

```
repeat  
    filter  
    subsample  
until min resolution reached
```



Whole pyramid is only $4/3$ the size of the original image!

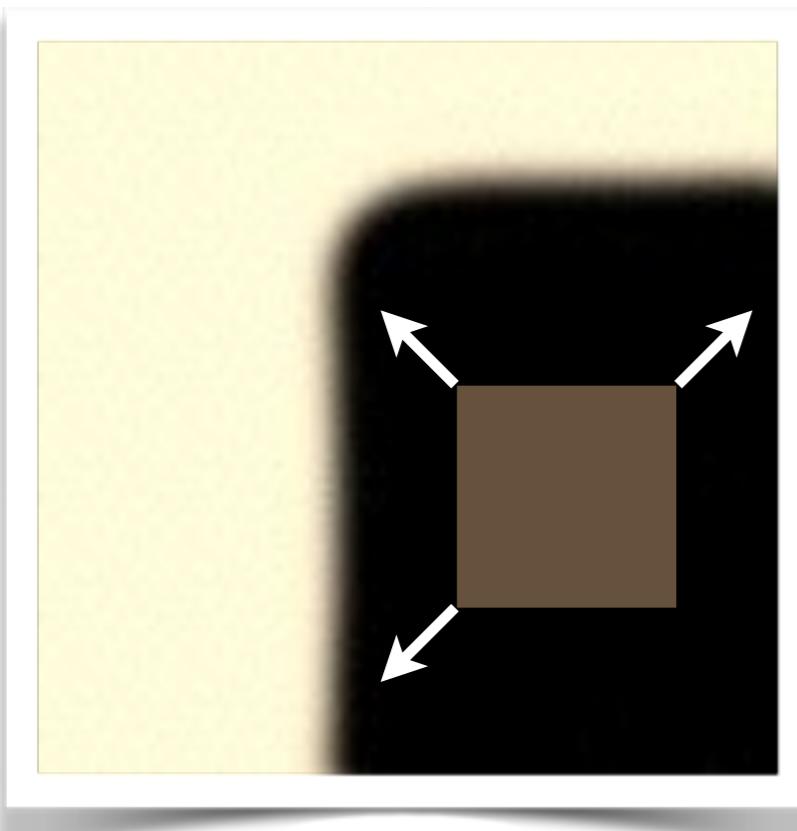
Constructing the Laplacian Pyramid



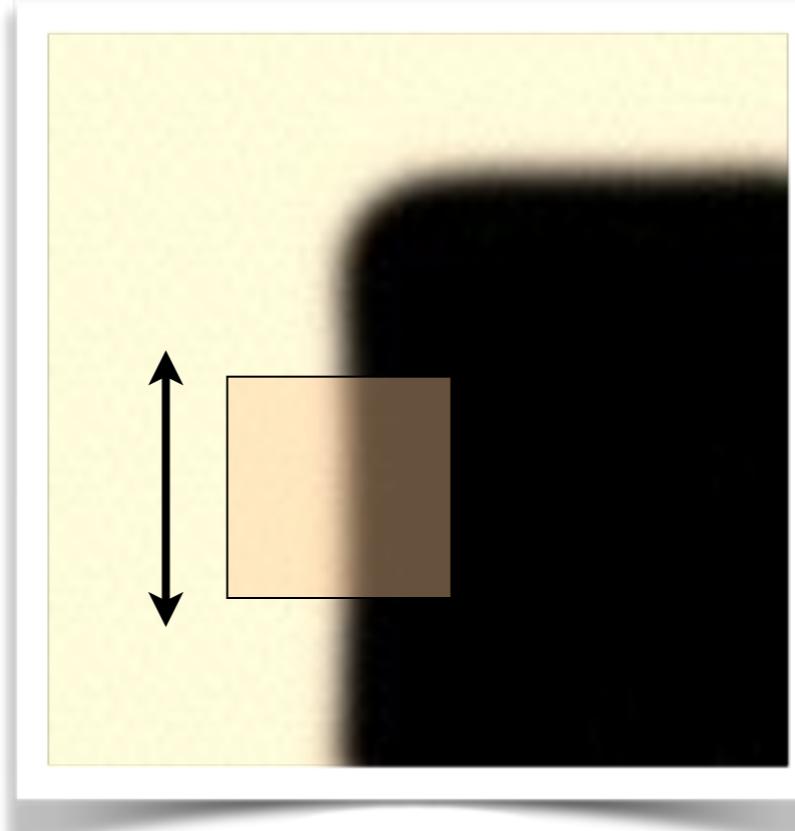
Feature Detector & Descriptor

Easily recognized by looking through a small window

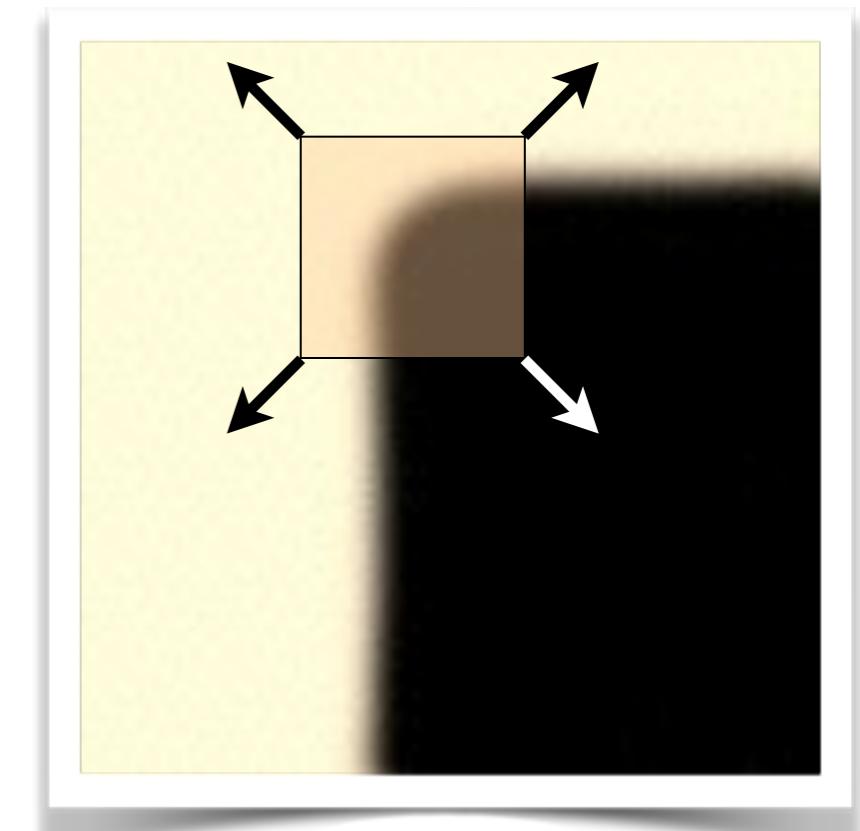
Shifting the window should give large change in intensity



“flat” region:
no change in all
directions



“edge”:
no change along the edge
direction



“corner”:
significant change in all
directions

Error function approximation

Change of intensity for the shift $[u, v]$:

$$E(u, v) = \sum_{x, y} w(x, y) [I(x + u, y + v) - I(x, y)]^2$$

Second-order Taylor expansion of $E(u, v)$ about $(0, 0)$
(bilinear approximation for small shifts):

$$E(u, v) \approx E(0, 0) + [u \ v] \begin{bmatrix} E_u(0, 0) \\ E_v(0, 0) \end{bmatrix} + \frac{1}{2} [u \ v] \begin{bmatrix} E_{uu}(0, 0) & E_{uv}(0, 0) \\ E_{uv}(0, 0) & E_{vv}(0, 0) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

first derivative

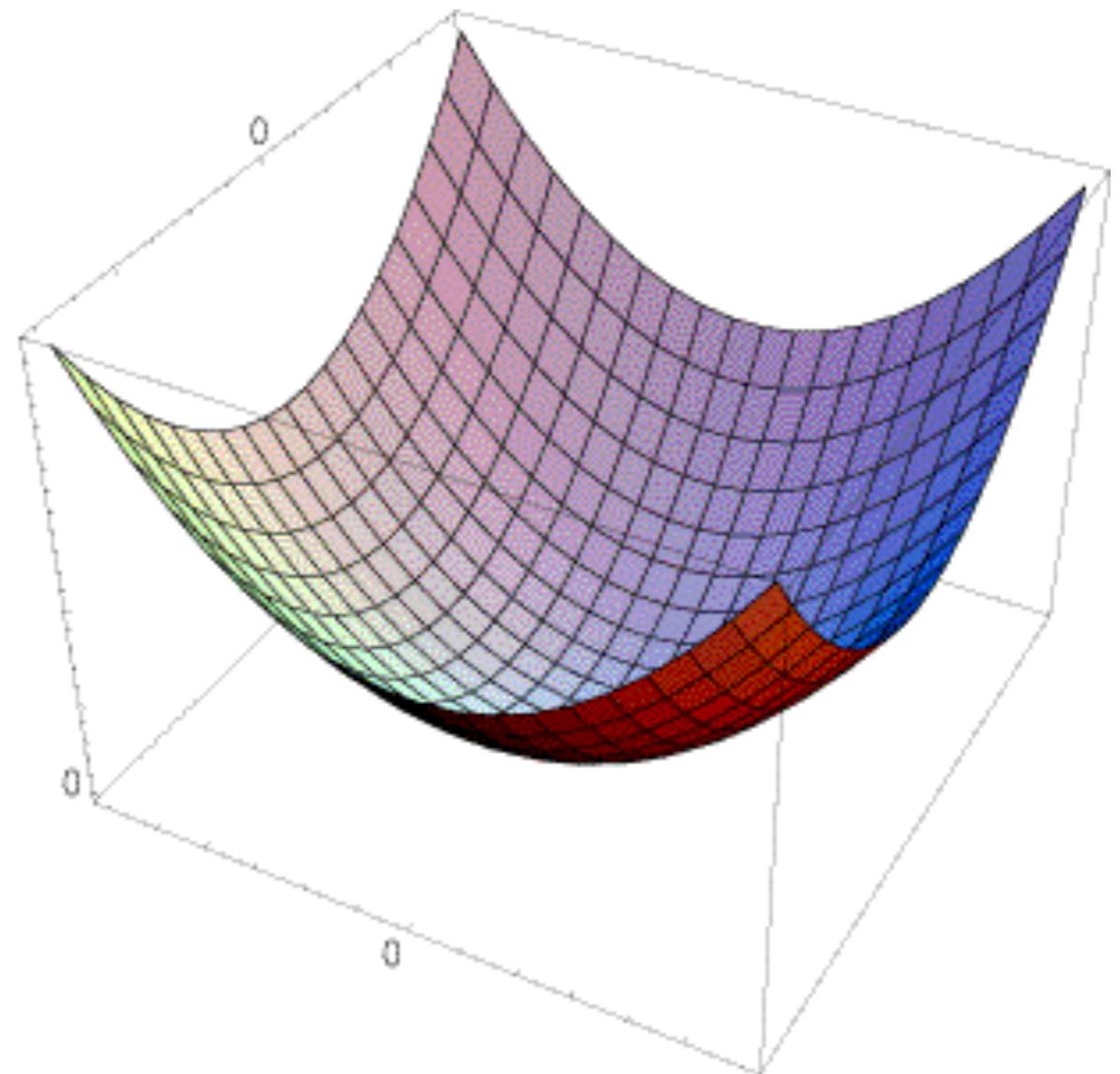
second derivative

Visualization of a quadratic

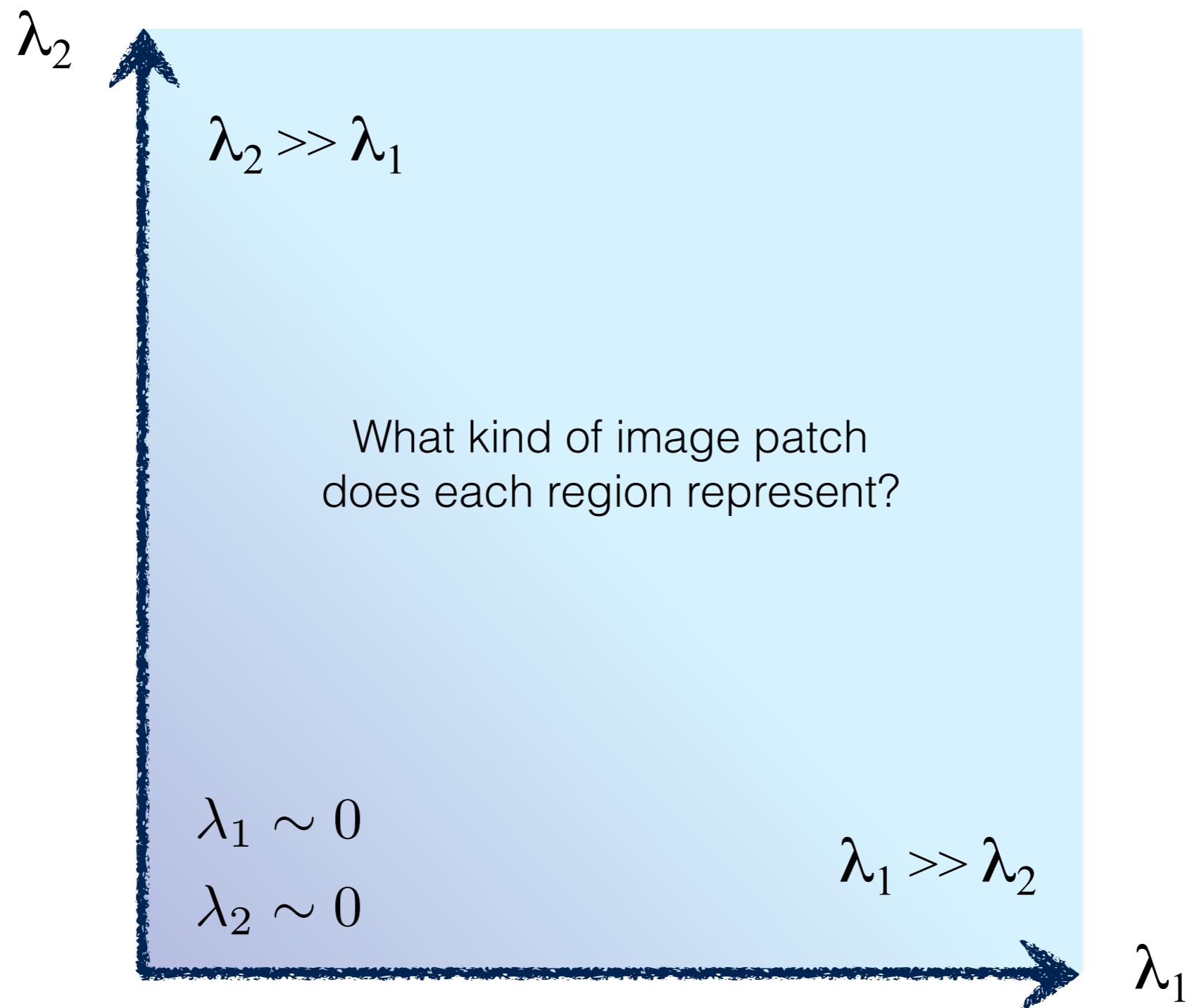
The surface $E(u,v)$ is locally approximated by a quadratic form

$$E(u,v) \approx [u \ v] M \begin{bmatrix} u \\ v \end{bmatrix}$$

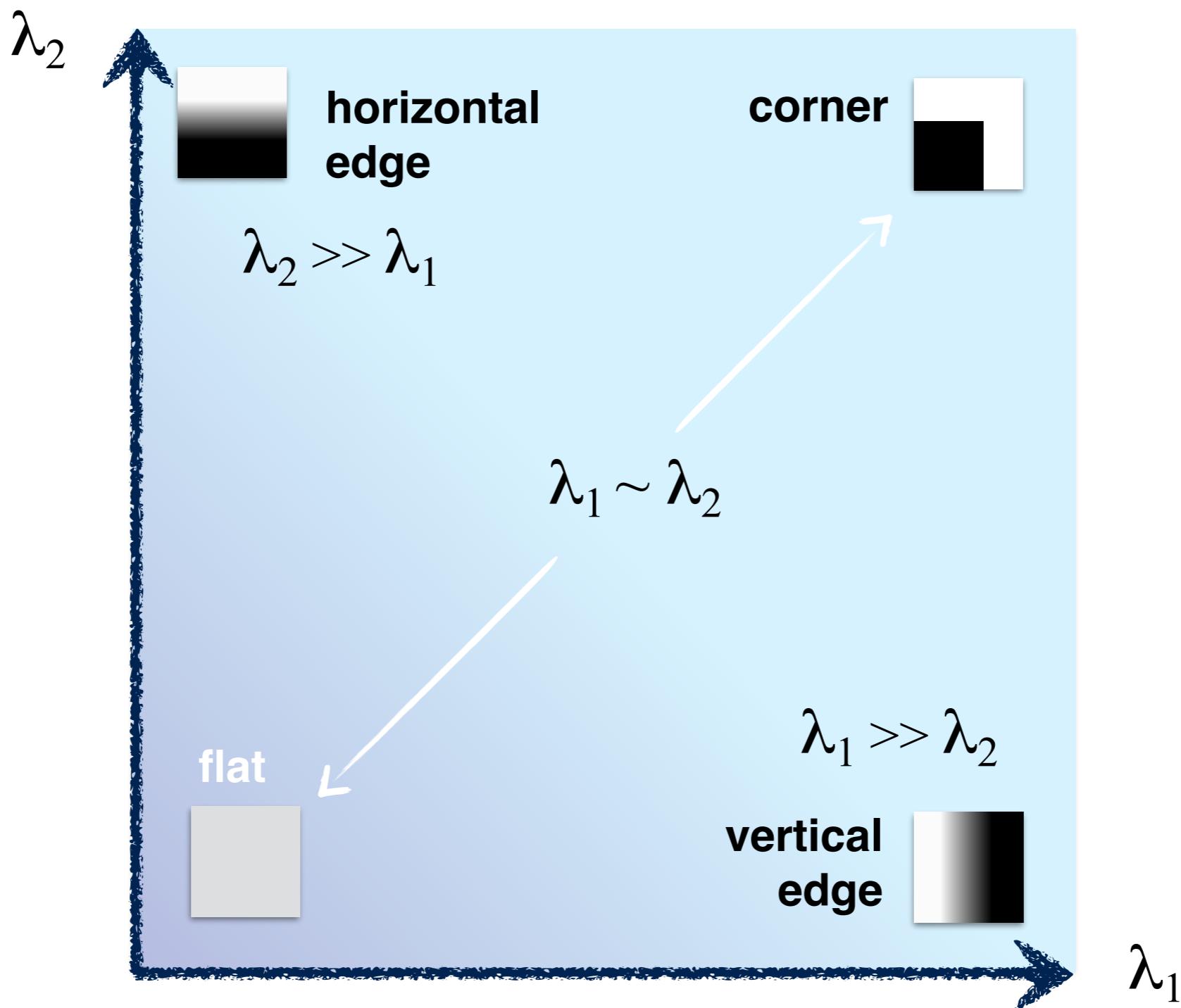
$$M = \Sigma \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$



interpreting eigenvalues



interpreting eigenvalues





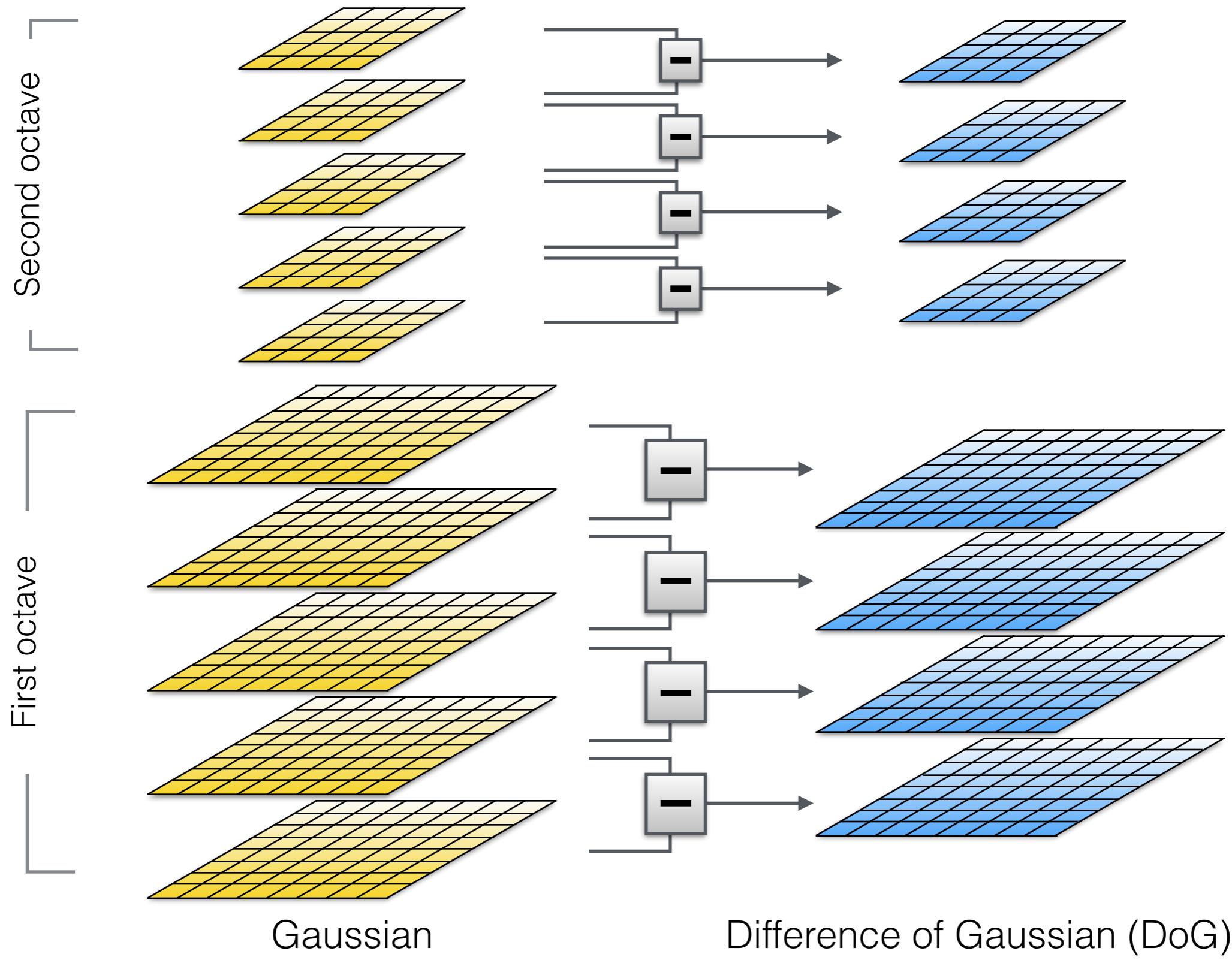
SIFT

(Scale Invariant Feature Transform)

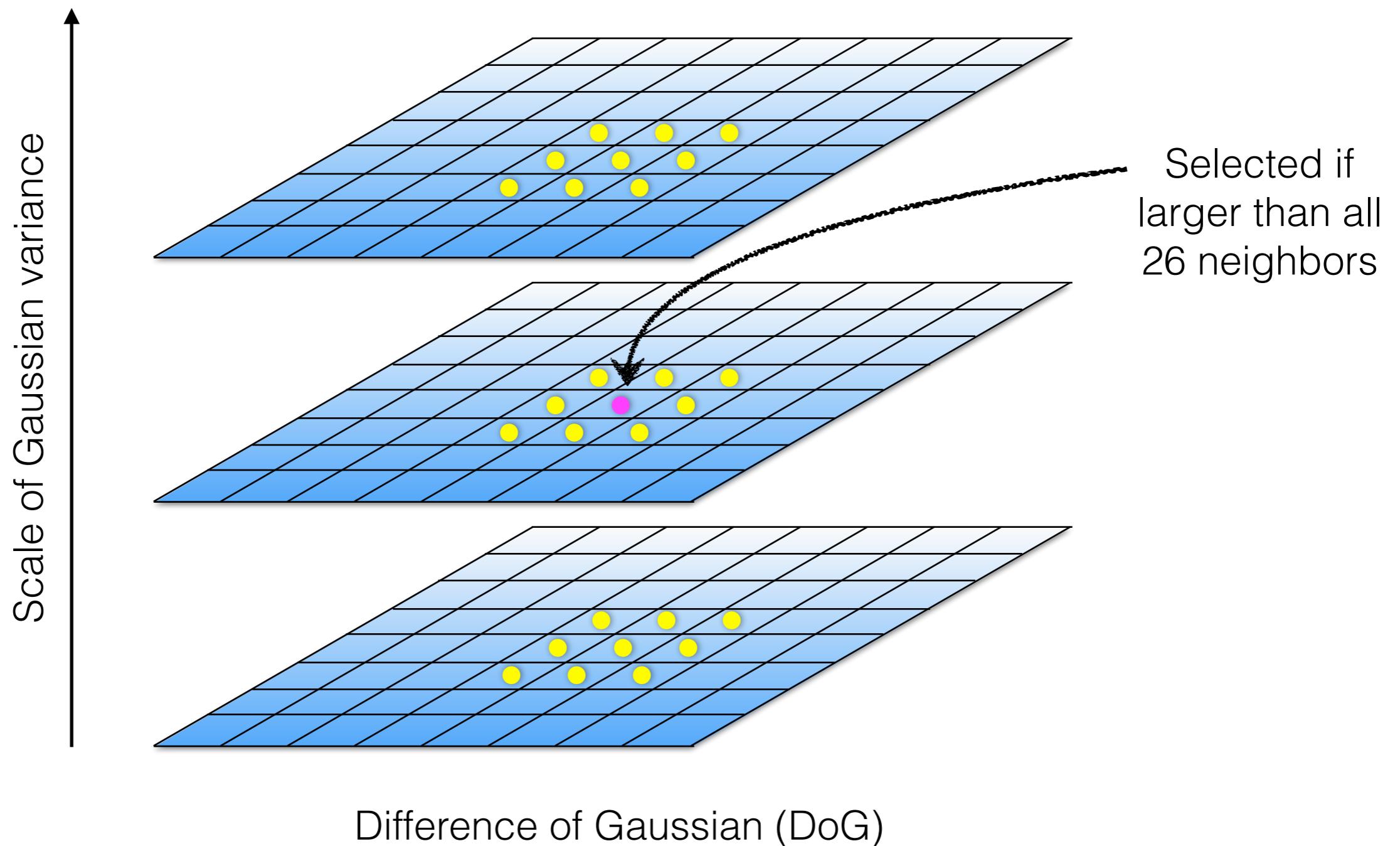
SIFT describes both a **detector** and **descriptor**

1. Multi-scale extrema detection
2. Keypoint localization
3. Orientation assignment
4. Keypoint descriptor

1. Multi-scale extrema detection



Scale-space extrema



2. Keypoint localization

2nd order Taylor series approximation of DoG scale-space

$$f(\mathbf{x}) = f + \frac{\partial f}{\partial \mathbf{x}}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \frac{\partial^2 f}{\partial \mathbf{x}^2} \mathbf{x}$$

$$\mathbf{x} = \{x, y, \sigma\}$$

Take the derivative and solve for extrema

$$\mathbf{x}_m = - \frac{\partial^2 f}{\partial \mathbf{x}^2}^{-1} \frac{\partial f}{\partial \mathbf{x}}$$

Additional tests to retain only strong features

3. Orientation assignment

$$m(x, y) = \sqrt{(L(x + 1, y) - L(x - 1, y))^2 + (L(x, y + 1) - L(x, y - 1))^2}$$

$$\theta(x, y) = \tan^{-1}((L(x, y + 1) - L(x, y - 1)) / (L(x + 1, y) - L(x - 1, y)))$$

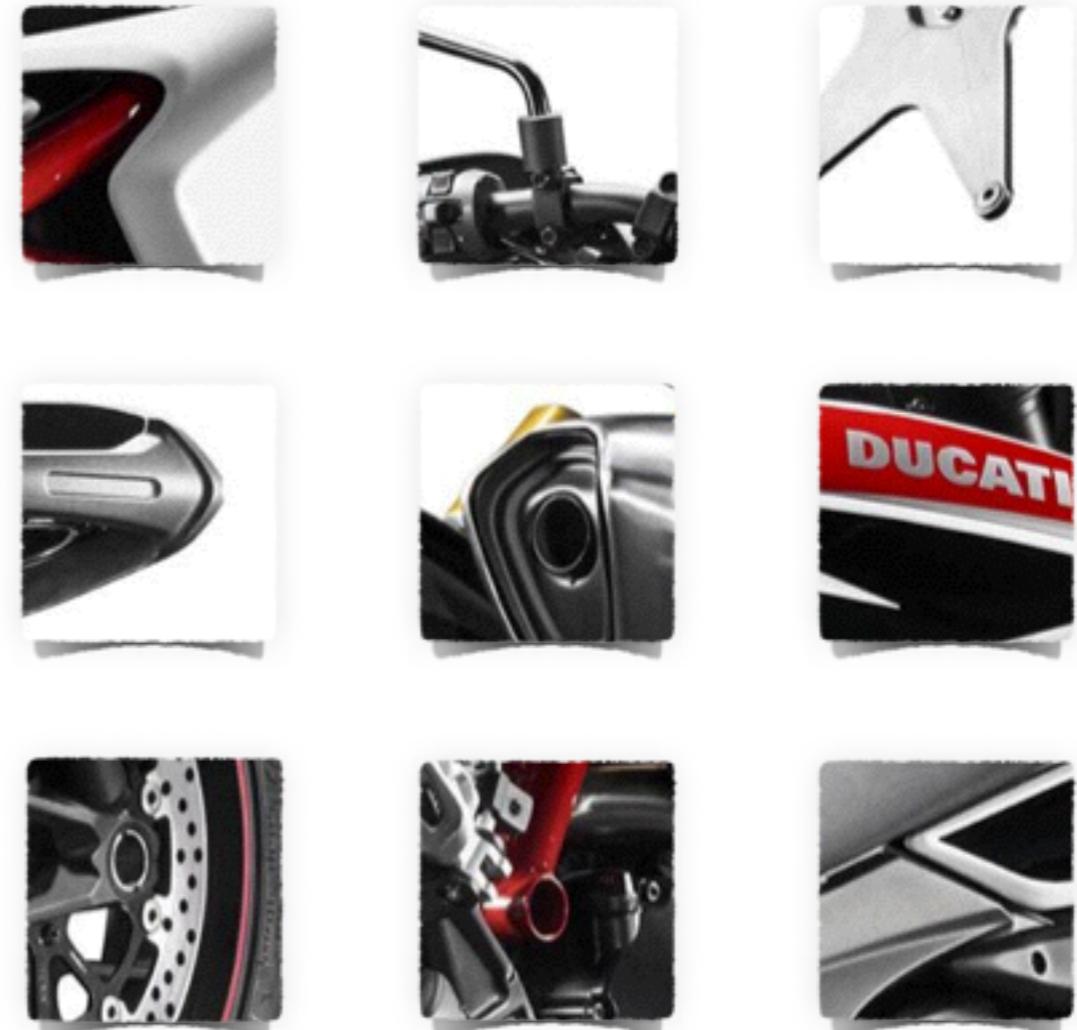
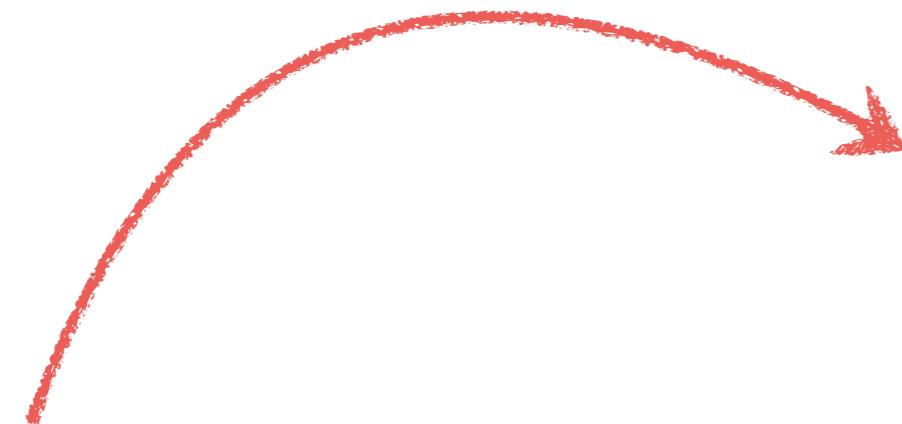
1. Compute gradients over image patch
2. Compute histogram over gradients
3. Take the histogram bin (orientation) with the greatest count
4. Rotate the image patch by that orientation
5. Compute descriptor

Bag-of-Words

Some local feature are very informative



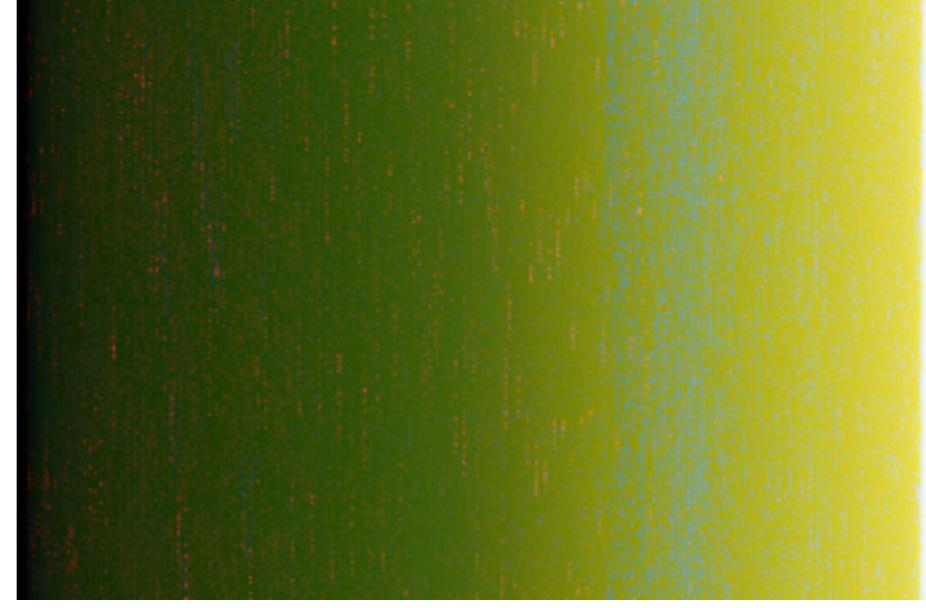
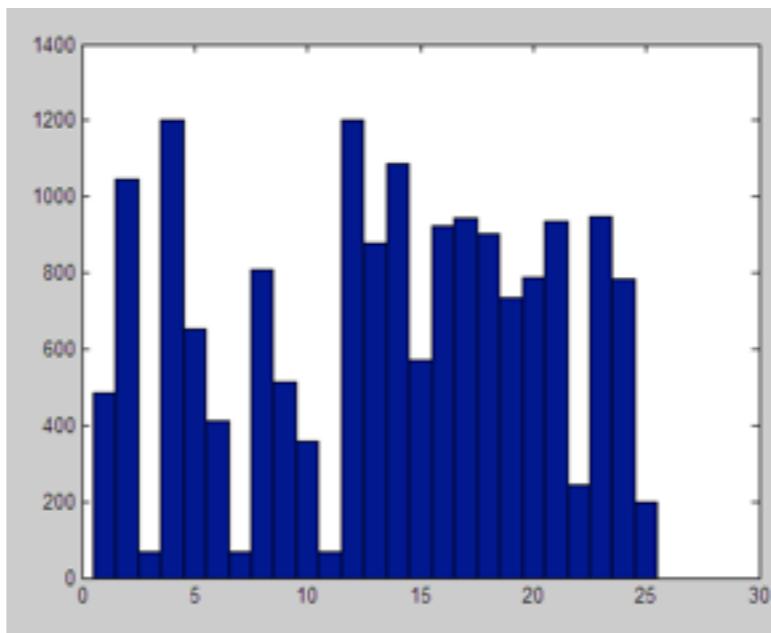
An object as



a collection of local features
(bag-of-features)

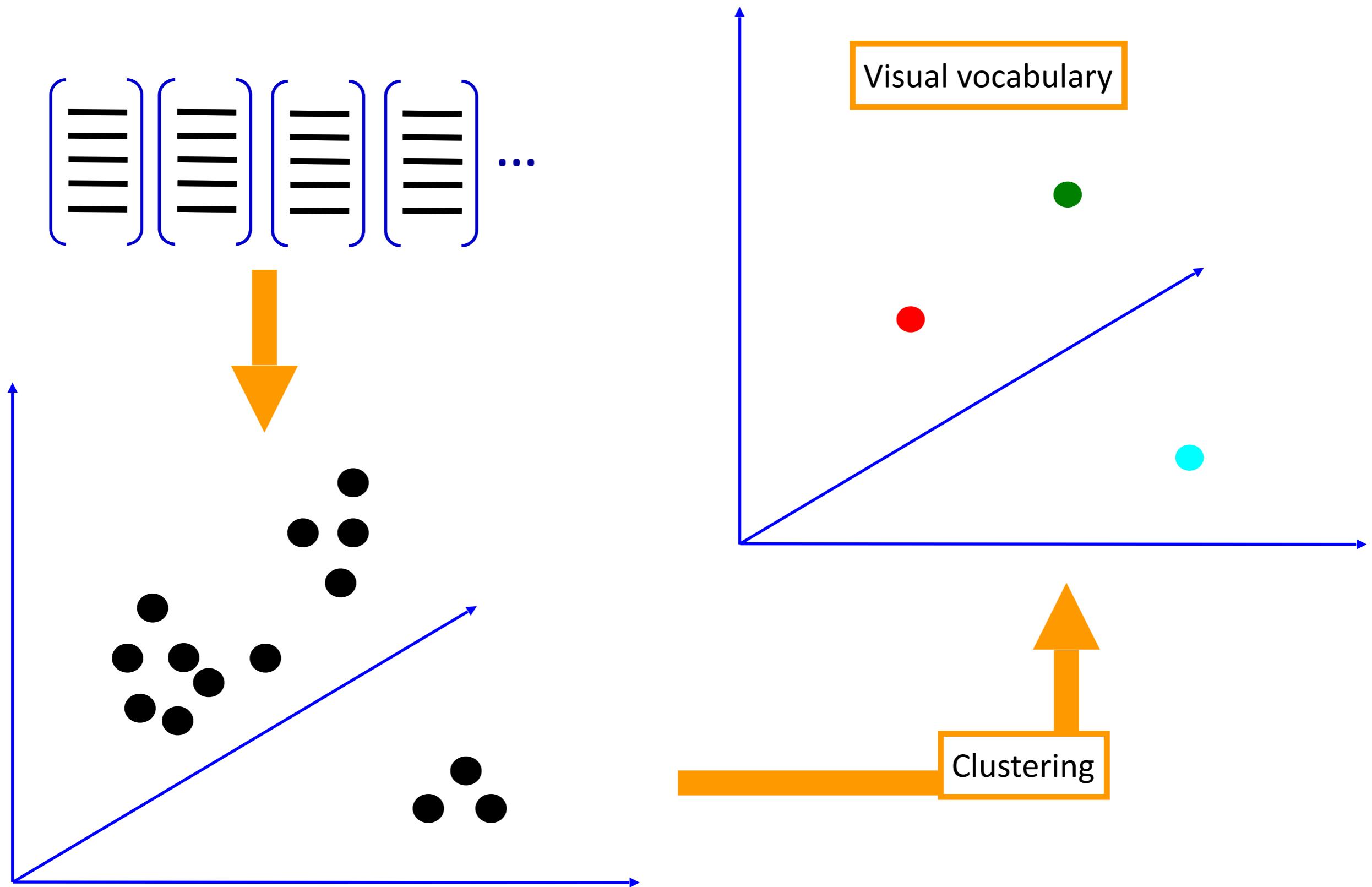
- deals well with occlusion
- scale invariant
- rotation invariant

But what about layout?



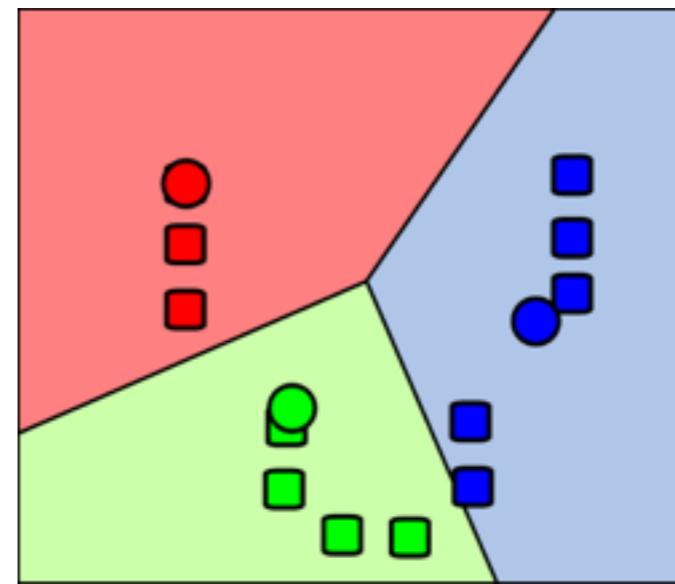
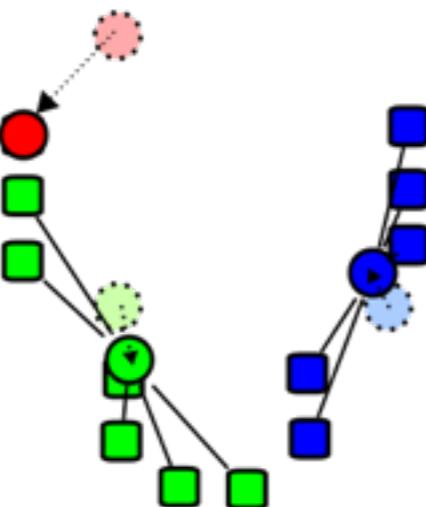
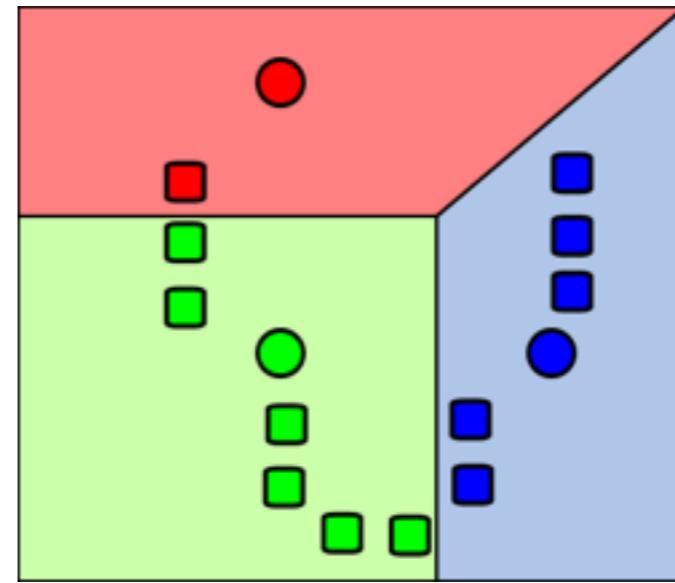
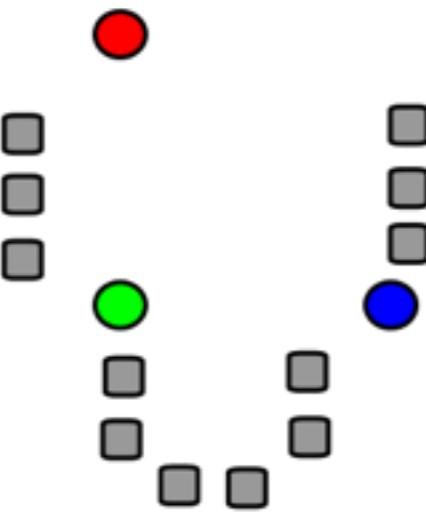
All of these images have the same color histogram
bag-of-words (histogram) representations remove spatial information

Do you know how to build a visual vocabulary?



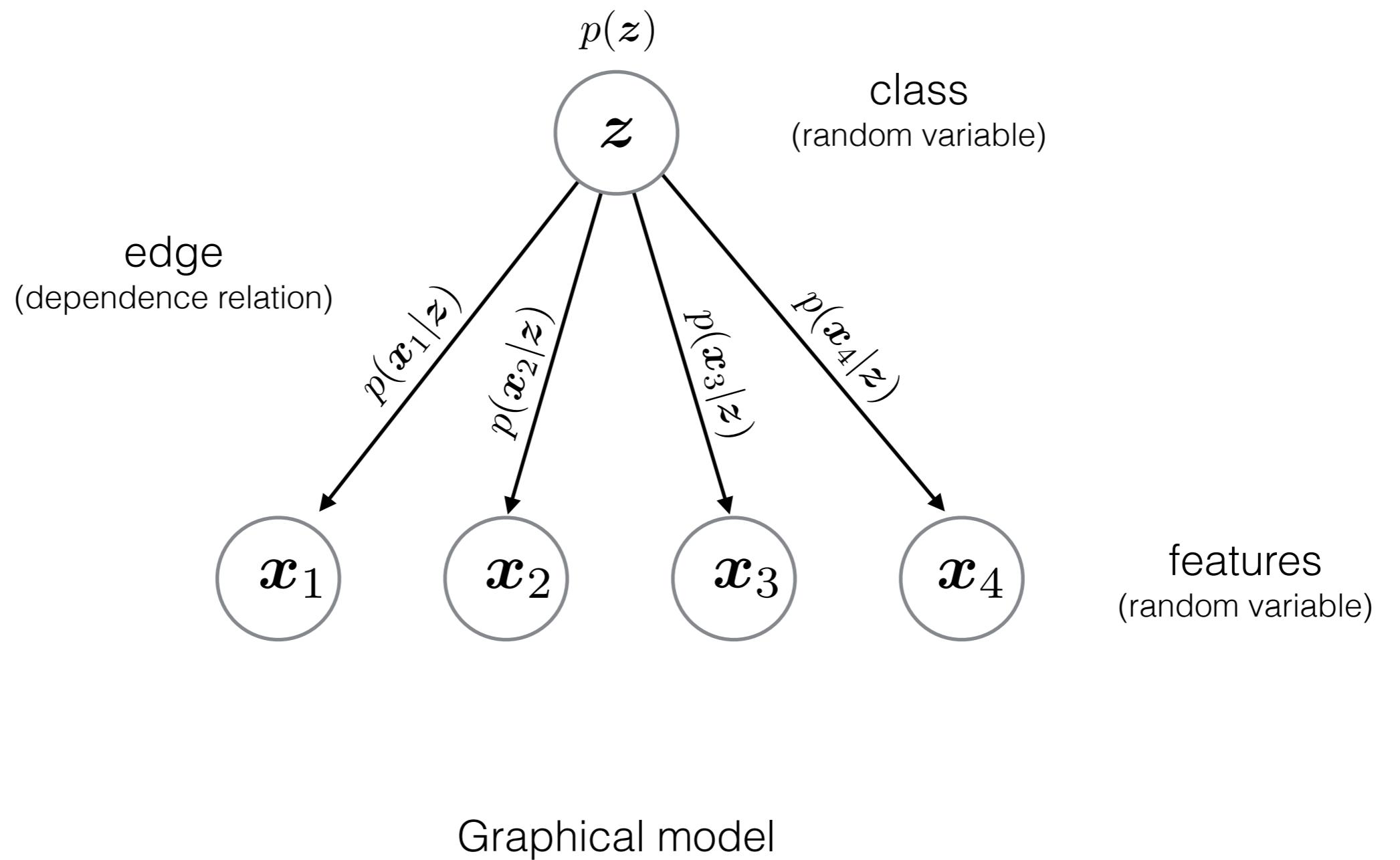
K-means Clustering

- Given k , the k-means algorithm consists of four steps:
- Select initial centroids at random.
- Assign each object to the cluster with the nearest centroid.
- Compute each centroid as the mean of the objects assigned to it.
- Repeat previous 2 steps until no change.



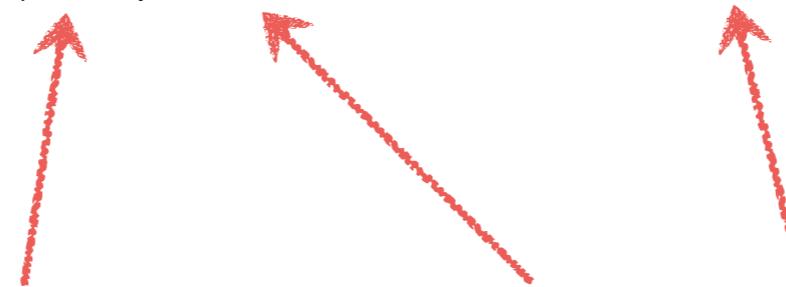
- How do you determine an optimal number of visual words?

Classification



This is called the posterior:
the probability of a class z given the observed features X

$$p(z|x_1, \dots, x_N)$$



For classification, z is a
discrete random variable
(e.g., car, person, building)

Each x is an observed feature
(e.g., visual words)

(it's a function that returns a single probability value)

The posterior can be decomposed according to
Bayes' Rule

$$p(A|B) = \frac{\text{likelihood} \quad \text{prior}}{p(B)}$$
$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

In our context...

$$p(z|x_1, \dots, x_N) = \frac{p(x_1, \dots, x_N | z)p(z)}{p(x_1, \dots, x_N)}$$

The naive Bayes' classifier is solving this optimization

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

MAP (maximum a posteriori) estimate

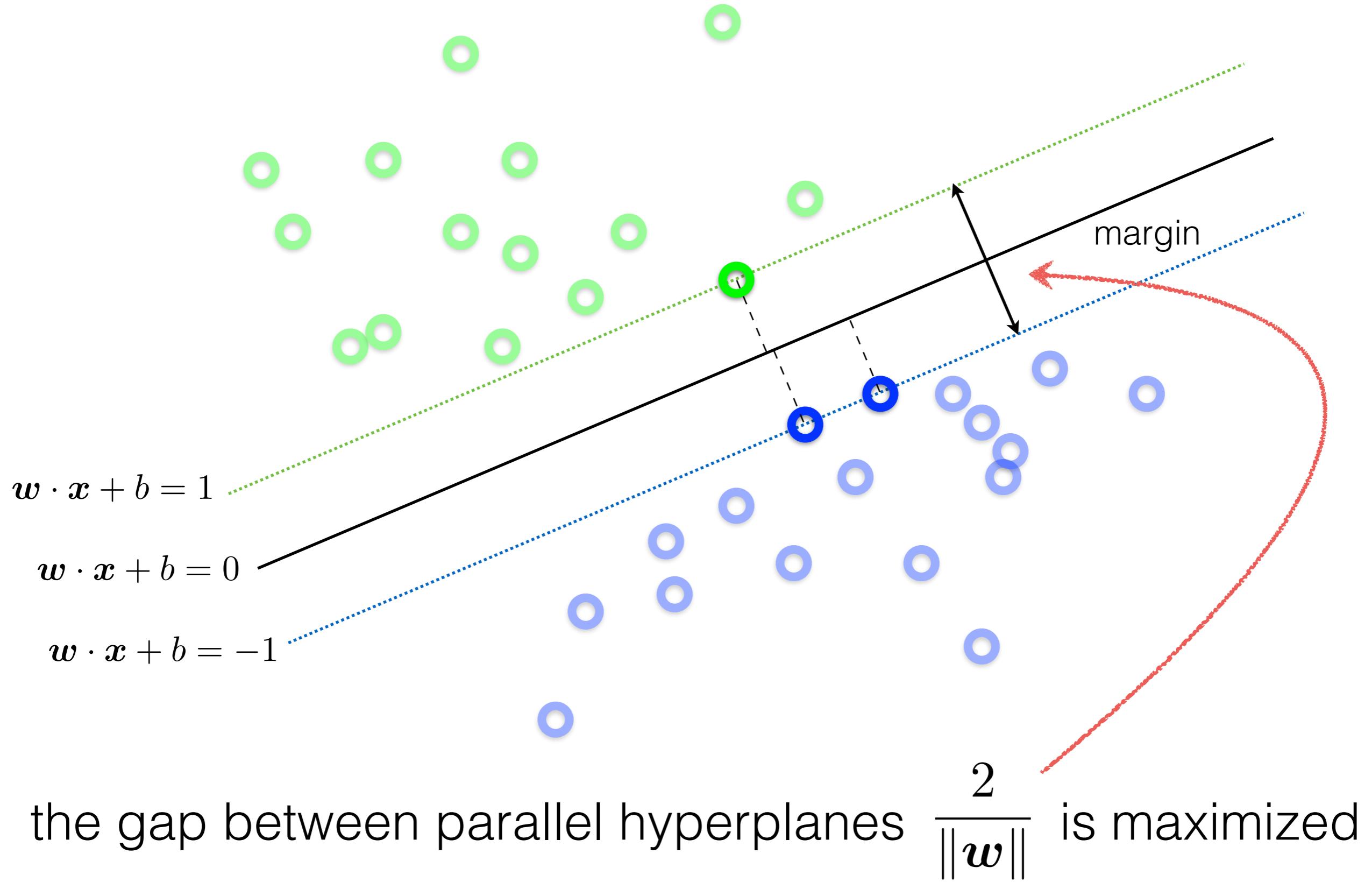
$$\hat{z} = \arg \max_{z \in \mathcal{Z}} \frac{p(\mathbf{X}|z)p(z)}{p(\mathbf{X})} \quad \text{Bayes' Rule}$$

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(\mathbf{X}|z)p(z) \quad \text{Remove constants}$$

A naive Bayes' classifier assumes all features are
conditionally independent

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N | z) &= p(\mathbf{x}_1 | z)p(\mathbf{x}_2, \dots, \mathbf{x}_N | z) \\ &= p(\mathbf{x}_1 | z)p(\mathbf{x}_2 | z)p(\mathbf{x}_3, \dots, \mathbf{x}_N | z) \\ &= p(\mathbf{x}_1 | z)p(\mathbf{x}_2 | z) \cdots p(\mathbf{x}_N | z) \end{aligned}$$

Find hyperplane \mathbf{w} such that ...



Can be formulated as a maximization problem

$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

$$\text{subject to } \mathbf{w} \cdot \mathbf{x}_i + b \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \text{ for } i = 1, \dots, N$$

Equivalently,

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \leq 0 \text{ for } i = 1, \dots, N$$

‘soft’ margin

objective

$$\min_{\mathbf{w}, \boldsymbol{\xi}} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

subject to

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{for } i = 1, \dots, N$$

- Every constraint can be satisfied if slack is large
- C is a regularization parameter
 - Small C: ignore constraints (larger margin)
 - Big C: constraints (small margin)
- Still QP problem (unique solution)