









### INPUT DATA

- . We need to get features!
- . For tabular data, this is "simple"
- . But what with more complex data?















### INPUT DATA - SOUND, IMAGES, VIDEO - ENCODING - MLP - LIMITATIONS

VU

VU

- . The input dimensions are very big
- . Too big for an MLP
- . Too many weights
- . Would not converge
- . would not fit in GPU memory
- . Especially when you also need to keep gradient information

INPUT DATA - SOUND, IMAGES, VIDEO - ENCODING

- . The features in this kind of data are **not** independent
- . They have locality

15

. But, an MLP does not remember this ordering



### GOAL FOR TODAY

- . We want to be able to do deep learning on this kind of data
- . Steps
- . 1D
  - . Build intuition
- . 2D
- . Do the same for images

10

























### 1D - SOUND - FILTERS

- . The features in the soundwave are not independent
- . Nearby features are more important as far away ones
- . This idea can be used in filters



VU



31













































## . To understand how the cor definition. For a 1D input so

VU

PART TWO-b: conv1D

### **1D CONVOLUTION - FORMAL DEFINITION**

. To understand how the conv. kernels are learned, we require a formal definition. For a 1D input sequence  $\mathbf{x} \in \mathbb{R}^n$  and a filter  $\mathbf{k} \in \mathbb{R}^{2m+1}$  the convolution operation is:

$$\mathbf{h}(t) = (\mathbf{x} * \mathbf{k})(t) = \sum_{\tau = -m}^{m} \mathbf{x}(t - \tau) \cdot \mathbf{k}(\tau)$$



### **1D CONVOLUTION - FORMAL DEFINITION**

- . How do we learn the filter k?
- . If **k** is learned, we will update the filter weights based on some loss  $\mathcal{L}_{\mathbf{h}} = \sum_{t=0}^{n} \mathcal{L}_{\mathbf{h}}(t)$  depending on the conv. response at all places, e.g., cross-correlation for classification.
- . The gradient utilized to update a kernel weight  $\mathbf{k}(\tau_0)$  is given by:



## $\begin{aligned} \mathbf{h}(t) &= (\mathbf{x} \ast \mathbf{k})(t) = \sum_{\substack{\tau = -m \\ \tau = -m \\ \partial \mathcal{L}_{\mathbf{h}}}}^{m} \mathbf{x}(t-\tau) \cdot \mathbf{k}(\tau) & \text{Always zero,} \\ & \underbrace{\partial \mathcal{L}_{\mathbf{h}}}_{\partial \mathbf{k}(\tau_0)} = \frac{\partial \mathcal{L}_{\mathbf{h}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{k}(\tau_0)} = \sum_{\substack{\tau = -m \\ t = 0 \\ d \mathbf{h}(t)}}^{m} \underbrace{\partial \mathcal{L}_{\mathbf{h}}(t)}_{\tau = -m} \sum_{\substack{\tau = -m \\ d \mathbf{h}(t)}}^{m} \mathbf{x}(t-\tau) \cdot \underbrace{\partial \mathbf{k}(\tau)}_{\partial \mathbf{k}(\tau_0)} & \underbrace{\partial \mathbf{k}(\tau)}_{\text{tau}} \\ & \text{The update of the weights takes into consideration ALL} \\ & \text{ELEMENTS OF THE INPUT SEQUENCE INTO ACCOUNT!} \end{aligned}$ . The weights are shared across the entire input ( MLPs learn independent weights at every position: $\mathbf{h}(t) = \mathbf{W}(t, :) \cdot \mathbf{x}(t)$ )

VU

VU

### **1D CONVOLUTION - FORMAL DEFINITION**

### . Advantages:

- Since weights are shared for every position, Convolutional Networks (CNNs) are much MUCH! MUCH! smaller than MLPs. -> PARAMETER EFFICIENCY
- Convolutions can learn a powerful pattern recognizers, e.g., for silence, based on "silences" appearing everywhere in the input (MLPs must learn an independent "silence recognizer" for every position) -> DATA EFFICIENCY
- Convolutions can recognize a "silence pattern" regardless of where it appears (MLPs must have seen silence at a given position before in order to recognize it)
   > GENERALIZATION IMPROVEMENTS

**IMPORTANT:** 2 and 3 are a consequence of convolution being translation equivariant.

### **1D - CONVOLUTIONS**

### Up till now:

- . We can create filters
- MLP in disguise
- . We can convolve them over the input which creates an output vector

### Coming next:

- . We need multiple filters
- . What do we do at the start and end of the data?
- . What at the next layer?
- . How do we get the dimension down?
- 60



VU

### **1D - CONVOLUTIONS**

### Up till now:

- . We can create filters
- . MLP in disguise
- . We can convolve them over the input which creates an output vector

### Coming next:

- . We need multiple filters
- . What at the next layer?
- . What do we do at the start and end of the data?
- . How do we get the dimension down?
- 63

### **1D - CONVOLUTIONS - FILTERS AT FURTHER LAYERS**

- What do we do at the next layer?
- . We have the output volume of the previous layer and we will just define
- a convolution operator over that!





### 1D - CONVOLUTIONS - FILTERS AT FURTHER LAYERS

What do we do at the next layer?

. We have the output volume of the previous layer and we will just define

a convolution operator over that!



### **1D - CONVOLUTIONS - FILTERS AT FURTHER LAYERS**

What do we do at the next layer?

- . We have the output volume of the previous layer and we will just define  $% \mathcal{A}_{\mathcal{A}}$
- a convolution operator over that!



### **1D - CONVOLUTIONS - FILTERS AT FURTHER LAYERS**

What do we do at the next layer?

- . We have the output volume of the previous layer and we will just define
- a convolution operator over that!





What do we do at the next layer?

- . We have the output volume of the previous layer and we will just define
- a convolution operator over that!



### **1D - CONVOLUTIONS**

### Up till now:

- . We can create filters
- . MLP in disguise
- . We can convolve them over the input which creates an output vector

VU

Coming next:

VU

- . We need multiple filters
- . What at the next layer?
- . What do we do at the start and end of the data?
- . How do we get the dimension down?

### **1D - CONVOLUTIONS - PADDING**

69

. Near the boundaries of the data, we cannot apply the convolution as we normally do:



### **1D - CONVOLUTIONS - PADDING**

. Near the boundaries of the data, we cannot apply the convolution as we normally do:



### **1D - CONVOLUTIONS - PADDING**

. Near the boundaries of the data, we cannot apply the convolution as we normally do:



• Ignore the boundaries

reduction of dimension!

information at the

VU

### **1D - CONVOLUTIONS - PADDING**

. Near the boundaries of the data, we cannot apply the convolution as we normally do:



### **1D - CONVOLUTIONS - PADDING**

73

. Near the boundaries of the data, we cannot apply the convolution as we normally do:



## **1D - CONVOLUTIONS - PADDING** . Near the boundaries of the data, we cannot apply the convolution as we normally do:





























# 2D - IMAGES - REPRESENTATION - 3D TENSOR The 2 dimensional color image becomes a 3 dimensional tensor! Image: the state of the

# A 3 dimensional color video becomes a 4 dimensional tensor!







2D - IMAGES - CONVOLUTION - PADDING			
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Filter W1 (3:3:3)       We start with a 5x5 in         v1 (1;1:0)       We start with a 5x5 in         1 1       We have 3 channels         1 1       We have 3 channels         1 1       We want to use 2 filter         1 1       We want to use 2 filter         1 1       We add padding to set         0 1 1       We add padding to set         with convolving near       With convolving near	nage ers, these are ional olve issues the border	
0         1         0         2         0         0         0         •           0         0         0         0         0         0         •         •           0         0         0         0         0         0         •         •           0         0         0         0         0         0         •         •           0         2         0         0         1         0         •         •           0         1         1         2         1         0         •         •           0         0         0         2         1         0         •         •         •           0         0         1         1         2         2         0         •         •         •           0         0         1         1         2         2         •         •         •         •           0         0         1         1         2         2         •         •         •         •         •	▼ ation/image source: https://cs231n.github.io/convolutional-networks/#fc	VU∕	

### **2D - IMAGES - CONVOLUTION - OUTPUT VOLUME** Filter W1 (3x3x3) Output Volume (3x3x2) Input Volume (+pad 1) (7x7x3) Filter W0 (3x3x3) x[:,:,0] 0 0 0 0 0 0 0 w1[:,:,0] 0 1 1 -100 o[:,:,0] 5 0 -2 We start with a 5x5 image 0 1 1 0 2 1 0 -1 -1 1 1 1 -1 -2 -3 0 0 2 1 0 2 0 0 0 1 -1 -1 1 -1 -4 -9 -1 We have 3 channels 0 2 2 1 2 2 0 o[:,:,1] 8 -1 7 w0[:,:,1] -1 0 -1 w1[:,:,1] 1 -1 0 0 0 2 2 1 2 0 1 -1 0 0 1 -1 1 1 1 4 4 8 We want to use 2 filters, these are 0 2 1 0 0 2 0 4 9 7 -1 1 0 0 0 0 0 0 0 0 w1[:,:,2] 0 1 -1 themselves 3 dimensional w0[:,:,2 0 -1 -1 0 0 0 0 0 0 0 1 1 .1 0 1 1 0 1 2 1 0 2 0 -1 -1 -1 1 0 0 We add padding to solve issues 0 2 0 0 0 2 0 0 0 2 0 1 0 0 Bias b0 (1x1x1) Bias b1 (1x1x1) with convolving near the border 0 1 2 0 1 2 0 b0[:,:,0] b1[:,:,0] 0 1 0 2 0 1 0 0 0 0 0 0 0 0 We convolve with a stride of 2 0 0 0 0 0 0 0 0 2 0 0 0 1 0 Try to understand why the output 0 0 0 0 1 0 0 0 1 1 0 2 1 0 volume has these dimensions. 0 0 0 0 2 1 0

animation/image source: https://cs231n.github.io/convolutional-networks/#fc

0 0 1 1 2 2 0

0000000

2D - IMAGES - REPRESENTATION				
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Output Volume (3x3x2)         0 = 1         2 = 30         2 = 30         4 = 9         6 = 1         8 = 1         8 = 1         8 = 1         8 = 1         8 = 1         8 = 1         8 = 1         9 = 7         We have 3 channels         We want to use 2 filters, these are themselves 3 dimensional         We add padding to solve issues with convolving near the border		
0         0	on/image source: htt	We convolve with a stride of 2 See the animation on the site ps://cs231n.github.econvoluenal-networks##c		

### CONVOLUTION - EQUIVARIANCE ANALYSIS

We saw that convolutions are **equivariant** to translations. Naturally it holds for ConvNDs as well:



### CONVOLUTION - EQUIVARIANCE ANALYSIS

What about other transformations?, e.g., rotation, scaling, ... Are ConvNDs rotation and scale equivariant?

The convolution is **not** a **rotation or scale** equivariant mapping.

### CONVOLUTION - EQUIVARIANCE ANALYSIS

What about other transformations?, e.g., rotation, scaling, ... Are ConvNDs **rotation and scale** equivariant?

The convolution is **not** a **rotation or scale** equivariant mapping. But a network can learn rotated / scaled versions of the same filter.



**BUT** approx. equivariant for scales / rotations learned by the network.

VU



\* Let us know if you are interested in writing your thesis in this topic ;)







