

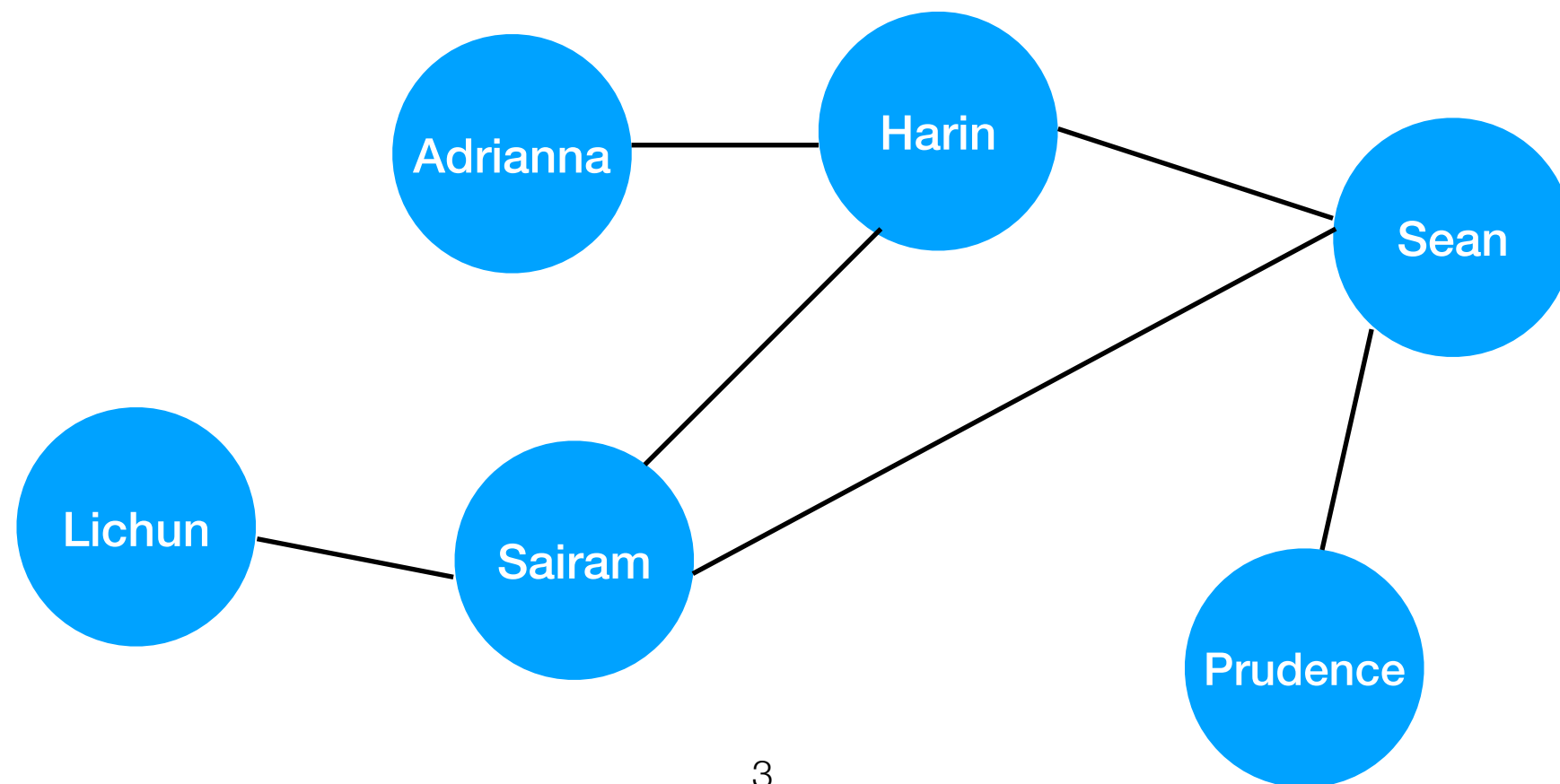
# CS/DS 541: Class 21

Jacob Whitehill

# Machine learning on graphs

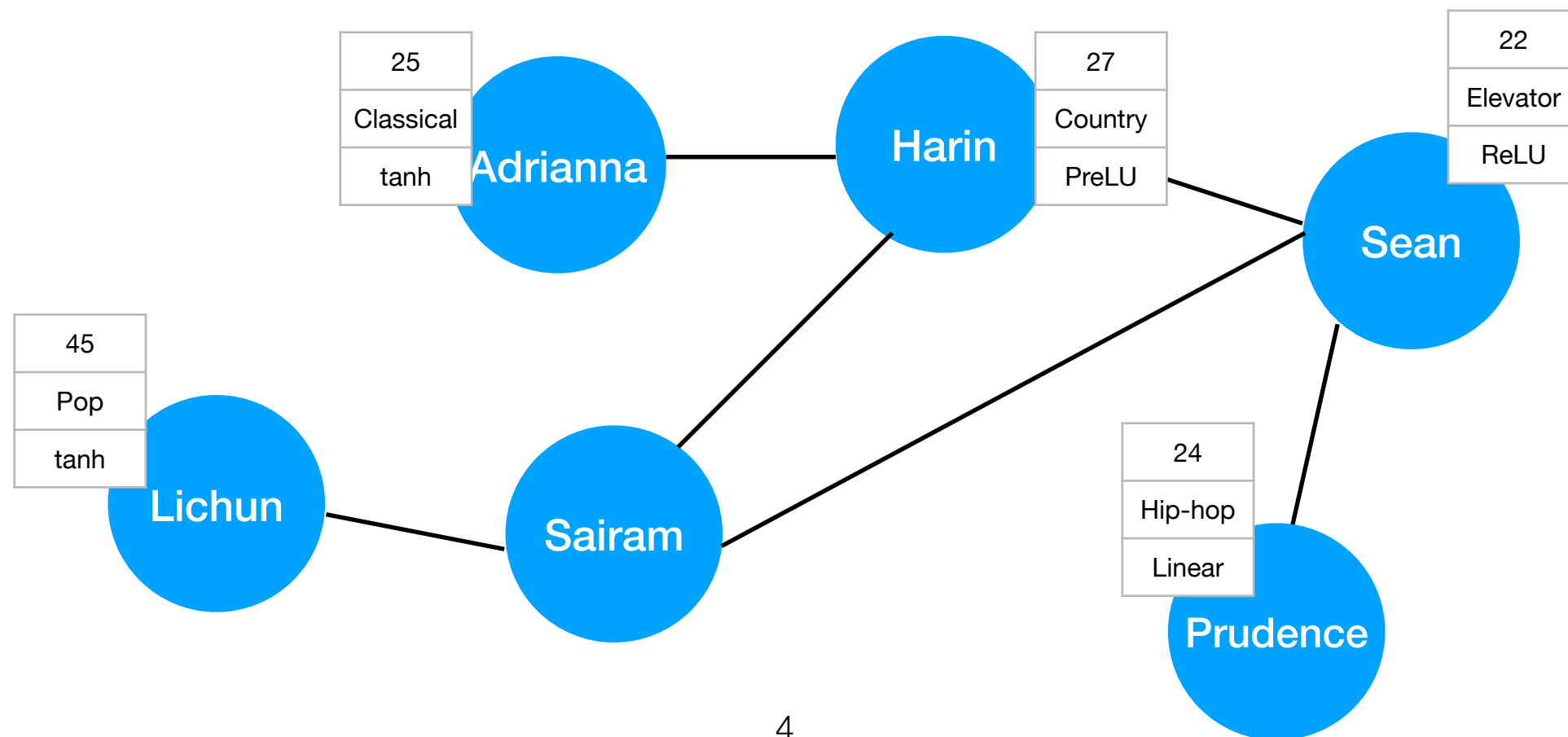
# Machine learning on graphs

- Consider an undirected graph where each node is associated with a feature vector.
- Example: social network of CS/DS 541.



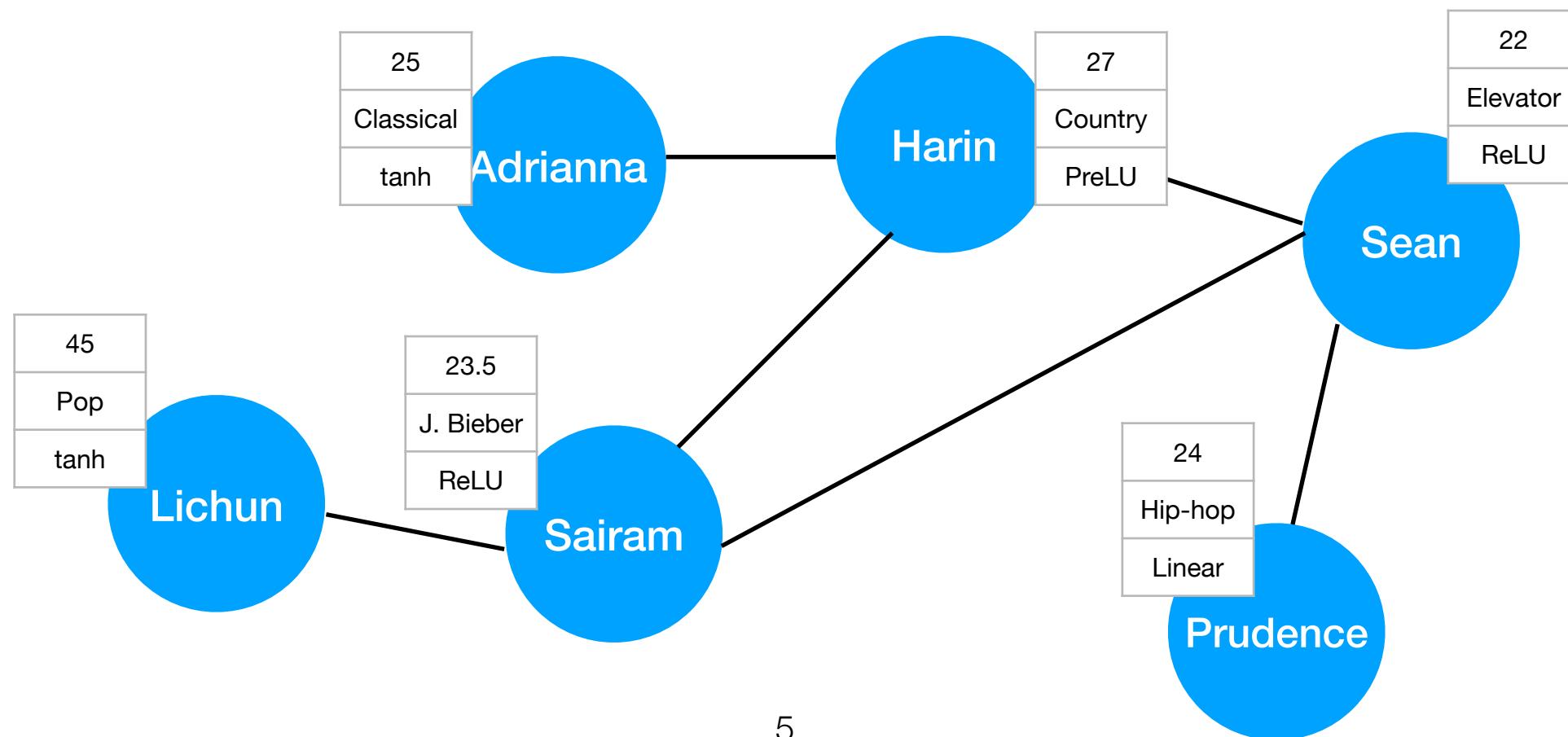
# Machine learning on graphs

- Consider an undirected graph where each node is associated with a feature vector.
- Example: social network of CS/DS 541.



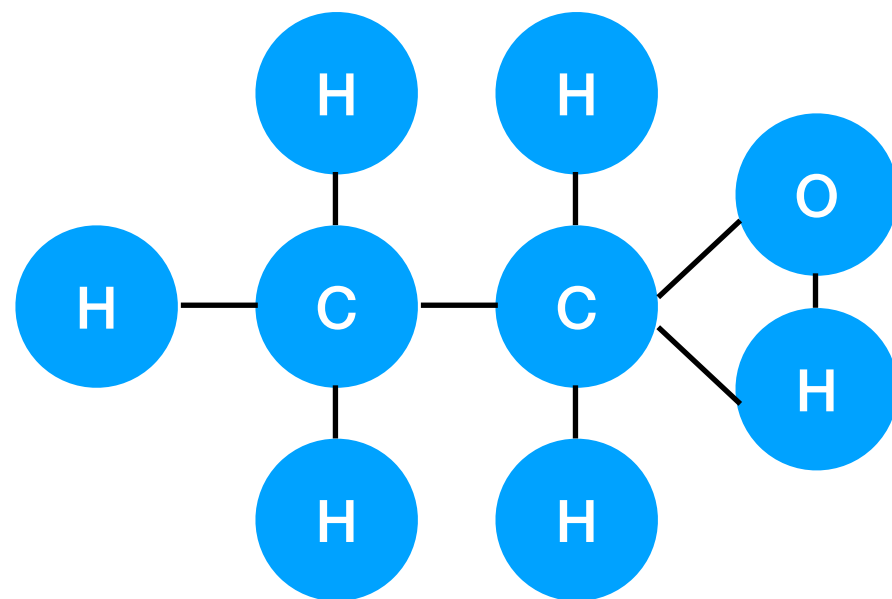
# Machine learning on graphs

- Based on the graph topology and set of observed feature vectors, can we infer the unobserved feature vector?
- This is a semi-supervised ML problem.



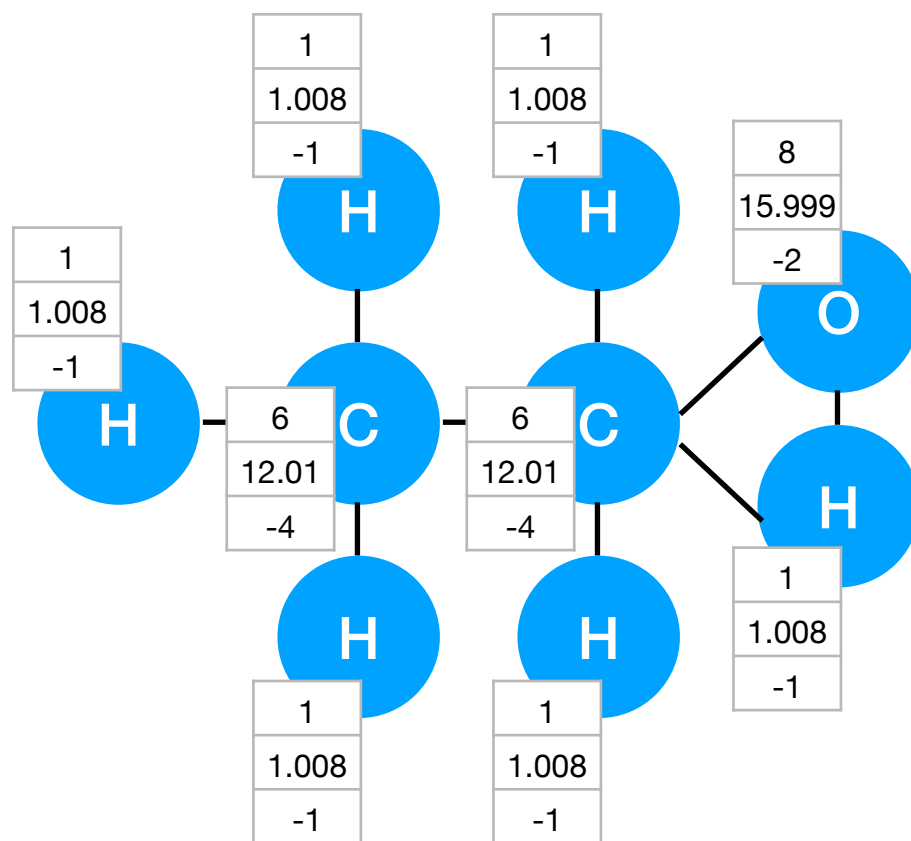
# Machine learning on graphs

- Example: computational chemistry.



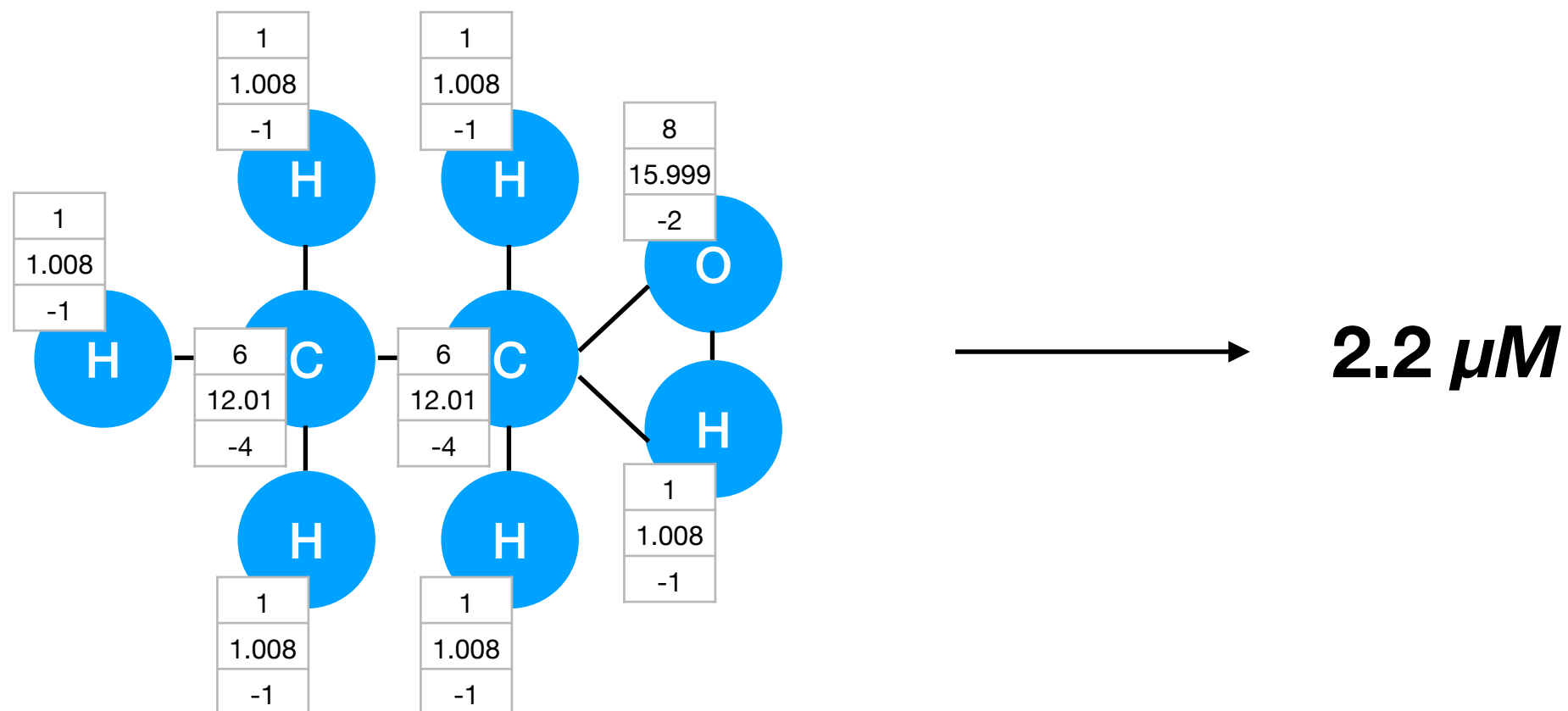
# Machine learning on graphs

- Example: computational chemistry.



# Machine learning on graphs

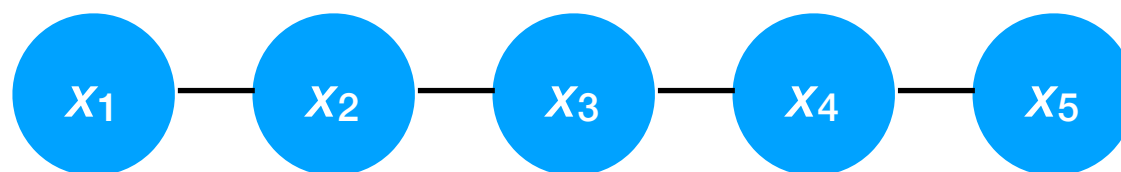
- Can we predict the binding affinity of a molecule to a particular receptor given its molecular structure?





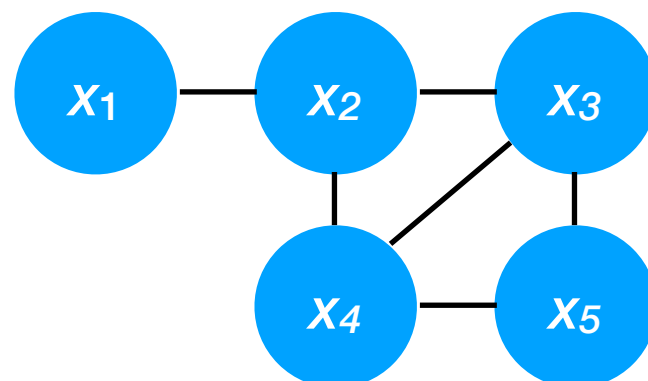
# Machine learning on graphs

- In contrast to vectors and images, the topology of graphs is **variable**.
- E.g., while each element of an  $n$ -vector is connected to 1-2 neighbors



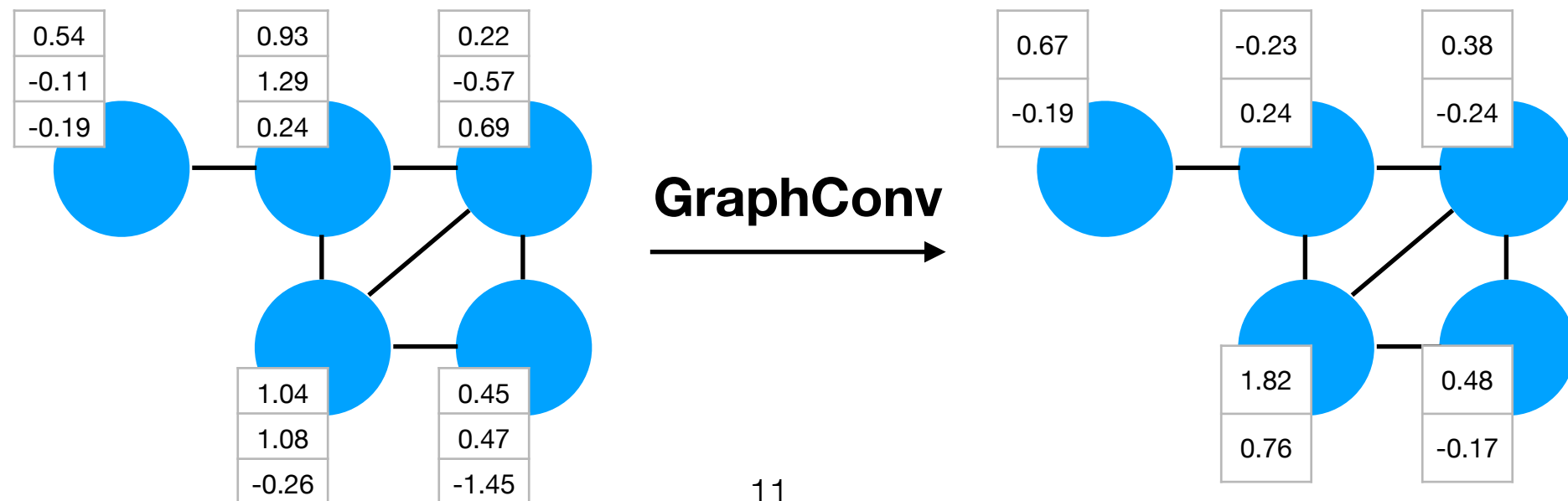
# Machine learning on graphs

- In contrast to vectors and images, the topology of graphs is **variable**.
- E.g., while each element of an  $n$ -vector is connected to 1-2 neighbors, each graph node has up to  $n-1$  neighbors.



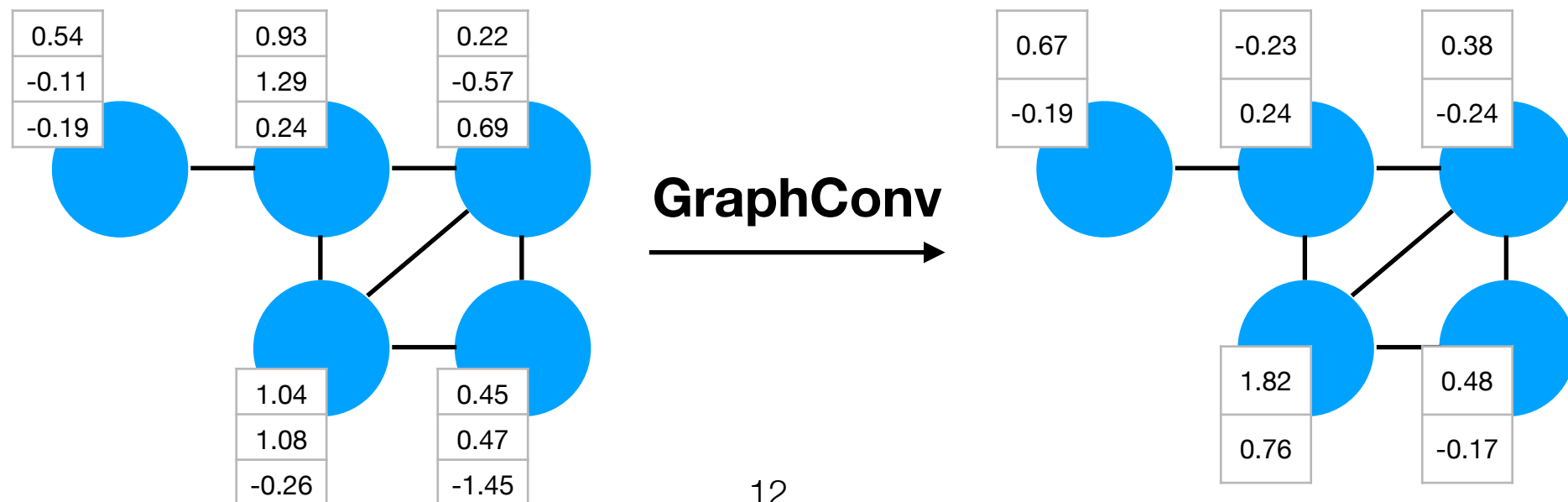
# Machine learning on graphs

- With GCNs, the feature vector at each node is transformed (by each GraphConv layer) based on the feature vectors of nearby nodes.



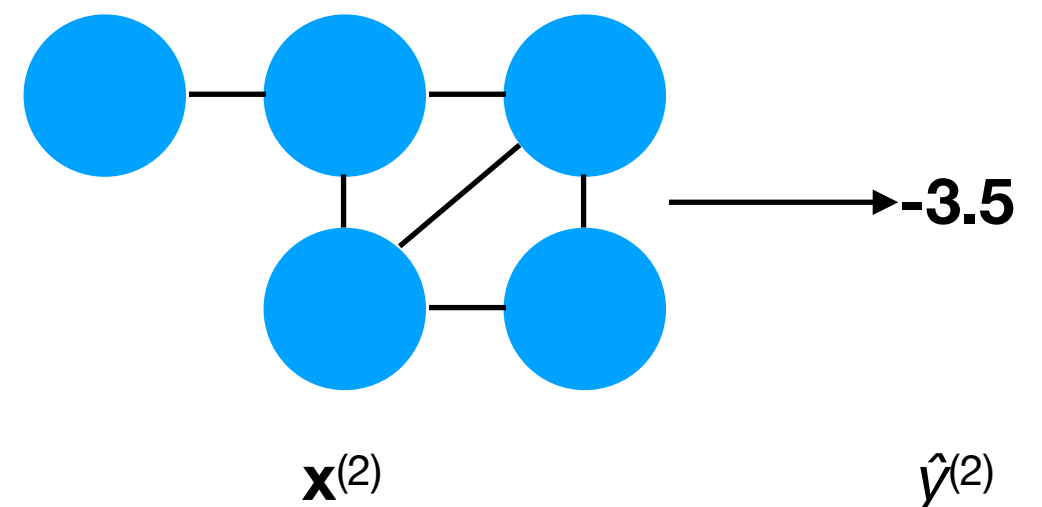
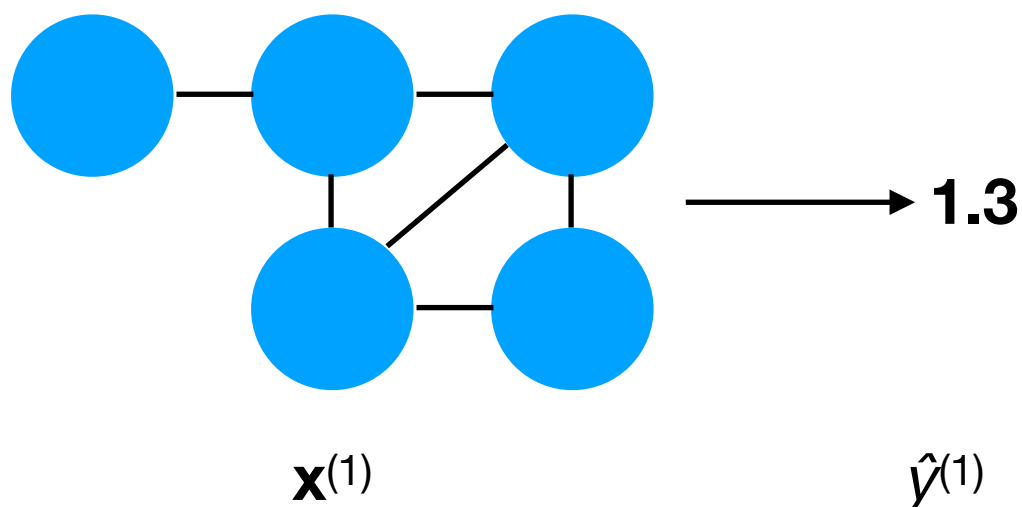
# Machine learning on graphs

- With GCNs, the feature vector at each node is transformed (by each GraphConv layer) based on the feature vectors of nearby nodes.
- With *true* graph convolution, “nearby” actually means *all* nodes.
- However, most GCNs actually perform *approximate* convolution where the set of influencing nodes is local.



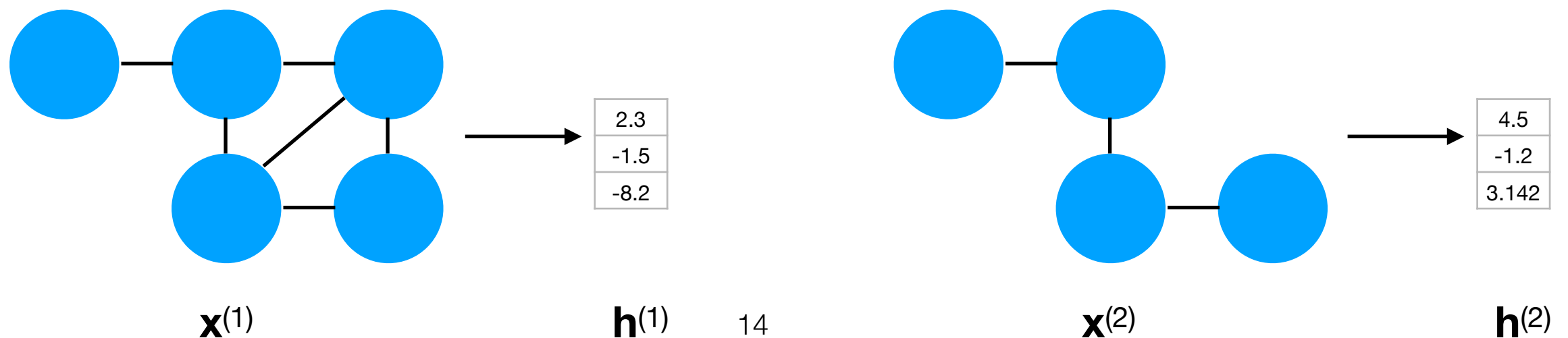
# Machine learning on graphs

- When doing ML on graphs, we usually either:
  1. Restrict the topology of all input graphs to be the same (but with different feature vectors):
    - Supervised: apply a tailored aggregation function that depends on a specific graph topology.
    - Semi-supervised: classify the unlabeled nodes.



# Machine learning on graphs

- ...or:
  2. Apply pooling to convert the graph into a fixed-length feature vector for downstream processing.



# GCN: theory

# Convolution

- What are two nice properties of convolution?



# Convolution

- What are two nice properties of convolution?
  1. Apply the same weights to everywhere in the input regardless of location:
    - Weight-sharing
    - Locality
  2. Convolution theorem.

# Convolution theorem

- According to the **convolution theorem**:

$$\mathbf{x} * \mathbf{w} = \mathcal{F}^{-1}(\mathcal{F}(\mathbf{x}) \odot \mathcal{F}(\mathbf{w}))$$

where  $*$  represents convolution, and  $\mathcal{F}$  is the Fourier transform.

# Convolution theorem

- According to the **convolution theorem**:

$$\mathbf{x} * \mathbf{w} = \mathcal{F}^{-1}(\mathcal{F}(\mathbf{x}) \odot \mathcal{F}(\mathbf{w}))$$

where  $*$  represents convolution, and  $\mathcal{F}$  is the Fourier transform.

- In other words, convolution can be implemented by:
  1. Applying a Fourier transform to  $\mathbf{x}$  and  $\mathbf{w}$ ;
  2. Multiplying the transformed vectors element-wise; and
  3. Applying an inverse Fourier transform to the result.

# Convolution theorem

- Demo.

# Fourier transform of functions

- The Fourier transform of a vector  $\mathbf{x}$  can be computed by multiplying by a **Fourier basis matrix  $\mathbf{V}$**  (though the FFT is faster):

$$\mathcal{F}(\mathbf{x}) = \mathbf{V}^\top \mathbf{x}$$

# Fourier transform of functions

- Where does the basis  $\mathbf{V}$  come from?

$$\mathbf{v}_k = \frac{1}{\sqrt{n}} \begin{bmatrix} e^{j2\pi 0k/n} \\ \vdots \\ e^{j2\pi (n-1)k/n} \end{bmatrix}$$

- The Fourier bases are eigenfunctions of the Laplacian operator for real functions.

# Fourier transform of graphs

- Let graph function  $x(i)$  be a function from graph nodes to scalar values; equivalently, we can let  $\mathbf{x} \in \mathbb{R}^m$ .
- For graphs, the Fourier basis matrix  $\mathbf{V}$  consists of the eigenvectors of the **normalized graph Laplacian matrix**  $\mathbf{L}$ :

$$\mathbf{L} = \mathbf{I} - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$$

**A:** adjacency matrix

**D:** degree matrix

$$\tilde{\mathbf{x}} = \mathcal{F}(\mathbf{x}) = \mathbf{V}^\top \mathbf{x}$$

$$\mathbf{x} = \mathcal{F}^{-1}(\tilde{\mathbf{x}}) = \mathbf{V} \tilde{\mathbf{x}}$$

# Graph convolution

- With this machinery, we can now implement convolution of a graph function  $\mathbf{x}$  with a filter  $\mathbf{w}$ :

$$\mathcal{F}^{-1}(\mathcal{F}(\mathbf{w}) \odot \mathcal{F}(\mathbf{x}))$$



# Graph convolution

- With this machinery, we can now implement convolution of a graph function  $\mathbf{x}$  with a filter  $\mathbf{w}$ :

$$\mathcal{F}^{-1}(\mathcal{F}(\mathbf{w}) \odot \mathcal{F}(\mathbf{x})) = \mathbf{V}(\mathbf{V}^\top \mathbf{w} \odot \mathbf{V}^\top \mathbf{x})$$

# Graph convolution

- With this machinery, we can now implement convolution of a graph function  $\mathbf{x}$  with a filter  $\mathbf{w}$ :

$$\begin{aligned}\mathcal{F}^{-1}(\mathcal{F}(\mathbf{w}) \odot \mathcal{F}(\mathbf{x})) &= \mathbf{V}(\mathbf{V}^\top \mathbf{w} \odot \mathbf{V}^\top \mathbf{x}) \\ &= \mathbf{V}\theta\mathbf{V}^\top \mathbf{x}\end{aligned}$$

where we define diagonal matrix  $\theta$  to be equivalent to multiplying  $\mathbf{V}^\top \mathbf{x}$  element-wise by  $\mathbf{V}^\top \mathbf{w}$ .

# Approximation

- Here,  $\theta$  is an  $n$ -dimensional vector of parameters; this performs a true graph convolution.
- However, this also implies that the convolution output, for each node in the graph, depends on the entire graph -- not just a local subgraph.
- It is also slow, since we must multiply by potentially large  $(n \times n)$  matrices.
- In Kipf & Welling (2017), they instead make an approximation...

# Approximation

- Here,  $\theta$  is an  $n$ -dimensional vector of parameters; this performs a true graph convolution.
- However, this also implies that the convolution output, for each node in the graph, depends on the entire graph -- not just a local subgraph.
- It is also slow, since we must multiply by potentially large  $(n \times n)$  matrices.
- In Kipf & Welling (2017), they instead make an approximation (and wrap it with  $\sigma$ ) that yields the result:

$$\mathbf{h} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} \right)$$

# Approximation

- **X**:  $n \times k$  feature matrix (for a graph with  $n$  nodes, and  $k$  feature channels).
- **W**:  $k \times m$  filter matrix
- **A**: graph adjacency matrix
- **D**: graph degree matrix

Linear transformation from  
 $k$  features to  $m$  features

$$\mathbf{h} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} \right)$$

# Approximation

- **X**:  $n \times k$  feature matrix (for a graph with  $n$  nodes, and  $k$  feature channels).
- **W**:  $k \times m$  filter matrix
- **A**: graph adjacency matrix
- **D**: graph degree matrix

Weighted sum over the  
features of neighboring nodes

$$\mathbf{h} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} \right)$$

# Approximation

- **$\mathbf{X}$** :  $n \times k$  feature matrix (for a graph with  $n$  nodes, and  $k$  feature channels).
- **$\mathbf{W}$** :  $k \times m$  filter matrix
- **$\mathbf{A}$** : graph adjacency matrix
- **$\mathbf{D}$** : graph degree matrix

$$\mathbf{h} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} \right)$$

# Relationship to FC layer

- Special cases:
  - If  $\mathbf{A}=\mathbf{I}$  (i.e., each node is adjacent only to itself), then the Graph Convolution (GC) layer simplifies to just a regular Fully Connected (FC) layer.

$$\mathbf{h} = \sigma \left( \mathbf{XW} \right)$$



# Relationship to FC layer

- Special cases:
  - If  $\mathbf{A}=\mathbf{I}$  (i.e., each node is adjacent only to itself), then the Graph Convolution (GC) layer simplifies to just a regular Fully Connected (FC) layer.
  - If  $\mathbf{A}=\mathbf{1}/n$  (i.e., each node is adjacent to every node in the graph), then the GC layer performs average pooling.

$$\mathbf{h} = \sigma \left( \mathbf{1}/n \mathbf{XW} \right)$$

# Connection between graph convolution and image convolution

# Convolution is linear

- Recall that convolution is linear and can thus be implemented by multiplication with a matrix **W**, e.g.:

1
2
-2
0
3

**Image**

1
2
3

**Filter**

-1
-2
7

**Feature map**

$$\begin{bmatrix} 1 & 2 & 3 & 0 & 0 \\ 0 & 1 & 2 & 3 & 0 \\ 0 & 0 & 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ -2 \\ 0 \\ 3 \end{bmatrix} = \begin{bmatrix} -1 \\ -2 \\ 7 \end{bmatrix} \quad \mathbf{Wx} = \mathbf{h}$$

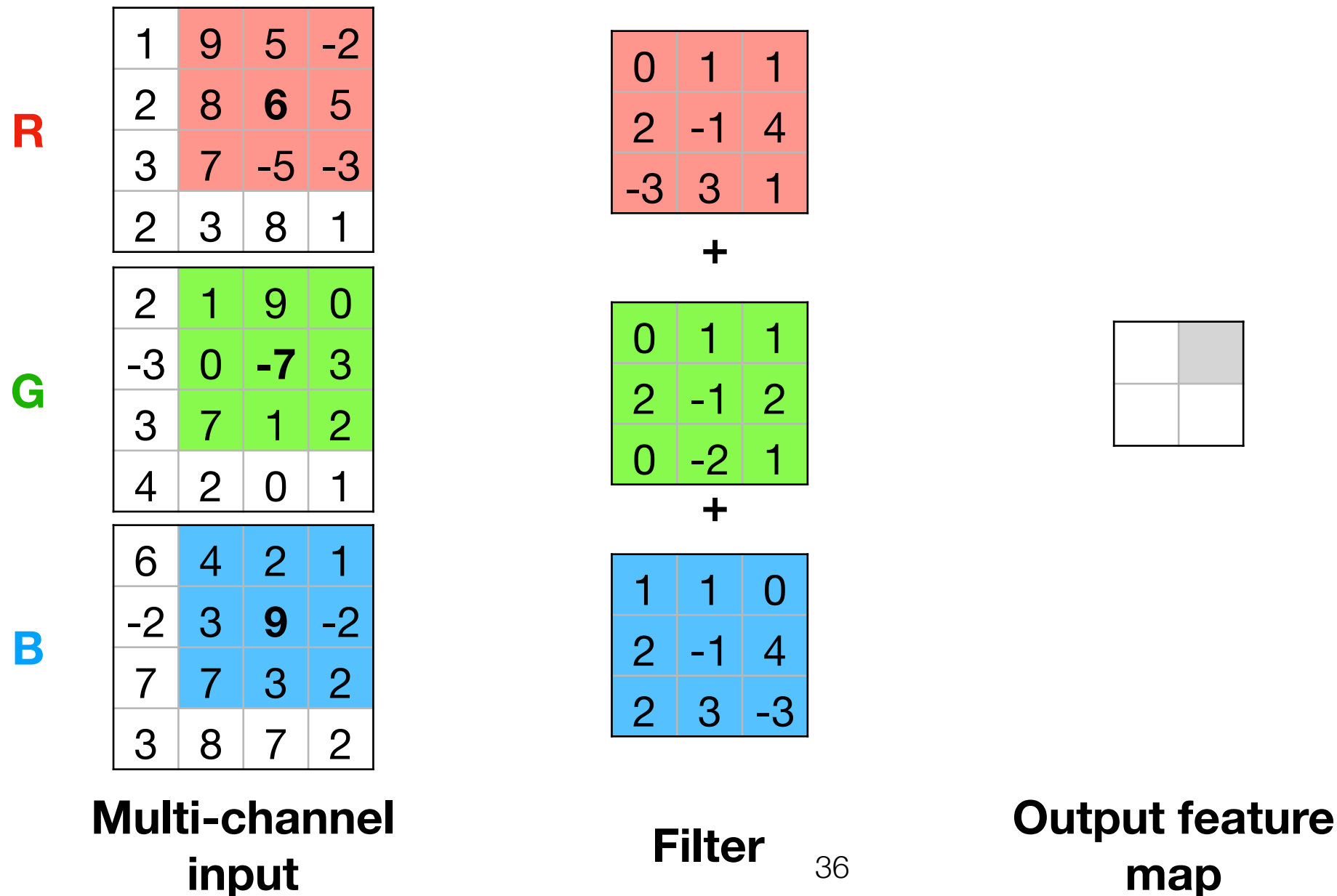
**W**

**x**

**h**

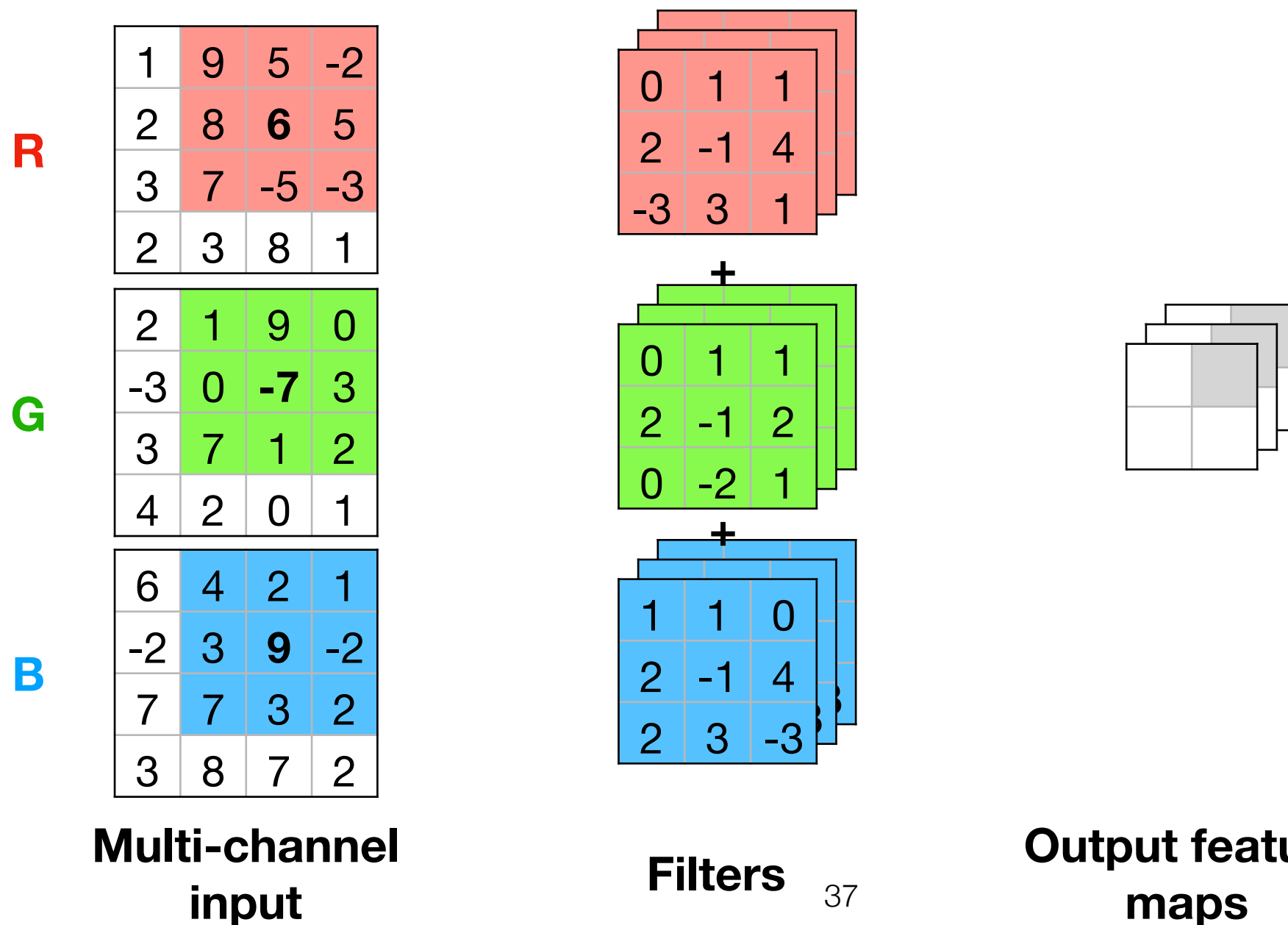
# Review of 3D convolution

- Recall how 3D convolution (3x3 filter) on a stack of input images works.



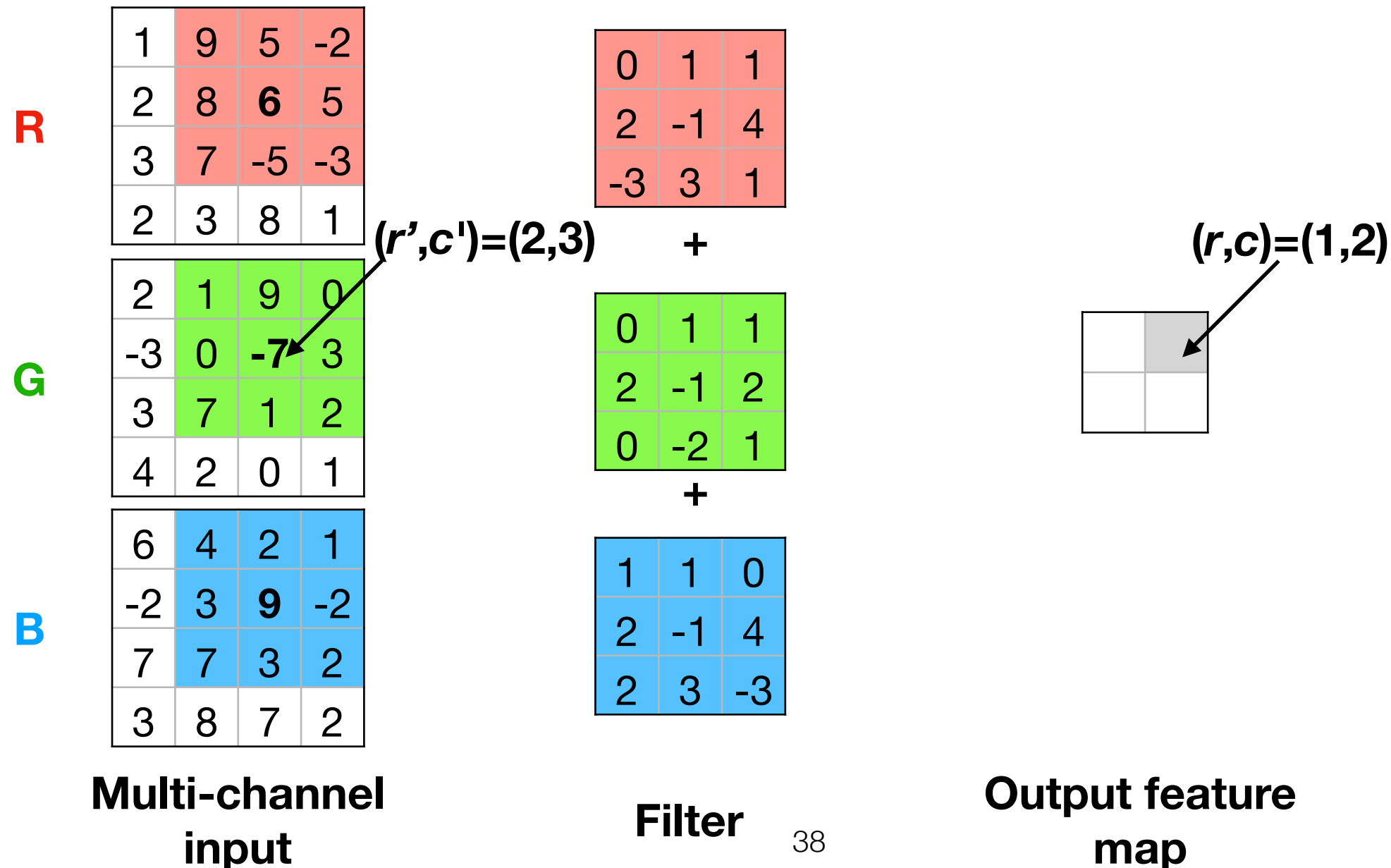
# Review of 3D convolution

- We can also apply multiple such filters to obtain a stack of output feature maps.



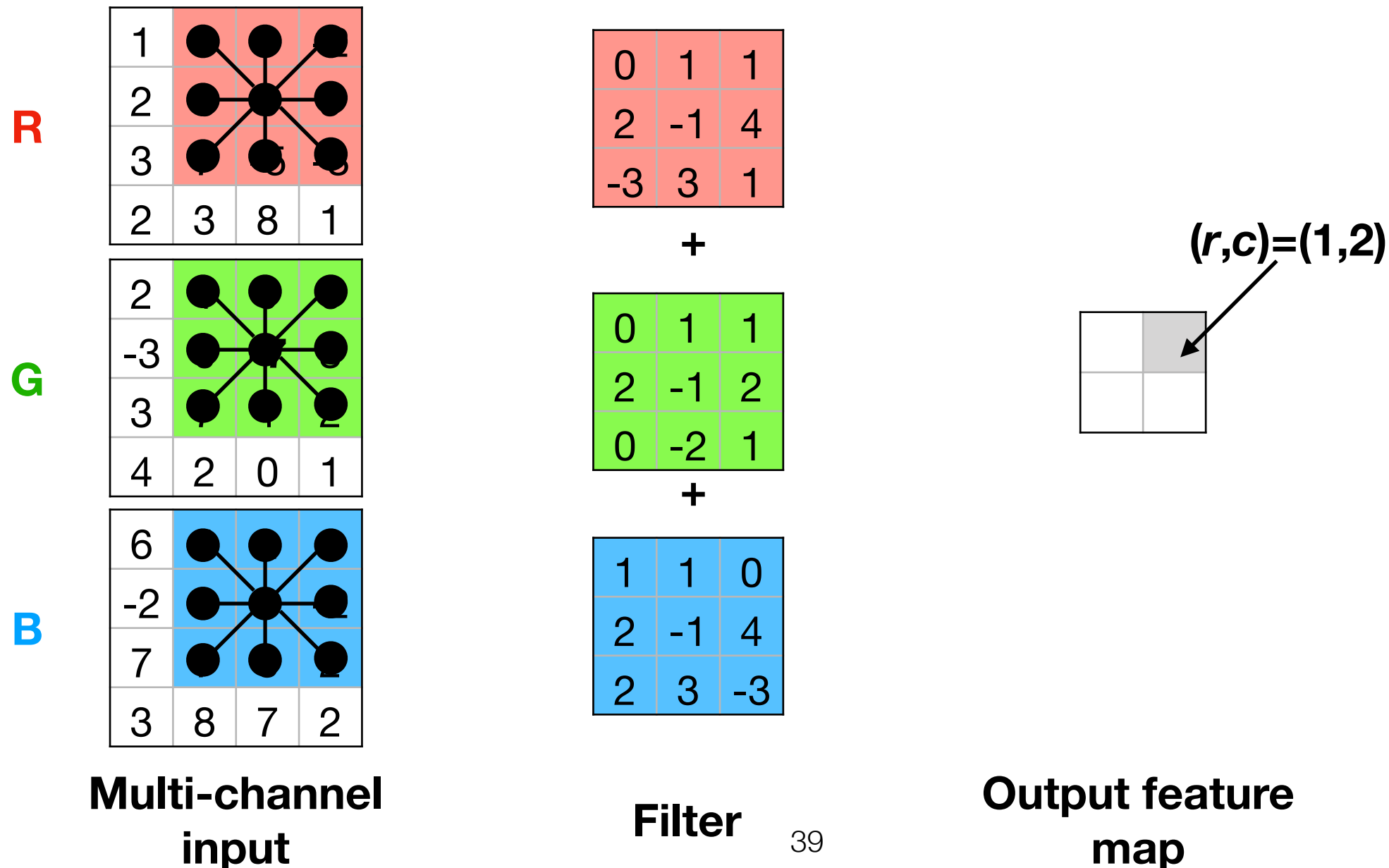
# Review of 3D convolution

- The value of each element  $(r,c)$  of the output feature map depends on the set of  $(k*k)$  **neighbors** of the corresponding element  $(r',c')$  in the input image.



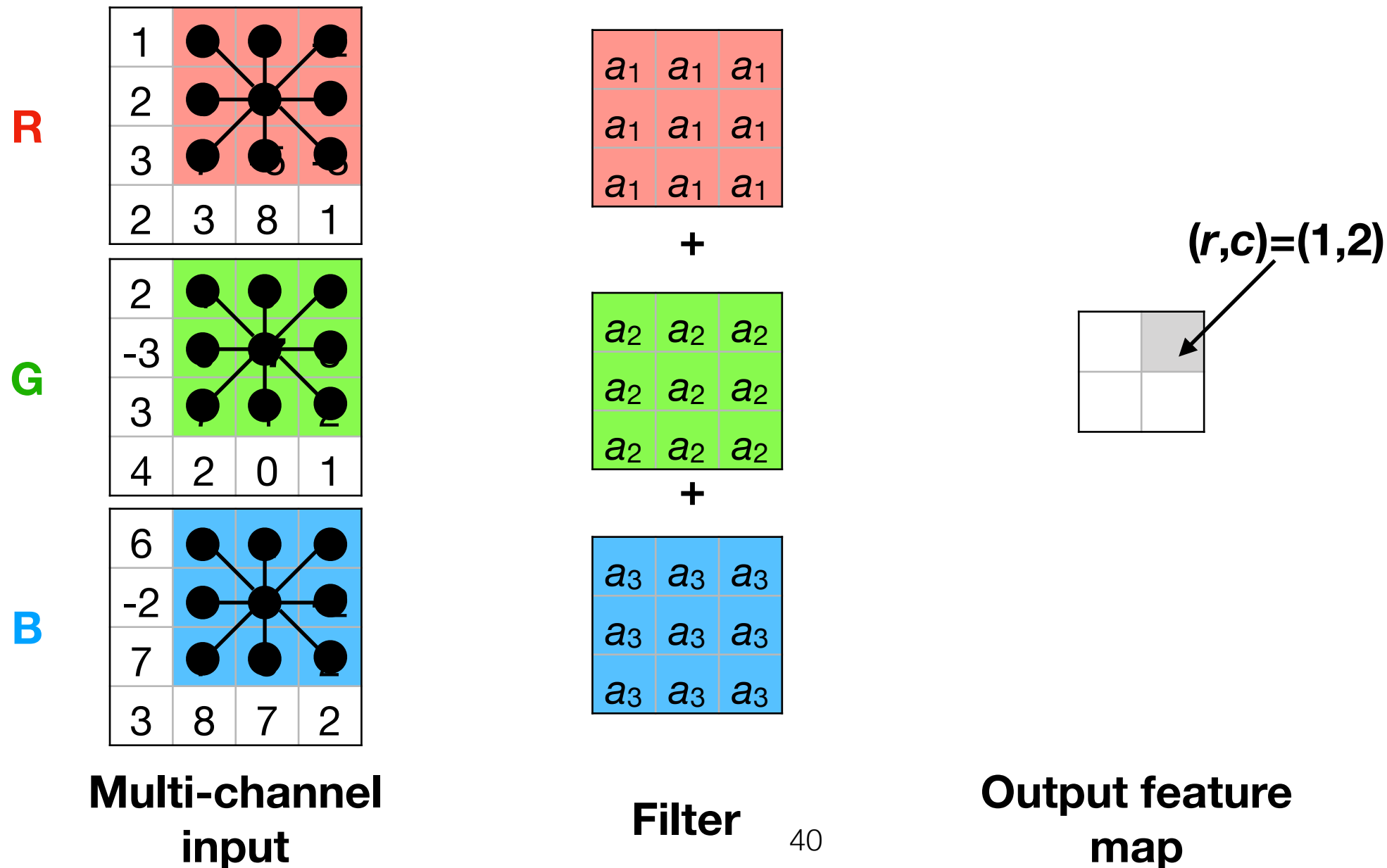
# Review of 3D convolution

- The value of each element  $(r,c)$  of the output feature map depends on the set of  $(k*k)$  **neighbors** of the corresponding element  $(r',c')$  in the input image.



# Special case

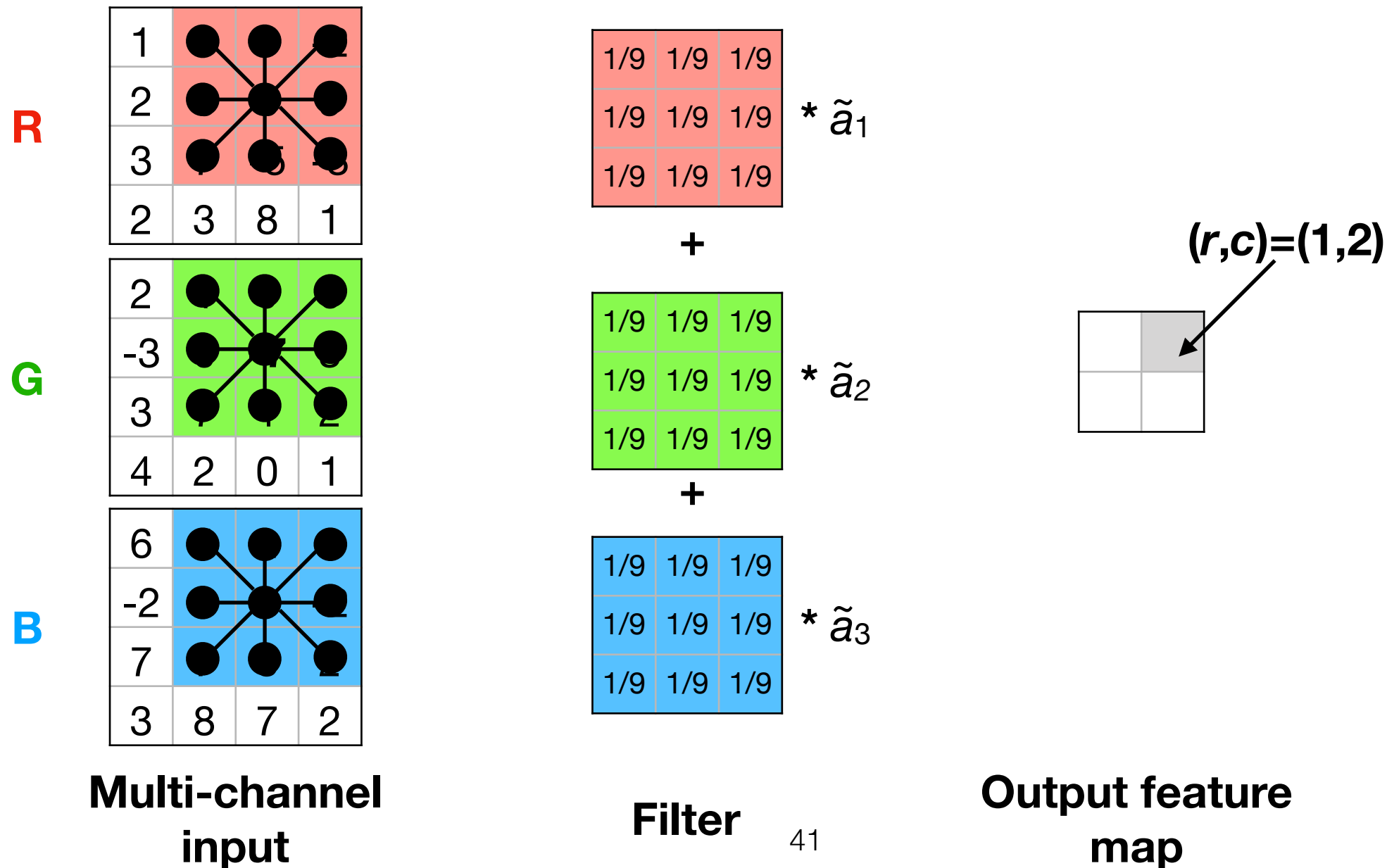
- Let's consider a special case where all the elements of each channel of the convolution filter are equal, i.e., there are just 3 learned parameters  $a_1$ ,  $a_2$ ,  $a_3$ .





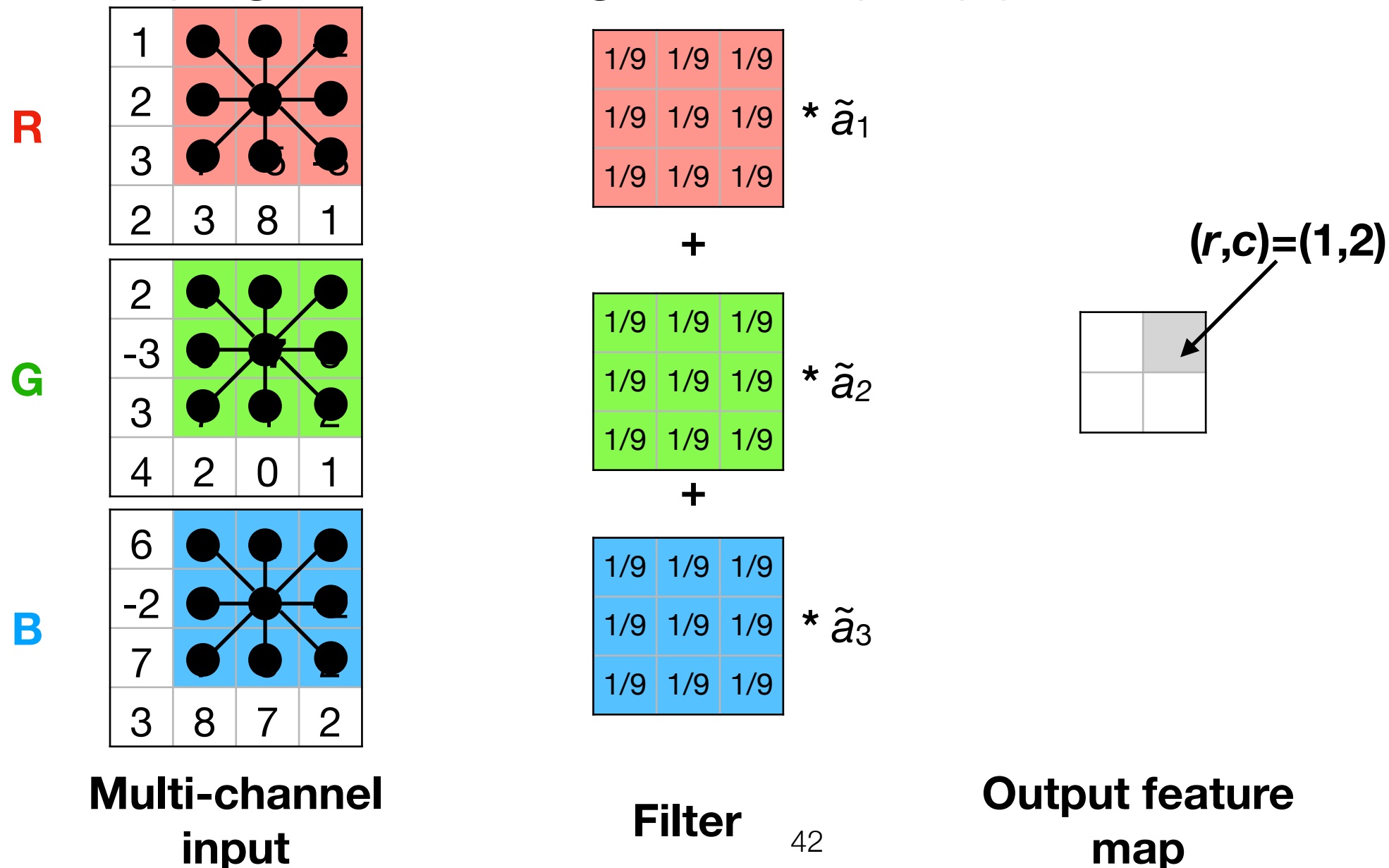
# Special case

- Now let's reparameterize it:



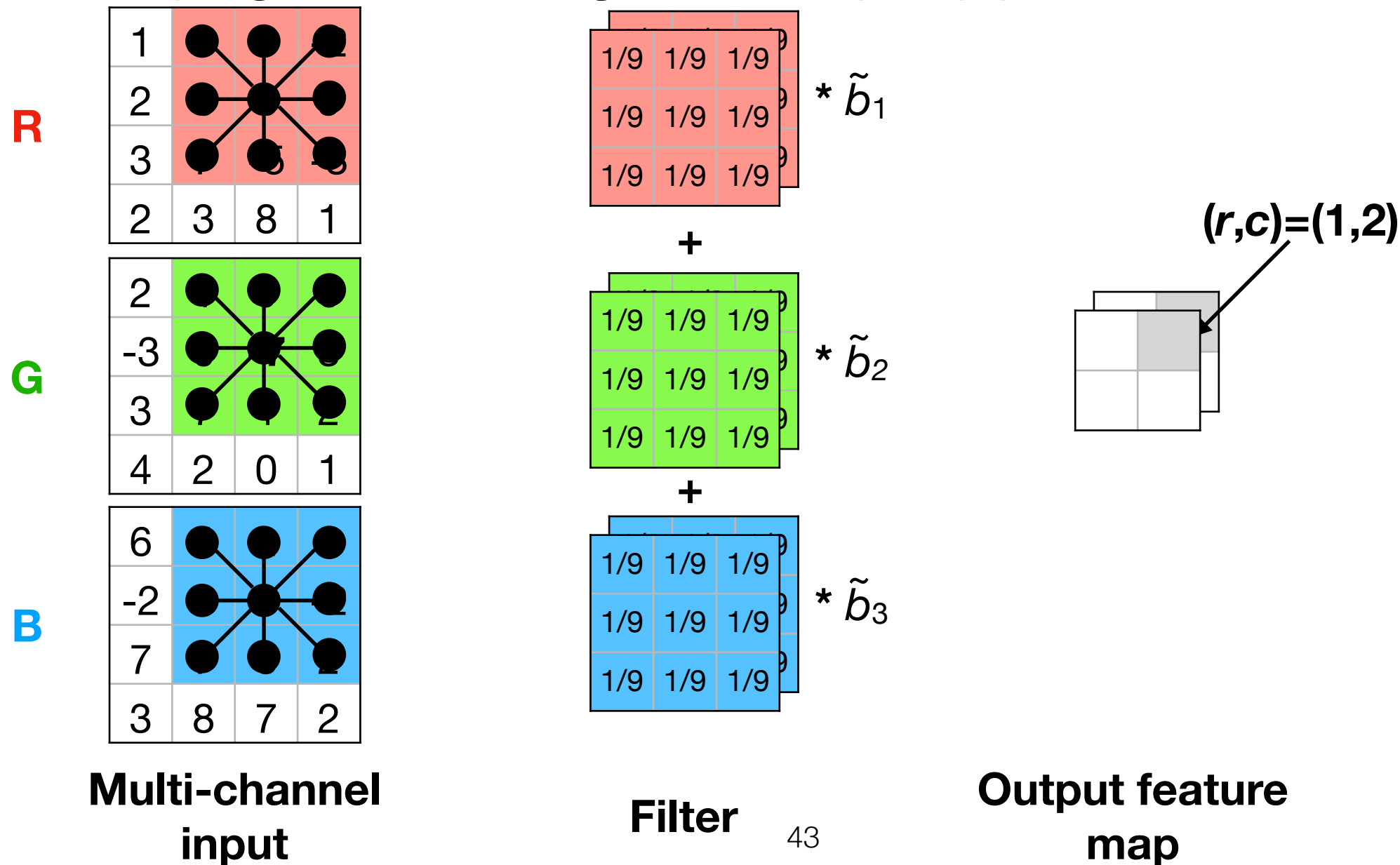
# Special case

- In this case, the output feature map at  $(r,c)$  equals:  
 $\tilde{a}_1 * (\text{Avg value of neighbors of } (r',c')^R) +$   
 $\tilde{a}_2 * (\text{Avg value of neighbors of } (r',c')^G) +$   
 $\tilde{a}_3 * (\text{Avg value of neighbors of } (r',c')^B)$



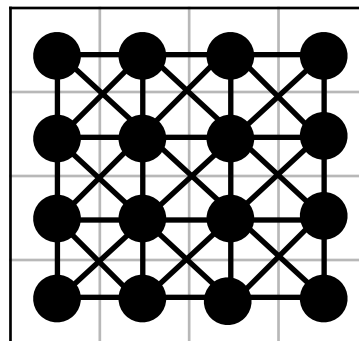
# Special case

- For a second filter, the output feature map at  $(r,c)$  would equal:  
 $\tilde{b}_1 * (\text{Avg value of neighbors of } (r',c')^R) +$   
 $\tilde{b}_2 * (\text{Avg value of neighbors of } (r',c')^G) +$   
 $\tilde{b}_3 * (\text{Avg value of neighbors of } (r',c')^B)$



# Images as graphs

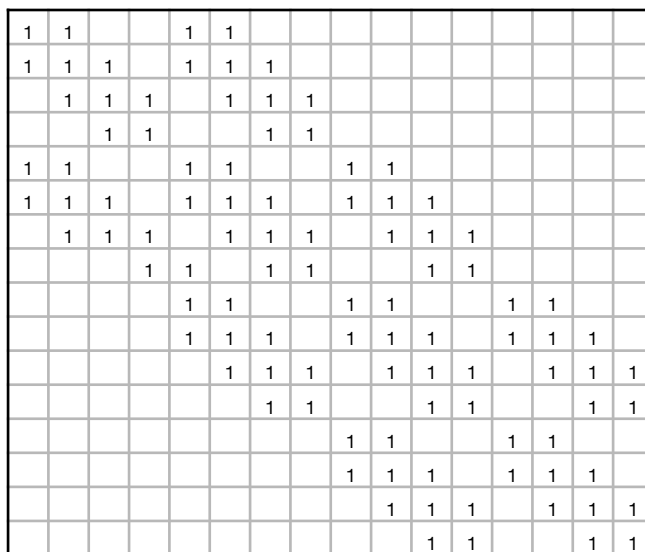
- Let us now represent an image as an undirected graph with:
  - One node per pixel.
  - An edge between each pair of adjacent pixels.



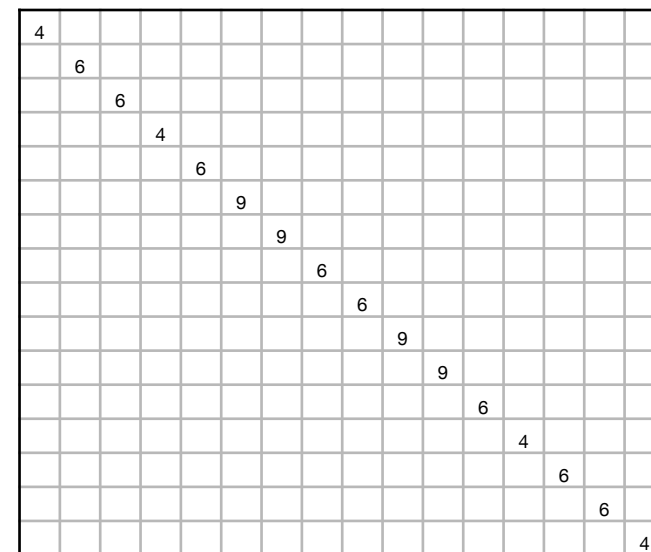
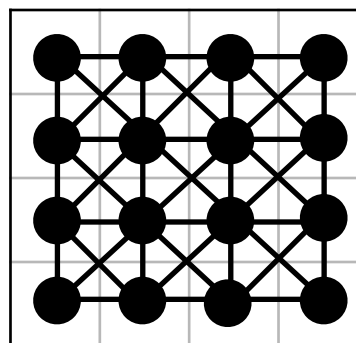
# Images as graphs

- We can represent this graph using an **adjacency matrix**.
- In particular, for an image with  $n$  pixels ( $n=16$  in this example), let  $\mathbf{A} \in \{0, 1\}^{n \times n}$  where  $A_{ij}=1$  iff pixel  $i$  neighbors pixel  $j$ .
- Also define diagonal matrix  $\mathbf{D}$  to count the total number of neighbors of each node  $i$ :  $D_{ii} = \sum_j A_{ij}$

If we include self-connections ( $A_{ii}=1$ ) then  $D_{ii}$  is 9 for interior pixels, 6 for edges, and 4 for corners in this example.)



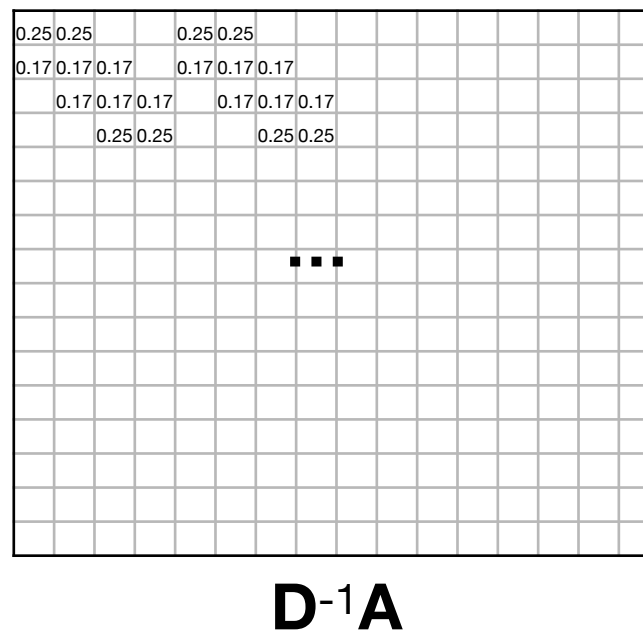
**A**



**D**

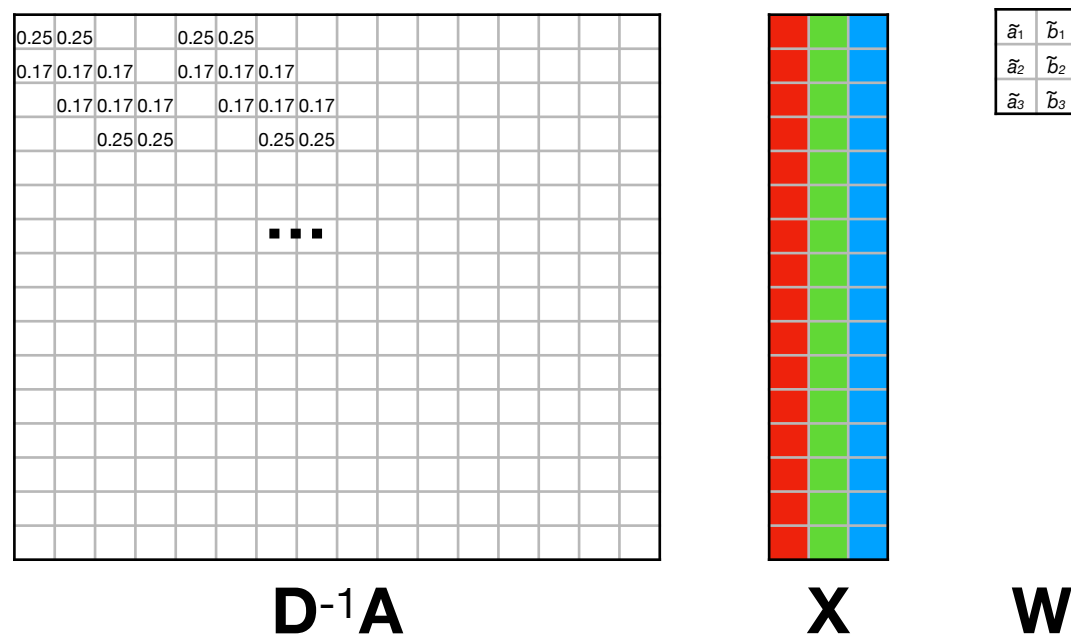
# Images as graphs

- By computing  $\mathbf{D}^{-1}\mathbf{A}$ , we can compute a normalized adjacency matrix:



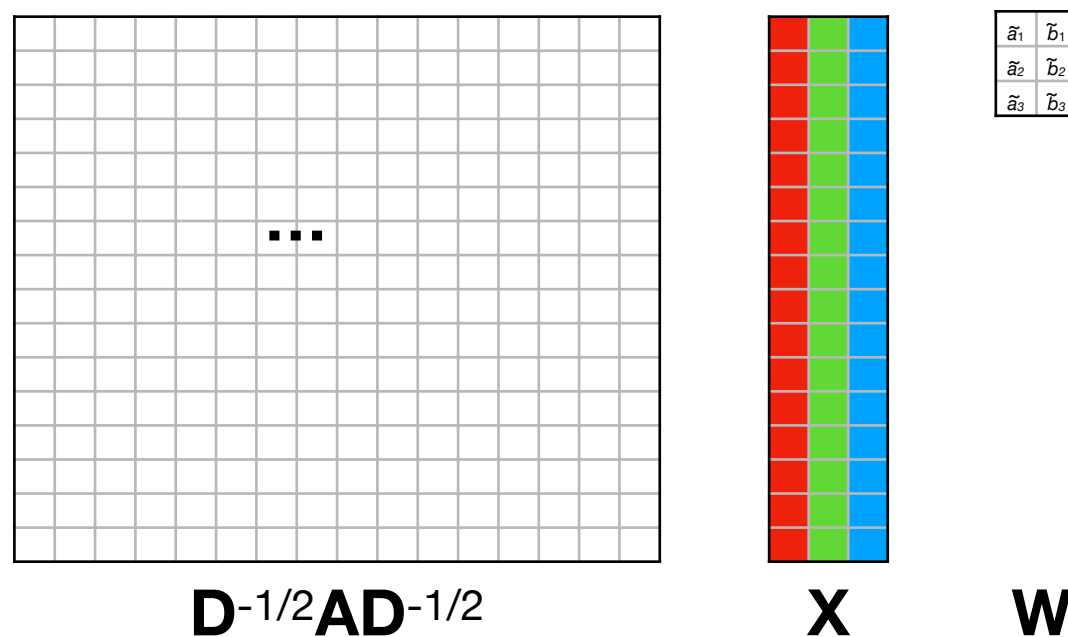
# Images as graphs

- We can now express this convolution as the product of the normalized adjacency matrix  $\mathbf{D}^{-1}\mathbf{A}$  with the input features  $\mathbf{X}$  and the filter weights  $\mathbf{W}$ :



# Images as graphs

- We can now express this convolution as the product of the normalized adjacency matrix  $\mathbf{D}^{-1}\mathbf{A}$  with the input features  $\mathbf{X}$  and the filter weights  $\mathbf{W}$ :

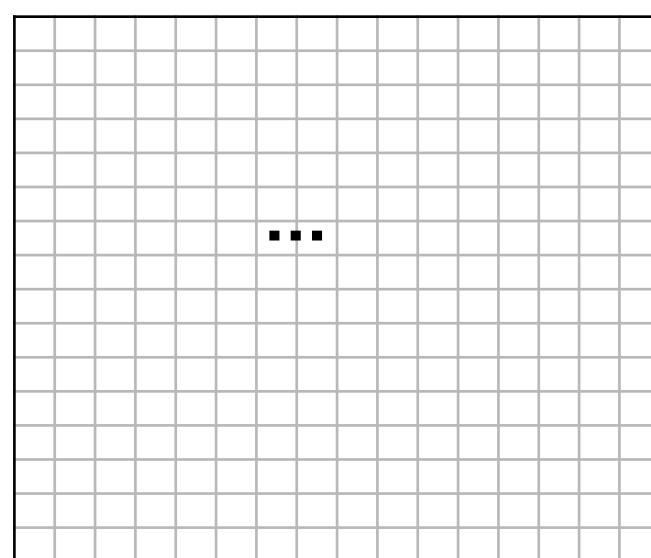


- A related matrix is given by  $\mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ .



# Images as graphs

- We can now express this convolution as the product of the normalized adjacency matrix  $\mathbf{D}^{-1}\mathbf{A}$  with the input features  $\mathbf{X}$  and the filter weights  $\mathbf{W}$ :



$\mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$



$\mathbf{X}$

$\tilde{a}_1$	$\tilde{b}_1$
$\tilde{a}_2$	$\tilde{b}_2$
$\tilde{a}_3$	$\tilde{b}_3$

$\mathbf{W}$

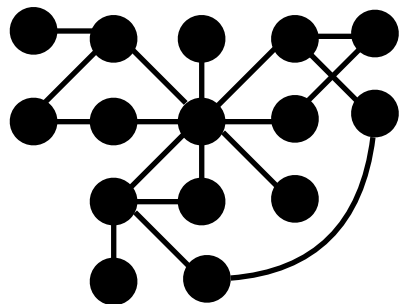
$$\mathbf{h} = \sigma(\mathbf{D}^{-\frac{1}{2}}\mathbf{A}\mathbf{D}^{-\frac{1}{2}}\mathbf{X}\mathbf{W})$$

- A related matrix is given by  $\mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ .
- We can then apply a non-linear activation function  $\sigma$ .

# Images as graphs

- We can apply the same procedure to general graphs.
- Moreover, the number of neighbors of each node does not have to be the same.

$$\mathbf{h} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} \right)$$



- Intuitively, we convolve the features of each graph node with a fixed filter (with equal value for all spatial neighbors) and renormalize according to the neighborhood structure.
- This is the basic operation of a **graph convolutional neural network** (GCN) (as proposed by Kipf & Welling 2017).

# Exercises

[https://cs230.stanford.edu/files/cs230exam\\_fall20\\_soln.pdf](https://cs230.stanford.edu/files/cs230exam_fall20_soln.pdf)

# Exercise 1

**(2 points)** You are training a large feedforward neural network (100 layers) on a binary classification task, using a sigmoid activation in the final layer, and a mixture of *tanh* and *ReLU* activations for all other layers. You notice your weights to your a *subset* of your layers stop updating after the first epoch of training, even though your network has not yet converged. Deeper analysis reveals the *gradients* to these layers completely, or almost completely, go to zero very early on in training. Which of the following fixes could help? (You also note that your loss is still within a reasonable order of magnitude).

- (i) Increase the size of your training set
- (ii) Switch the ReLU activations with leaky ReLUs everywhere
- (iii) Add Batch Normalization before every activation
- (iv) Increase the learning rate

# Exercise 2

**(2 points)** Which of the following would you consider to be valid activation functions (elementwise non-linearities) to train a neural net in practice?

(i)  $f(x) = -\min(2, x)$

(ii)  $f(x) = 0.9x + 1$

(iii)  $f(x) = \begin{cases} \min(x, .1x) & | x \geq 0 \\ \min(x, .1x) & | x < 0 \end{cases}$

(iv)  $f(x) = \begin{cases} \max(x, .1x) & | x \geq 0 \\ \min(x, .1x) & | x < 0 \end{cases}$

# Exercise 3

**(2 points)** During backpropagation, as the gradient flows backward through a *sigmoid* non-linearity, the gradient will always:

- (i) Increase in magnitude, maintain polarity
- (ii) Increase in magnitude, reverse polarity
- (iii) Decrease in magnitude, maintain polarity
- (iv) Decrease in magnitude, reverse polarity

# Exercise 4

- 1 (2 points) You are benchmarking runtimes for layers commonly encountered in CNNs. Which of the following would you expect to be the fastest (in terms of floating point operations)?
- (i) Conv layer (convolution operation + bias addition)
  - (ii) Max pooling
  - (iii) Average pooling
  - (iv) Batch Normalization

# Exercise 5

**(2 points)** You come across a nonlinear function that passes 1 if its input is nonnegative, else evaluates to 0, i.e.

$$f(x) = \begin{cases} 1 & | x \geq 0 \\ 0 & | x < 0 \end{cases}$$

A friend recommends you use this non-linearity in your convolutional neural network with the Adam optimizer. Would you follow their advice? Why or why not?